WMT 2024

# Ninth Conference on Machine Translation

## Proceedings of the Conference

November 15-16, 2024

Order copies of this and other ACL proceedings from:

# Introduction

The Ninth Conference on Machine Translation (WMT 2024) took place on Friday, November 15 and Saturday, November 16, 2024, immediately following the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP 2024) in Miami, Florida, USA.

This is the ninth time WMT has been held as a conference. The first time WMT was held as a conference was at ACL 2016 in Berlin, Germany, the second time at EMNLP 2017 in Copenhagen, Denmark, the third time at EMNLP 2028 in Brussels, Belgium, the fourth time at ACL 2019 in Florence, Italy, the fifth time at EMNLP-2020, which was held as an online event due to the COVID-19 pandemic, the sixth time at EMNLP 2021 at Punta Cana, Dominican Republic, the seventh time at EMNLP 2022 in Abu Dhabi, United Arab Emirates, and the eight time at EMNLP 2023 in Singapore. Prior to being a conference, WMT was held 10 times as a workshop. WMT was held for the first time at HLT-NAACL 2006 in New York City, USA. In the following years the Workshop on Statistical Machine Translation was held at ACL 2007 in Prague, Czech Republic, ACL 2008, Columbus, Ohio, USA, EACL 2009 in Athens, Greece, ACL 2010 in Uppsala, Sweden, EMNLP 2011 in Edinburgh, Scotland, NAACL 2012 in Montreal, Canada, ACL 2013 in Sofia, Bulgaria, ACL 2014 in Baltimore, USA, EMNLP 2015 in Lisbon, Portugal.

The focus of our conference is to bring together researchers from the area of machine translation and invite selected research papers to be presented at the conference.

Prior to the conference, in addition to soliciting relevant papers for review and possible presentation, we conducted 12 shared tasks. These consisted of 8 translation tasks: General Translation, Translation into Low-Resource Languages of Spain, Low-Resource Indic Language Translation, Chat Translation, Biomedical Translation, MultiIndic22MT, Non-Repetitive Translation, English-to-Lowres Multi-Modal Translation, three evaluation tasks: Metrics, MT Test Suites, Quality Estimation, and finally the Open Language Data Initiative.

The results of all shared tasks were announced at the conference, and these proceedings also include overview papers for the shared tasks, summarizing the results, as well as providing information about the data used and any procedures that were followed in conducting or scoring the tasks. In addition, there are short papers from each participating team that describe their underlying system in greater detail.

Like in previous years, we have received a far larger number of submissions than we could accept for presentation. WMT 2024 has received 54 full research paper submissions (not counting withdrawn submissions). In total, WMT 2024 featured 25 full research paper presentations and 96 shared task presentations.

The invited talk entitled "What makes MT research special in the LLM age?was given by Ricardo Rei and Nuno M. Guerreiro from Unbabel, Portugal.

We would like to thank the members of the Program Committee for their timely reviews. We also would like to thank the participants of the shared task and all the other volunteers who helped with the evaluations.

Barry Haddow, Tom Kocmi, Philipp Koehn, and Christof Monz

Co-Organizers

# Organizing Committee

**Chairs**

Barry Haddow, University of Edinburgh
Tom Kocmi, Cohere
Philipp Koehn, Johns Hopkins University
Christof Monz, University of Amsterdam

# Program Committee

**Program Committee**

Sadaf Abdul Rauf, Fatima Jinnah Women Unversity
Antonios Anastasopoulos, George Mason University
Duygu Ataman, New York University
Eleftherios Avramidis, German Research Center for Artificial Intelligence (DFKI)
Seth Aycock, University of Amsterdam
Parnia Bahar, Apple
Rachel Bawden, Inria
Alexandra Birch, University of Edinburgh
Nikolay Bogoychev, University of Edinburgh
Sheila Castilho, Dublin City University
Chenhui Chu, Kyoto University
Ann Clifton, Spotify
Josep Crego, CHAPSVISION
Steve Deneefe, RWS Language Weaver
Michael Denkowski, Amazon
Marion Di Marco, TUM
Shuoyang Ding, Amazon
Miguel Domingo, Universitat Politècnica de València
Kevin Duh, Johns Hopkins University
Koel Dutta Chowdhury, Saarland Informatics Campus,Saarland University
Hiroshi Echizen'ya, Hokkai-Gakuen University
Miquel Esplà-Gomis, Universitat d'Alacant
Marcello Federico, Amazon
Atsushi Fujita, National Institute of Information and Communications Technology
Jesús González-Rubio, WebInterpret
Isao Goto, Ehime University
Thamme Gowda, Microsoft
Jeremy Gwinnup, Air Force Research Laboratory
Jindřich Helcl, Charles University in Prague
John Henderson, Mechanical Learning
Amr Hendy, Microsoft
Kenji Imamura, National Institute of Information and Communications Technology
Vivek Iyer, The University of Edinburgh
Josef Jon, Charles University
Shahram Khadivi, eBay
Mateusz Klimaszewski, Warsaw University of Technology
Mateusz Krubiński, Snowflake
Roland Kuhn, National Research Council of Canada
Gaurav Kumar, Bloomberg LP
Anoop Kunchukuttan, Microsoft AI and Research
Wen Lai, Technical University of Munich
Surafel M. Lakew, Amazon.com, Inc
Ekaterina Lapshinova-Koltunski, University of Hildesheim
Samuel Larkin, National Research Council Canada
Gregor Leusch, eBay
Baohao Liao, University of Amsterdam

Jean Maillard, Meta AI
Bhavitvya Malik, University of Edinburgh
Yan Meng, University of Amsterdam
Antonio Valerio Miceli Barone, The University of Edinburgh
Kenton Murray, Johns Hopkins University
Jan Niehues, Karlsruhe Institut of Technology
Xing Niu, Amazon AI
Tsuyoshi Okita, Kyushu institute of technology
Daniel Ortiz-Martinez, University of Barcelona
Santanu Pal, Wipro
Jianhui Pang, University of Macau
Stephan Peitz, Apple
Sergio Penkale, Lingo24
Matt Post, Microsoft
Vikas Raunak, Microsoft
Ricardo Rei, Unbabel/INESC-ID
Matiss Rikters, AIST
Elizabeth Salesky, Johns Hopkins University
Rico Sennrich, University of Zurich
Patrick Simianer, Lilt
Felix Stahlberg, Google Research
David Stap, University of Amsterdam
Katsuhito Sudoh, Nara Women's University
Felipe Sánchez-Martínez, Universitat d'Alacant
Aleš Tamchyna, Memsource
Shaomu Tan, University of Amsterdam
Gongbo Tang, Beijing Language and Culture University
Jörg Tiedemann, University of Helsinki
Evgeniia Tokarchuk, University of Amsterdam
Antonio Toral, University of Groningen
Masao Utiyama, NICT
Dusan Varis, Charles University, Institute of Formal and Applied Linguistics
David Vilar, Google
Martin Volk, University of Zurich
Ekaterina Vylomova, University of Melbourne
Longyue Wang, Tencent AI Lab
Taro Watanabe, Nara Institute of Science and Technology
Di Wu, University of Amsterdam
Joern Wuebker, Lilt, Inc.
Tong Xiao, Northeastern University
François Yvon, ISIR CNRS and Sorbonne Université
Zhong Zhou, Carnegie Mellon University
Vilém Zouhar, ETH Zurich, Charles University

# Table of Contents

vii

viii

xi

xiii

# Program

**Friday, November 15, 2024**

08:45 - 09:00    *Opening Remarks*

09:00 - 10:30    *Session 1 — Shared Task Overview Papers I*

*Findings of the WMT24 General Machine Translation Shared Task: The LLM Era Is Here but MT Is Not Solved Yet*
Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinthór Steingrímsson and Vilém Zouhar

*Are LLMs Breaking MT Metrics? Results of the WMT24 Metrics Shared Task*
Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-Kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva and Alon Lavie

*Findings of the Quality Estimation Shared Task at WMT 2024: Are LLMs Closing the Gap in QE?*
Chrysoula Zerva, Frederic Blain, José G. C. De Souza, Diptesh Kanojia, Sourabh Deoghare, Nuno M. Guerreiro, Giuseppe Attanasio, Ricardo Rei, Constantin Orasan, Matteo Negri, Marco Turchi, Rajen Chatterjee, Pushpak Bhattacharyya, Markus Freitag and André Martins

*Findings of the WMT 2024 Shared Task of the Open Language Data Initiative*
Jean Maillard, Laurie Burchell, Antonios Anastasopoulos, Christian Federmann, Philipp Koehn and Skyler Wang

*Results of the WAT/WMT 2024 Shared Task on Patent Translation*
Shohei Higashiyama

*Findings of the WMT 2024 Biomedical Translation Shared Task: Test Sets on Abstract Level*
Mariana Neves, Cristian Grozea, Philippe Thomas, Roland Roller, Rachel Bawden, Aurélie Névéol, Steffen Castle, Vanessa Bonato, Giorgio Maria Di Nunzio, Federica Vezzani, Maika Vicente Navarro, Lana Yeganova and Antonio Jimeno Yepes

10:30 - 11:00    *Coffee Break*

11:00 - 12:00    *Session 1 — Shared Task Overview Papers I*

11:00 - 12:00    *General Translation Task*

*MSLC24 Submissions to the General Machine Translation Task*
Samuel Larkin, Chi-Kiu Lo and Rebecca Knowles

*IOL Research Machine Translation Systems for WMT24 General Machine Translation Shared Task*
Wenbo Zhang

xvi

**Friday, November 15, 2024 (continued)**

*FLORES+ Translation and Machine Translation Evaluation for the Erzya Language*
Isai Gordeev, Sergey Kuldin and David Dale

*Spanish Corpus and Provenance with Computer-Aided Translation for the WMT24 OLDI Shared Task*
Jose Cols

12:30 - 13:30    *Patent Translation Task*

*Efficient Terminology Integration for LLM-based Translation in Specialized Domains*
Sejoon Kim, Mingi Sung, Jeonghwan Lee, Hyunkuk Lim and Jorge Gimenez Perez

*Rakuten's Participation in WMT 2024 Patent Translation Task*
Ohnmar Htun and Alberto Poncelas

12:30 - 13:30    *Biomedical Translation Task*

*The SETU-ADAPT Submission for WMT 24 Biomedical Shared Task*
Antonio Castaldo, Maria Zafar, Prashanth Nayak, Rejwanul Haque, Andy Way and Johanna Monti

14:00 - 15:00    *Session 4 — Invited Talk by Ricardo Rei and Nuno M. Guerreiro. What Makes MT Research Special in the LLM Age?"*

15:00 - 15:30    *Coffee Break*

15:30 - 17:00    *Session 5 — Featured Research Papers Oral Presentations*

*Translating Step-by-Step: Decomposing the Translation Process for Improved Translation Quality of Long-Form Texts*
Eleftheria Briakou, Jiaming Luo, Colin Cherry and Markus Freitag

*Is Preference Alignment Always the Best Option to Enhance LLM-Based Translation? An Empirical Analysis*
Hippolyte Gisserot-Boukhlef, Ricardo Rei, Emmanuel Malherbe, Céline Hudelot, Pierre Colombo and Nuno M. Guerreiro

*On Instruction-Finetuning Neural Machine Translation Models*
Vikas Raunak, Roman Grundkiewicz and Marcin Junczys-Dowmunt

**Saturday, November 16, 2024**

09:00 - 10:30    *Session 6 — Shared Task Overview Papers II*

*Findings of WMT 2024 Shared Task on Low-Resource Indic Languages Translation*
Partha Pakray, Santanu Pal, Advaitha Vetagiri, Reddi Krishna, Arnab Kumar Maji, Sandeep Dash, Lenin Laitonjam, Lyngdoh Sarah and Riyanka Manna

*Findings of WMT 2024's MultiIndic22MT Shared Task for Machine Translation of 22 Indian Languages*
Raj Dabre and Anoop Kunchukuttan

*Findings of WMT2024 English-to-Low Resource Multimodal Translation Task*
Shantipriya Parida, Ondřej Bojar, Idris Abdulmumin, Shamsuddeen Hassan Muhammad and Ibrahim Said Ahmad

*Findings of the WMT 2024 Shared Task Translation into Low-Resource Languages of Spain: Blending Rule-Based and Neural Systems*
Felipe Sánchez-Martínez, Juan Antonio Perez-Ortiz, Aaron Galiano Jimenez and Antoni Oliver

*Findings of the WMT 2024 Shared Task on Discourse-Level Literary Translation*
Longyue Wang, Siyou Liu, Chenyang Lyu, Wenxiang Jiao, Xing Wang, Jiahao Xu, Zhaopeng Tu, Yan Gu, Weiyu Chen, Minghao Wu, Liting Zhou, Philipp Koehn, Andy Way and Yulin Yuan

*Findings of the WMT 2024 Shared Task on Chat Translation*
Wafaa Mohammed, Sweta Agrawal, Amin Farajian, Vera Cabarrão, Bryan Eikema, Ana C Farinha and José G. C. De Souza

*Findings of the WMT 2024 Shared Task on Non-Repetitive Translation*
Kazutaka Kinugawa, Hideya Mino, Isao Goto and Naoto Shirai

10:30 - 11:00    *Coffee Break*

11:00 - 12:00    *Session 7 — Shared Task Posters III*

11:00 - 12:00    *Low-Resource Indic Language Translation Task*

*A3-108 Controlling Token Generation in Low Resource Machine Translation Systems*
Saumitra Yadav, Ananya Mukherjee and Manish Shrivastava

**Saturday, November 16, 2024 (continued)**

*Samsung R&D Institute Philippines @ WMT 2024 Indic MT Task*
Matthew Theodore Roque, Carlos Rafael Catalan, Dan John Velasco, Manuel Antonio Rufino and Jan Christian Blaise Cruz

*DLUT-NLP Machine Translation Systems for WMT24 Low-Resource Indic Language Translation*
Chenfei Ju, Junpeng Liu, Kaiyu Huang and Degen Huang

*SRIB-NMT's Submission to the Indic MT Shared Task in WMT 2024*
Pranamya Patil, Raghavendra Hr, Aditya Raghuwanshi and Kushal Verma

*MTNLP-IIITH: Machine Translation for Low-Resource Indic Languages*
Abhinav P M, Ketaki Shetye and Parameswari Krishnamurthy

*Exploration of the CycleGN Framework for Low-Resource Languages*
Sören Dreano, Derek Molloy and Noel Murphy

*The SETU-ADAPT Submissions to the WMT24 Low-Resource Indic Language Translation Task*
Neha Gajakos, Prashanth Nayak, Rejwanul Haque and Andy Way

*SPRING Lab IITM's Submission to Low Resource Indic Language Translation Shared Task*
Advait Joglekar, Hamees Ul Hasan Sayed and Srinivasan Umesh

*Machine Translation Advancements of Low-Resource Indian Languages by Transfer Learning*
Bin Wei, Zheng Jiawei, Zongyao Li, Zhanglin Wu, Jiaxin Guo, Daimeng Wei, Zhiqiang Rao, Shaojun Li, Yuanchang Luo, Hengchao Shang, Jinlong Yang, Yuhao Xie and Hao Yang

*NLIP_Lab-IITH Low-Resource MT System for WMT24 Indic MT Shared Task*
Pramit Sahoo, Maharaj Brahma and Maunendra Sankar Desarkar

*Yes-MT's Submission to the Low-Resource Indic Language Translation Shared Task in WMT 2024*
Yash Bhaskar and Parameswari Krishnamurthy

11:00 - 12:00    *MultiIndic22MT Task*

*System Description of BV-SLP for Sindhi-English Machine Translation in MultiIndic22MT 2024 Shared Task*
Nisheeth Joshi, Pragya Katyayan, Palak Arora and Bharti Nathani

*WMT24 System Description for the MultiIndic22MT Shared Task on Manipuri Language*
Ningthoujam Justwant Singh, Kshetrimayum Boynao Singh, Ningthoujam Avichandra Singh, Sanjita Phijam and Thoudam Doren Singh

*NLIP-Lab-IITH Multilingual MT System for WAT24 MT Shared Task*
Maharaj Brahma, Pramit Sahoo and Maunendra Sankar Desarkar

11:00 - 12:00    *English-to-Lowres Multi-Modal Translation Task*

*DCU ADAPT at WMT24: English to Low-resource Multi-Modal Translation Task*
Sami Haq, Rudali Huidrom and Sheila Castilho

*English-to-Low-Resource Translation: A Multimodal Approach for Hindi, Malayalam, Bengali, and Hausa*
Ali Hatami, Shubhanker Banerjee, Mihael Arcan, Bharathi Chakravarthi, Paul Buitelaar and John Mccrae

*OdiaGenAI's Participation in WMT2024 English-to-Low Resource Multimodal Translation Task*
Shantipriya Parida, Shashikanta Sahoo, Sambit Sekhar, Upendra Jena, Sushovan Jena and Kusum Lata

*Arewa NLP's Participation at WMT24*
Mahmoud Ahmad, Auwal Khalid, Lukman Aliyu, Babangida Sani and Mariya Abdullahi

*Multimodal Machine Translation for Low-Resource Indic Languages: A Chain-of-Thought Approach Using Large Language Models*
Pawan Rajpoot, Nagaraj Bhat and Ashish Shrivastava

*Chitranuvad: Adapting Multi-lingual LLMs for Multimodal Translation*
Shaharukh Khan, Ayush Tarun, Ali Faraz, Palash Kamble, Vivek Dahiya, Praveen Pokala, Ashish Kulkarni, Chandra Khatri, Abhinav Ravi and Shubham Agarwal

*Brotherhood at WMT 2024: Leveraging LLM-Generated Contextual Conversations for Cross-Lingual Image Captioning*
Siddharth Betala and Ishan Chokshi

12:30 - 13:30      *Session 8 — Shared Task Posters IV*

12:30 - 13:30      *Translation into Low-Resource Languages of Spain Task*

*TIM-UNIGE Translation into Low-Resource Languages of Spain for WMT24*
Jonathan Mutal and Lucía Ormaechea

*TAN-IBE Participation in the Shared Task: Translation into Low-Resource Languages of Spain*
Antoni Oliver

*Enhaced Apertium System: Translation into Low-Resource Languages of Spain Spanish–Asturian*
Sofía García

*Universitat d'Alacant's Submission to the WMT 2024 Shared Task on Translation into Low-Resource Languages of Spain*
Aaron Galiano Jimenez, Víctor M. Sánchez-Cartagena, Juan Antonio Perez-Ortiz and Felipe Sánchez-Martínez

*Samsung R&D Institute Philippines @ WMT 2024 Low-resource Languages of Spain Shared Task*
Dan John Velasco, Manuel Antonio Rufino and Jan Christian Blaise Cruz

*Back to the Stats: Rescuing Low Resource Neural Machine Translation with Statistical Methods*
Menan Velayuthan, Dilith Jayakody, Nisansa De Silva, Aloka Fernando and Surangika Ranathunga

*Hybrid Distillation from RBMT and NMT: Helsinki-NLP's Submission to the Shared Task on Translation into Low-Resource Languages of Spain*
Ona De Gibert, Mikko Aulamo, Yves Scherrer and Jörg Tiedemann

*Robustness of Fine-Tuned LLMs for Machine Translation with Varying Noise Levels: Insights for Asturian, Aragonese and Aranese*
Martin Bär, Elisa Forcada Rodríguez and Maria Garcia-Abadillo

*Training and Fine-Tuning NMT Models for Low-Resource Languages Using Apertium-Based Synthetic Corpora*
Aleix Sant, Daniel Bardanca, José Ramom Pichel Campos, Francesca De Luca Fornaciari, Carlos Escolano, Javier Garcia Gilabert, Pablo Gamallo, Audrey Mash, Xixian Liao and Maite Melero

*Vicomtech@WMT 2024: Shared Task on Translation into Low-Resource Languages of Spain*
David Ponce, Harritxu Gete and Thierry Etchegoyhen

**Saturday, November 16, 2024 (continued)**

*SJTU System Description for the WMT24 Low-Resource Languages of Spain Task*
Tianxiang Hu, Haoxiang Sun, Ruize Gao, Jialong Tang, Pei Zhang, Baosong Yang and Rui Wang

*Multilingual Transfer and Domain Adaptation for Low-Resource Languages of Spain*
Yuanchang Luo, Zhanglin Wu, Daimeng Wei, Hengchao Shang, Zongyao Li, Jiaxin Guo, Zhiqiang Rao, Shaojun Li, Jinlong Yang, Yuhao Xie, Zheng Jiawei, Bin Wei and Hao Yang

*TRIBBLE - TRanslating IBerian languages Based on Limited E-resources*
Igor Kuzmin, Piotr Przybyła, Euan Mcgill and Horacio Saggion

12:30 - 13:30    *Discourse-Level Literary Translation Task*

*CloudSheep System for WMT24 Discourse-Level Literary Translation*
Lisa Liu, Ryan Liu, Angela Tsai and Jingbo Shang

*Final Submission of SJTULoveFiction to Literary Task*
Haoxiang Sun, Tianxiang Hu, Ruize Gao, Jialong Tang, Pei Zhang, Baosong Yang and Rui Wang

*Context-aware and Style-related Incremental Decoding Framework for Discourse-Level Literary Translation*
Yuanchang Luo, Jiaxin Guo, Daimeng Wei, Hengchao Shang, Zongyao Li, Zhanglin Wu, Zhiqiang Rao, Shaojun Li, Jinlong Yang and Hao Yang

*NovelTrans: System for WMT24 Discourse-Level Literary Translation*
Yuchen Liu, Yutong Yao, Runzhe Zhan, Yuchu Lin and Derek F. Wong

*LinChance-NTU for Unconstrained WMT2024 Literary Translation*
Kechen Li, Yaotian Tao, Hongyi Huang and Tianbo Ji

12:30 - 13:30    *Chat Translation Task*

*Improving Context Usage for Translating Bilingual Customer Support Chat with Large Language Models*
Jose Pombal, Sweta Agrawal and André Martins

*Optimising LLM-Driven Machine Translation with Context-Aware Sliding Windows*
Xinye Yang, Yida Mu, Kalina Bontcheva and Xingyi Song

16:30 - 17:00      *Coffee Break*

17:00 - 18:00      *Session 11 — Research Paper Poster Session II*

*Post-edits Are Preferences Too*
Nathaniel Berger, Stefan Riezler, Miriam Exel and Matthias Huck

*Translating Step-by-Step: Decomposing the Translation Process for Improved Translation Quality of Long-Form Texts*
Eleftheria Briakou, Jiaming Luo, Colin Cherry and Markus Freitag

*Scaling Laws of Decoder-Only Models on the Multilingual Machine Translation Task*
Gaëtan Caillaut, Mariam Nakhlé, Raheel Qader, Jingshu Liu and Jean-Gabriel Barthélemy

*Shortcomings of LLMs for Low-Resource Translation: Retrieval and Understanding Are Both the Problem*
Sara Court and Micha Elsner

*Introducing the NewsPaLM MBR and QE Dataset: LLM-Generated High-Quality Parallel Data Outperforms Traditional Web-Crawled Data*
Mara Finkelstein, David Vilar and Markus Freitag

*Is Preference Alignment Always the Best Option to Enhance LLM-Based Translation? An Empirical Analysis*
Hippolyte Gisserot-Boukhlef, Ricardo Rei, Emmanuel Malherbe, Céline Hudelot, Pierre Colombo and Nuno M. Guerreiro

*Quality or Quantity? On Data Scale and Diversity in Adapting Large Language Models for Low-Resource Translation*
Vivek Iyer, Bhavitvya Malik, Pavel Stepachev, Pinzhen Chen, Barry Haddow and Alexandra Birch

*Efficient Technical Term Translation: A Knowledge Distillation Approach for Parenthetical Terminology Translation*
Myung Jiyoon, Jihyeon Park, Jungki Son, Kyungro Lee and Joohyung Han

*Assessing the Role of Imagery in Multimodal Machine Translation*
Nicholas Kashani Motlagh, Jim Davis, Jeremy Gwinnup, Grant Erdmann and Tim Anderson

*Error Span Annotation: A Balanced Approach for Human Evaluation of Machine Translation*
Tom Kocmi, Vilém Zouhar, Eleftherios Avramidis, Roman Grundkiewicz, Marzena Karpinska, Maja Popović, Mrinmaya Sachan and Mariya Shmatova

# Findings of the WMT24 General Machine Translation Shared Task: The LLM Era is Here but MT is Not Solved Yet

**Tom Kocmi**
Microsoft

**Eleftherios Avramidis**
DFKI

**Rachel Bawden**
Inria, Paris, France

**Ondřej Bojar**
Charles University

**Anton Dvorkovich**
Dubformer

**Christian Federmann**
Microsoft

**Mark Fishel**
University of Tartu

**Markus Freitag**
Google

**Thamme Gowda**
Microsoft

**Roman Grundkiewicz**
Microsoft

**Barry Haddow**
University of Edinburgh

**Marzena Karpinska**
UMass Amherst

**Philipp Koehn**
Johns Hopkins University

**Benjamin Marie**
The Kaitchup

**Christof Monz**
University of Amsterdam

**Kenton Murray**
JHU

**Masaaki Nagata**
NTT

**Martin Popel**
Charles University

**Maja Popović**
DCU & IU

**Mariya Shmatova**
Dubformer

**Steinþór Steingrímsson**
The Árni Magnússon Institute

**Vilém Zouhar**
ETH Zürich

## Abstract

This overview paper presents the results of the General Machine Translation Task organised as part of the 2024 Conference on Machine Translation (WMT). In the general MT task, participants were asked to build machine translation systems for any of 11 language pairs, to be evaluated on test sets consisting of three to five different domains. In addition to participating systems, we collected translations from 8 different large language models (LLMs) and 4 online translation providers. We evaluate system outputs with professional human annotators using a new protocol called Error Span Annotations (ESA).

## 1 Introduction

The Ninth Conference on Machine Translation (WMT24)[1] was held at EMNLP 2024 and hosted a number of shared tasks on various aspects of machine translation (MT). This conference built on 18 previous editions as a workshop or a conference (Koehn and Monz, 2006; Callison-Burch et al., 2007, 2008, 2009, 2010, 2011, 2012; Bojar et al., 2013, 2014, 2015, 2016, 2017, 2018; Barrault et al., 2019, 2020; Akhbardeh et al., 2021; Kocmi et al., 2022, 2023).

The goal of the General Machine Translation shared task is to explore the translation capabilities of current systems in diverse settings. We assess MT systems' ability to handle a broad range of translation and language use. How to test general MT performance is a research question in itself. Countless phenomena could be evaluated, the most important being:

- variety of domain (news, medicine, IT, patents, legal, social, gaming, etc.)
- style of text (formal or spoken language, fiction, technical reports, etc.)
- non-standard user-generated content (grammatical errors, code-switching, abbreviations, etc.)
- source modalities (text, speech, image)

Evaluating all phenomena is nearly impossible and creates numerous unforeseen problems. Therefore, we decided to simplify the problem and tackle only a selection of the phenomena.

We choose to evaluate different domains, this year focusing on the following ones: news, social/user-generated content, speech, literary, and educational. They were chosen to represent topics with different content styles and to be understandable for humans without specialist in-domain knowledge, thus not requiring specialized translators or human raters for evaluation. Due to limited access to monolingual data across all languages,

---

[1] www2.statmt.org/wmt24

the test set for each language direction contains at most four of the domains (Czech-Ukrainian uses different domains).

We evaluate a diverse set of languages pairs:

Czech→Ukrainian,
Japanese→Chinese – *new*,
English→Chinese,
English→Czech,
English→German,
English→Hindi,
English→Icelandic  – *new*,
English→Japanese,
English→Russian,
English→Spanish (Latin America) – *new*,
English→Ukrainian,

We newly test an audio modality as an additional source in the speech domain. Participants in this domain were provided with audio files and automatic speech recognized (ASR) text. Submission could use the original audio as an additional cleaner source modality instead of the provided ASR text.

In contrast to previous years, we adopt the Error Span Annotation protocol (Kocmi et al., 2024b), ESA for evaluation. This protocol, described in Section 6, combines aspects of DA (Graham et al., 2013) and MQM (Lommel et al., 2014).

In a shift towards document-level evaluation, we no longer provide source texts segmented into individual sentences. Instead, we keep all paragraphs intact and evaluated together.

Finally, this year's shared task included an increased number of test suites (Section 8) under the motto "Help us break the LLMs", focusing on revealing issues in the LLM translations from different perspectives, including a range of linguistic phenomena, idiomatic expressions and proper names, complex sentence structures, multiple domains, translation isochrony, speaker-listener gender resolution, prompt injection attacks, and gender-diverse, queer-inclusive content.

All General MT task submissions, sources, references and human judgements are available in the dedicated Github repository.[2] The interactive visualization and comparison of differences between systems can be browsed online on an interactive leaderboard[3] using MT-ComparEval (Klejch et al., 2015; Sudarikov et al., 2016).

The structure of the paper is as follows. We describe the process of collecting, cleaning and translating the test sets in Section 2 followed by a summary of the permitted training data and pretrained models for the constrained track in Section 3. We list all submitted systems in Section 4. Automatic evaluation is described in Section 5. The human evaluation approach of ESA is described in Section 6. The main results can be found in Section 7 and their extended version in Appendix D. Finally, Section 8 describes the test suites and summarises their conclusions.

**Findings of the WMT2024 General MT Task**
Across the evaluation conditions, we observe the following:

- The best systems for English→Spanish produced close to flawless translations making it the easiest language pair (Section 6.4).

- The speech domain is the most challenging domain (likely due to the ASR) while the other three domains (news, literary, social) are similarly difficult (Section 6.4).

- Human references are in the winning cluster in 7 out of 11 language pairs. For one of the remaining 4 pairs (English→Hindi), we know the reference quality was suboptimal. This suggest that ESA protocol works well in our setting.

- ESA produced 37% more clusters than DA+SQM while using only half the number of human annotations (Section 6.5).

- The best performing system in the open and constrained system category is IOL-Research (wins 10 language pairs in this category). The best performing participating system is Unbabel-Tower70B, which wins in 8 language pairs. And the best performing system in general is Claude-3.5-Sonnet winning in 9 language pairs.

- Automatic scores are biased; although Unbabel-Tower70B placed first across all languages in automatic ranking it didn't perform as the winning system across the board of human evaluation. This is likely because we used the same metric (COMET) for automatic ranking as the system used during MBR highlighting the issue of hill-climbing on automatic metrics.

- We got a total of 28 participants, which nearly 50% more than last year. Most of the participants use an LLM as a part of their system, generally by fine-tuning it.

- Quality estimation metrics with fixed score for perfect translation and interpretable delta are promising for checking the quality of standalone human references.

## 2 Test Data

In this section, we describe the data collection process (Section 2.1), and the production of human reference translations (Section 2.3).

### 2.1 Collecting test data

As in previous years, the test sets consist of unseen translations created specifically for the shared task and released publicly to be used as translation benchmarks. Our aim was to collect public domain or open-licence source data covering a range of domains, and we also focused on using as recent data as possible to limit possible contamination (particularly relevant when using LLMs).

We chose four main domains from which to collect data (news, literary, speech and social), although we were not able to collect data in all domains for all three source languages (no social domain data is provided for Japanese→Chinese and Czech→Ukrainian data was collected separately, comprising news data and four other separate domains). For all language pairs, the test sets are "source-original", meaning that the text was originally written in the source language, which is then manually translated into the target languages. This is important to avoid "translationese" in the source texts, which can have a negative impact on evaluation accuracy (Toral et al., 2018; Freitag et al., 2019; Läubli et al., 2020; Graham et al., 2020). We aimed for a certain number of *tokens*[4] in each domain rather than a certain number of *sentences* (as in previous years) to better balance the domains and also because the document-level focus this year allowed avoid manual sentence splitting. We aimed for approximately 10,000 tokens per domain, adjusting this figure in cases where not all domains could be covered (this is notably the case for Japanese→Chinese, where the other domains are up-sampled to account for us not being able to provide data in the social domain). Basic statistics of each subdomain are given in Table 1.[5]

Note that by default, when languages are mentioned in this section, this refers to the source language of the texts.

**News domain** This domain contains data prepared in the same way as in previous years (Kocmi et al., 2023). We collected news articles from January 2024 extracted from online news sites, preserving document boundaries. We expect that news domain text will generally be of high quality.

For Japanese, the total amount of text data was determined by the number of characters since Japanese does not put spaces between words. Using the WMT23 Japanese test set and its translation into English, we found the ratio of the number of Japanese characters to English words was 2 to 1. Since the English news test set consisted of 8K words, we started making a Japanese news test set with a goal of 16K characters. After discovering that the Japanese social domain was unavailable, we set this goal to 24K characters.

**Literary Domain** The English source texts were manually selected from Archive of Our Own,[6] focusing on recent, high-quality stories.[7] The stories were divided into 1000-word segments, ensuring the preservation of entire paragraphs. In total, we obtained data from four stories (8K words).[8]

For the Japanese source texts, we selected five novels recently made public on Aozora Bunko,[9] a website that digitizes and publishes Japanese literary works whose copyright has expired. To maintain consistency with the English dataset, we tokenized the Japanese novels using MeCab (Kudo, 2005) and divided them into segments of up to 1000 tokens, while preserving paragraph boundaries. The final size of the Japanese literary test set was 15 chunks (22K characters).

**Speech domain** The speech data corpus was compiled from a diverse range of YouTube videos licensed under Creative Commons. These sources encompassed various domains, including documentaries, instructional (DIY) videos, tutorials, travel blogs, and film content. For this part of the test set, segments from 166 videos were selected and processed through automated speech recognition (ASR) systems. For the English-language source

---

[4]For Japanese source texts, we choose to use a certain number of characters, since words are not space-separated.

[5]Texts are sentence-segmented and tokenised using the language-specific Spacy models (Honnibal and Montani, 2017) optimised for accuracy where available. For Czech, we use the multilingual Spacy model, as a language-specific

one is not available. Note that statistics, particularly for this language, are approximate.

[6]archiveofourown.org

[7]Texts were published between February and April 2024.

[8]For each, we select first two chunks of up to 1000 words.

[9]aozora.gr.jp

| Language pair | News | Literary | Speech | Social | Education | Official | Personal | Voice |
|---|---|---|---|---|---|---|---|---|
| | | | | #tokens | | | | |
| English→* | 9,268 | 9,601 | 9,611 | 9,829 | - | - | - | - |
| Japanese→Chinese | 14,896 | 14,541 | 11,025 | - | - | - | - | - |
| Czech→Ukrainian | 7,996 | - | - | - | 7,825 | 6,029 | 6,846 | 5,305 |
| | | | | #segs (% of total #segs for language pair) | | | | |
| English→* | 149 (14.9) | 206 (20.7) | 111 (11.1) | 531 (53.3) | - | - | - | - |
| Japanese→Chinese | 269 (37.3) | 316 (43.8) | 136 (18.9) | - | - | - | - | - |
| Czech→Ukrainian | 175 (7.6) | - | - | - | 1160 (50.1) | 243 (10.5) | 323 (13.9) | 415 (17.9) |
| | | | | #docs (#segments/doc) | | | | |
| English→* | 17 (8.8) | 8 (25.8) | 111 (1.0) | 34 (15.6) | - | - | - | - |
| Japanese→Chinese | 45 (6.0) | 15 (21.1) | 136 (1.0) | - | - | - | - | - |
| Czech→Ukrainian | 23 (7.6) | - | - | - | 166 (7.0) | 23 (10.6) | 29 (11.1) | 61 (6.8) |
| | | | | #sents (#sents/doc) | | | | |
| English→* | 333 (19.6) | 607 (75.9) | 685 (6.2) | 759 (22.3) | - | - | - | - |
| Japanese→Chinese | 634 (14.1) | 875 (58.3) | 332 (2.4) | - | - | - | - | - |
| Czech→Ukrainian | 439 (19.1) | - | - | - | 1166 (7.0) | 412 (17.9) | 571 (19.7) | 462 (7.6) |
| | | | | Type-token ratio of source texts | | | | |
| English→* | 0.30 | 0.23 | 0.24 | 0.27 | - | - | - | - |
| Japanese→Chinese | 0.22 | 0.20 | 0.19 | - | - | - | - | - |
| Czech→Ukrainian | 0.46 | - | - | - | 0.39 | 0.45 | 0.34 | 0.37 |

**Table 1:** Basic statistics concerning the subdomains of each test set. Statistics are calculated on the source side. Sentence segmentation and tokenisation are carried out automatically as described in Footnote 5.

material, we used the proprietary Dubformer engine developed in-house. Japanese-language content was processed using the Whisper ASR system (Radford et al., 2022).

For Japanese, We selected 136 segments from 56 YouTube videos. They include both monologues and dialogues, as well as a variety of speakers, both men and women, adults and children. Video content includes press conferences, interviews, cooking recipes, travel vlogs, DIY videos, tutorials, product reviews, etc. We decided the total amount of speech data based on the number of characters transcribed. We started creating the data with a target of 16K characters and eventually ended up with 18K characters.

**Social domain** The social domain data is sourced using the Mastodon Social API.[10] Mastodon is a federated social network that is compatible with the W3C standard ActivityPub (Webber et al., 2018). Users publish short-form content known as "toots", with the possibility of replying to other toots to form threads. We decided to use the original server, `mastodon.social` because of its large community and publicly available toots.

We collected data in the first four months of 2024, using the reported language ID label to target the source languages of interest. Unfortunately,

there were too few good quality posts for Czech and Japanese, and we therefore only release social domain data for English.

Given the document-level nature of the task this year, our aim was to collect threads comprising multiple toots. Our sourcing therefore involved regularly scraping random toots from the previous hour but also identifying and scraping any missing toots that made up threads only partially sourced (identified using the 'in_reply_to_id' attribute of already sourced toots). To avoid spam and uninformative toots, we removed empty toots, texts that appeared several times (probable spam), texts from accounts that produced a large number of toots overall (we set this to 100 for a total of 1.5M toots scraped) and from accounts we heuristically identified as being news outlets or bots (containing the keywords 'bot', 'news', 'weather', 'sports', 'feeds' or 'press' in their handle). We created threads from the individual toots and manually selected threads of interest from threads of minimum 2 and maximum 100 toots. Our selection was based on having a diverse range of topics and targeting those characteristic of non-standard user-generated content.

The selected documents contain between 5 and 76 segments of text, each segment corresponding either to a whole toot or a line of text within a toot (depending on whether the toot contained newlines, i.e. there is no distinction between newlines indi-

cating a boundary between two toots and a newline within a toot). Each segment can therefore contain one or several sentences, depending on the original composition of the toots.

**Czech and Ukrainian source texts**   Source texts for Czech→Ukrainian translation included the news domain as described above, Educational domain collected from online exercises and three domains (Personal, Official and Voice) from texts collected through Charles Translator.[11] The Charles Translator mobile app supports voice input, which is converted to text using Google ASR. The texts collected this way were marked as the Voice domain. The remaining Czech inputs from the Charles Translator service were classified either as Official (formal communication) or Personal (personal communication, usually between a Czech and Ukrainian speaker).

The texts were filtered and pseudonymized in the same way as in the last two years (Kocmi et al., 2022). For example we asked the annotators not to delete or fix noisy inputs as long as they are comprehensible. The only exception was the voice domain, where the source texts were post-edited to fix ASR errors, including punctuation and casing.

The Educational domain includes selected exercises from an online portal *Škola s nadhledem*[12] for elementary-school students from various subjects (chemistry, geography, Czech language, etc.). Some segments are not full sentences but short phrases. The reference translations for this domain were created by professional translators within the EdUKate project.

## 2.2   Comparison between Domains

Due to the change to document-level translation this year, for each language direction, we measured the amount of text per domain by counting tokens, aiming for approximately the same number of tokens per domain (see Table 1 for statistics of the different domains). In one sense, this means that the amount of textual content is roughly balanced per domain, as opposed to taking the same number of sentences per domain, which would result in domains with longer sentence lengths (e.g. news or literary) being over-represented with respect to domains with shorter sentences (e.g. social). However, it is worth noting that the nature of documents, in terms of their length and structure, differs greatly

depending on the domain. This can be exemplified at its most extreme by a comparison between the literary, social and speech domains for from-English language directions.

The literary domain has only 8 documents, each one containing a large number of segments (25.8 on average), with each segment containing an average of 75.9 sentences. A document represents an extract from a longer literary text and each segment represents a paragraph of text.

The speech domain is represented by a larger number of documents (111), each one containing a single segment, composed of an average of 6.2 sentences. Each document in this case corresponds to a short dialogue, provided without segmentation into dialogue turns.

The social domain is represented by a fair number of documents (34 in total), but the composition is very different from the other domains, as we made a choice to preserve the structure of the initial posts (new-line separated text is represented by multiple segments) and of the thread itself (separate posts are separate segments). This has the advantage of preserving post boundaries and formatting choices, but has the disadvantage of creating a large number of individual segments (531 in total, compared to 206 for the literary domain and 111 for the speech domain), each containing few sentences. This has two main consequences: (i) if segments are handled individually by systems, most sentences will be handled with little context, since the other sentences appear in separate segments, (ii) in terms of the overall number of segments evaluated in the human evaluation (see Section 6), the social domain represents over half of the total number of evaluated segments (53.3% compared to 20% for the literary domain and only 11.1% for the speech domain). This has consequences for the calculation of macro-average scores when computing human rankings, as discussed in Section 7.1. The formatting choice could be rethought for future years, although would have to take into account the particularities of non-standard text in order to not introduce extra noise (e.g. concatenating newline-separated sentences would have to take into account the potential lack of end-of-sentence punctuation, but it would also have to take into account instances where newlines are used with a single sentence for purely visualisation purposes. A possible solution would be to allow a linebreak symbol such as `<br/>` to appear in the segments.

---

[11]translator.cuni.cz
[12]skolasnadhledem.cz

## 2.3 Human References

The test sets were translated by professional translation agencies, according to the translation brief shown in Appendix C. Different partners sponsored each language pair and various translation agencies were therefore used, which could affect the differences and quality of translations.

The quality of human references is critical especially for reference-based metrics (Freitag et al., 2023), and getting high quality translations is challenging despite the use of professional translators. Therefore, we propose to use a quality estimation metric to assess the quality of translation. We need a metric whose score is interpretable in an absolute way, i.e. a metric that generates a fixed score for perfect translations (such as 0) and has an understandable delta (for example -1 means a single minor error as in MQM-based metrics). For that reason, we chose a GPT-4-based implementation of GEMBA-MQM (Kocmi and Federmann, 2023).

Table 2 shows the GEMBA scores for individual domains together with the ESA human cluster that was obtained a few months later in our official manual evaluation.

The two target languages with the lowest GEMBA scores were Russian and Hindi. The vendor providing Russian translations improved the initial quality of translations after being presented with the GEMBA scores. On the other hand, the vendor providing Hindi translators claimed that the translations were flawless.

When we compare the average GEMBA score to human rank in Table 2, we can see that human reference is ranked in the top cluster for all languages except of Hindi, Ukrainian, and Chinese. While the GEMBA score did not reflect lower quality of Ukrainian, its low score for Hindi was confirmed by ESA. This shows that using quality estimation metrics is a possible way of assessing the quality of human translations, although better approaches needs to be developed.

## 2.4 Test Suites

In addition to the test sets of the regular domains, the test sets given to the system participants were augmented with several *test suites*, i.e. custom-made test sets focusing on particular aspects of MT translation. The test suites were contributed and evaluated by test suite providers as part of a decentralized sub-task, detailed in Section 8. Across all language pairs of the shared task, test suites

| | Literary | News | Social | Speech | Avg. | Hum. |
|---|---|---|---|---|---|---|
| En.→Czech | -2.4 | -2.0 | -1.9 | -1.8 | -2.03 | 1 |
| En.→German$_A$ | -2.1 | -2.0 | -2.3 | -2.3 | -2.18 | 1 |
| En.→German$_B$ | -2.7 | -0.8 | -1.7 | -2.0 | -1.80 | 1 |
| En.→Spanish | -1.1 | -1.6 | -1.2 | -1.6 | -1.38 | 1 |
| En.→Hindi | -3.4 | -4.5 | -2.5 | -2.9 | -3.33 | 3 |
| En.→Icelandic | -2.6 | -0.8 | -1.9 | -1.4 | -1.68 | 1 |
| En.→Japanese | -1.7 | -1.6 | -1.7 | -1.7 | -1.68 | 1 |
| En.→Russian | -2.6 | -2.8 | -2.5 | -2.3 | -2.55 | 1 |
| En.→Ukrainian | -1.8 | -1.0 | -2.0 | -2.3 | -1.78 | 3 |
| En.→Chinese | -3.1 | -1.7 | -2.8 | -2.2 | -2.45 | 2 |

**Table 2:** GEMBA-MQM score for human references. The first four columns are scores for individual domains, the fifth column is the average. The last column is the human cluster assigned with ESA protocol. Czech→Ukrainian is not included because of different domains and source data.

contributed 718,598 source test segments (detailed numbers can be found in Table 9).

## 3 Training Data

Similar to the previous years, we provide a selection of parallel and monolingual corpora for model training. The provenance and statistics of the selected parallel datasets are provided in the appendix in Table 10 and Table 11. Specifically, our parallel data selection include large multilingual corpora such as Europarl-v10 (Koehn, 2005), Paracrawl-v9 (Bañón et al., 2020), CommonCrawl, NewsCommentary-v18.1, WikiTitles-v3, WikiMatrix (Schwenk et al., 2021), TildeCorpus (Rozis and Skadiņš, 2017), OPUS (Tiedemann, 2012), CCAligned (El-Kishky et al., 2020), UN Parallel Corpus (Ziemski et al., 2016), and language-specific corpora such as CzEng v2.0 (Kocmi et al., 2020), YandexCorpus,[13] ELRC EU Acts, JParaCrawl (Morishita et al., 2020), Japanese-English Subtitle Corpus (Pryzant et al., 2018), KFTT(Neubig, 2011), TED (Cettolo et al., 2012), and back-translated news.

Links for downloading these datasets were provided on the task web page. In addition, we have automated the data preparation pipeline using MTDATA (Gowda et al., 2021).[14] MTDATA downloads all the mentioned datasets, except CzEng v2.0, which required user authentication. This year's monolingual data include the following: News Crawl, News Discussions, News Commentary, CommonCrawl, Europarl-v10 (Koehn, 2005), Extended CommonCrawl (Conneau et al., 2020), Leipzig Corpora (Goldhahn et al., 2012), UberText and Legal Ukrainian.

[13] github.com/mashashma/WMT2022-data
[14] statmt.org/wmt24/mtdata

| System | Language pairs | Architecture | Strategy |
|---|---|---|---|
| AIST-AIRC (Rikters and Miwa, 2024) | en→de, en→ja | dec, enc-dec, MEGA | sentence |
| AMI (Jasonarson et al., 2024) | en→is | enc-dec | hybrid |
| BJFU-LPT | cs→uk | – | – |
| CUNI-DocTransformer (Hrabal et al., 2024) | en→cs | enc-dec | paragraph |
| CUNI-Transformer (Hrabal et al., 2024) | cs→uk, en→cs | enc-dec | sentence |
| CUNI-DS (Semin and Bojar, 2024) | en→ru | dec | sentence |
| CUNI-GA (Hrabal et al., 2024) | en→cs | enc-dec | sentence |
| CUNI-MH (Hrabal et al., 2024) | en→cs | dec | sentence |
| CUNI-NL (Hrabal et al., 2024) | en→de | dec | sentence |
| CycleL (Dreano et al., 2024) | **All language pairs** | CycleGAN | – |
| CycleL2 (Dreano et al., 2024) | en→cs, en→de, en→ru, en→zh | CycleGAN | – |
| DLUT-GTCOM (Zong et al., 2024) | en→ja, ja→zh | enc-dec | – |
| Dubformer | en→de, en→es, en→is, en→ru, en→uk | – | – |
| HW-TSC (Wu et al., 2024) | en→zh | hybrid | sentence |
| IKUN (Liao et al., 2024) | **All language pairs** | dec | sentence |
| IKUN-C (Liao et al., 2024) | **All language pairs** | dec | sentence |
| IOL-Research (Zhang, 2024) | **All language pairs** | dec | paragraph |
| MSLC (Larkin et al., 2024) | en→de, en→es, ja→zh | enc-dec | sentence |
| NTTSU (Kondo et al., 2024) | en→ja, ja→zh | hybrid | paragraph |
| NVIDIA-NeMo | All except cs→uk, en→is and ja→zh | dec | paragraph |
| Occiglot (Avramidis et al., 2024) | en→de, en→es | dec | – |
| SCIR-MT (Li et al., 2024) | en→cs | dec | – |
| Team-J (Kudo et al., 2024) | en→ja, ja→zh | hybrid | hybrid |
| TranssionMT | All except en→ja, en→zh and ja→zh | enc-dec | – |
| TSU-HITs (Mynka and Mikhaylovskiy, 2024) | en→cs, en→de, en→es, en→is, en→ru | ddm | sentence |
| Unbabel-Tower70B (Rei et al., 2024) | **All language pairs** | dec | paragraph |
| UvA-MT (Tan et al., 2024) | en→ja, en→zh, ja→zh | hybrid | hybrid |
| YandexSubtitles (Elshin et al., 2024) | en→ru | dec | paragraph |
| Aya23 (Aryabumi et al., 2024) | **All language pairs** | dec | paragraph |
| Claude-3.5 | **All language pairs** | dec | paragraph |
| CommandR+ | **All language pairs** | dec | paragraph |
| GPT-4 (OpenAI, 2024) | **All language pairs** | dec | paragraph |
| Gemini-1.5-Pro (Team, 2024a) | All except en→is | dec | paragraph |
| Llama3-70B (Team, 2024b) | **All language pairs** | dec | paragraph |
| Mistral-Large (Jiang et al., 2023) | **All language pairs** | dec | paragraph |
| Phi-3-Medium (Team, 2024c) | **All language pairs** | dec | paragraph |
| ONLINE-A | **All language pairs** | – | – |
| ONLINE-B | **All language pairs** | – | – |
| ONLINE-G | **All language pairs** | – | – |
| ONLINE-W | All except en→is and en→hi | – | – |

**Table 3:** Participating submissions in the General MT shared task. The top section covers the externally contributed submissions, the middle section lists the language models added by us and the lower section covers the online systems. Online system translations were not submitted by their respective companies but were obtained by us, and are therefore anonymized in a fashion consistent with previous editions of the task. Row coloring shows closed-track (dark gray), open-track (light gray) and constrained (white background) submissions. The Architecture column shows whether the submission used decoder-only language models (dec), sequence-to-sequence (enc-dec), hybrid between dec and enc-dec or other architectures. The Strategy column shows the approach used to handling paragraph-level test data: sentence-level training and translation (sentence), paragraph-level training and translation (paragraph), hybrid between both (hybrid). Some values are unknown (–) due to missing information or submission papers.

## 4 System Submissions

This year, we received a total of 105 primary submissions from 28 participants. The increase in number of participants from last year's 19 can be explained by the shift in the field and the ease with which LLMs can be fine-tuned. The increased number of primary submissions can be explained by the fact that most submissions are multilingual and therefore cover many translation directions.

In the same manner as previous years, we also collected translations from online MT systems for all language pairs. Online system outputs come from four public MT services and were anonymized as ONLINE-{A,B,G,W}, which resulted in further 42 system outputs. Finally, we added contrastive translations by 8 LLMs, which included closed commercial products (such as GPT-4) and open models (such as Llama3). This resulted in 95 more submissions, with the total number of submissions being 242.

All participating systems are listed in Table 3. Appendix B provides more detailed short descriptions of the submitted systems, as provided by the authors at submission time. Section 4.1 discusses the general trends in chosen architectures and approaches to paragraph-level translation. Section 4.2 presents details on LLM benchmark usage in the task. Section 4.3 describes the different tracks to which participants could submit outputs: constrained, open and closed track. Section 4.4 describes the submission system setup.

### 4.1 Architectures and Strategies

In addition to a reference to a description paper (if one was provided), the submission name and the list of language pairs covered, Table 3 includes columns for the architecture and strategy used to approach the task of paragraph-level translation. If we compare the frequency of usage of different architectures between the external participants (i.e,. excluding benchmarking LLMs and online systems), we can see that:

- 11 participants train decoder-only language models (*dec* in Table 3)

- 7 participants train encoder-decoder seq2seq transformer models (*enc-dec*)

- 4 participants use a hybrid of the decoder-only and encoder-decoder architectures (*hybrid*)

- 3 alternative architectures were used: MEGA (Ma et al., 2023) in AIST-AIRC, CycleGAN (Zhu et al., 2017) in CycleL and discrete diffusion models in TSU-HITs.

Not all description papers specified the strategy used to translate the test set paragraphs. Of those who did, 5 submissions approached it by explicitly training paragraph-level translation systems, while 11 submissions translated single sentences after sentence-splitting the paragraph. 3 submissions described a hybrid approach of, for example, translating single sentences but automatically post-editing at the paragraph level. Several papers do not mention the strategy at all. We plan to address this lack of information in future WMT editions by requesting that the information be provided at submission time.

Interestingly, the paragraph-level approach is not limited to a single architecture: for instance, the CUNI-DocTransformer team uses an encoder-decoder approach, but trains it on paragraph-level parallel data, which includes synthetic data. There are examples to the contrary: several submissions fine-tune a decoder-only language model, but apply it to translate single sentences (IKUN, AIST-AIRC, several CUNI submissions).

Finally, almost all submissions used an LLM as a part of their setup. The most common use is fine-tuning of a pretrained model, most often LLama. Other uses of LLMs are for generating or cleaning up training data with an LLM (Jasonarson et al., 2024) or using an LLM for automatic post-editing (Tan et al., 2024).

### 4.2 LLM Benchmark

Over the last year, many new LLMs claimed multilingual and translation capabilities. However, there is no systematic and reliable MT evaluation of the most popular LLMs using the same setup on blind test sets. We therefore decided to collect the translations of LLMs ourselves.

We design unified code for collecting the translations in an identical setup for all LLMs. We used a 3-shot approach, where three fixed examples are taken from the past WMT test sets. We set the temperature to zero to avoid introducing randomness into the process.[15]

We evaluated most of the popular LLMs, both closed-source models and those with open

---

| Language model | Input tok. | Output tok. | Cost |
|---|---|---|---|
| Aya23 | 4.4 M | 0.7 M | 4.1 $ |
| Claude-3.5 | 5.5 M | 1.0 M | 31.9 $ |
| CommandR-plus | 4.4 M | 0.7 M | 23.4 $ |
| Gemini-1.5-Pro | 3.9 M | 0.6 M | 40.3 $ |
| GPT-4 | 5.9 M | 1.0 M | 240.4 $ |
| Llama3-70B | 5.0 M | 0.7 M | 5.1 $ |
| Mistral-Large | 6.0 M | 1.1 M | 37.0 $ |
| Phi-3-Medium | 5.9 M | 1.1 M | 4.5 $ |

**Table 4:** Number of input and output tokens and estimated pricing for translating the full WMT24 test set without test suites. The Gemini model refused to translate Icelandic, and the estimate is therefore lower. Pricing for the open models Aya23 and Llama3 was estimated via together.ai.

weights. Specifically, we collect translations from Aya23, Claude-3.5-Sonnet, Command R+, GPT-4, Gemini-1.5-Pro, Llama3-70B, Mistral-Large, Nvidia-NeMo and Phi-3-Medium. As most of the models do not claim multilingual capabilities for all languages covered, we looked into the original reports for these LLMs to see which languages are claimed to be supported. We check if both source and target language are mentioned or evaluated in any of their multilingual settings. We mark LLMs that do not officially claim a support for a given language with the § symbol in the tables. However, to avoid selection bias, we collect translations for all languages for all LLMs, even those not officially claimed to be supported.

We collect all translations via the API of the respective services, and all data was collected during the submission week. Table 4 shows the number of input and output tokens as segmented via the models' internal tokenizers. The estimated cost is for the whole test set without test suites. Note that the prices for more recent GPT models are significantly lower.

### 4.3 Constrained, Open, and Closed Tracks

We distinguish three types of MT systems participating in the shared task: constrained, open and closed systems. The main idea is to level the field for different setups. For the constrained setup, we only allow specific training data and pretrained models from a specified list. Open systems are those developed using publicly available data or models. The final group of closed systems corresponds to all other systems that are built at least partly with a non-replicable setup.

- **Constrained systems** are those using only the specifically allowed training data (see Section 3) and the following pretrained models: Llama-2-

7B, Llama-2-13B, Mistral-7B, mBART, BERT, RoBERTa, XLM-RoBERTa, sBERT, LaBSE. Constrained systems may use any publicly available metric that was evaluated in past WMT Metrics shared tasks (e.g. COMET or Bleurt) and any basic linguistic tools (e.g. taggers, parsers, morphology analyzers).

- **Open systems** (marked in tables with a light gray background) are limited to using software, data and models that are freely available for research purposes, so that the subsequent work could be replicated by a research group.

- **Closed systems** (marked with dark gray) correspond to all the remaining (fully automatic) systems, with no limitations imposed on their training data (all ONLINE systems and LLMs released without binaries fit into this category).

### 4.4 OCELoT

We used the open-source OCELoT platform[16] to collect system submissions again this year. As in previous years, only registered and verified teams with correct contact information were allowed to submit their system outputs and each verified team was limited to 7 submissions per test set. Submissions on leaderboards with BLEU (Papineni et al., 2002) and CHRF (Popović, 2015) scores from SacreBLEU (Post, 2018) were displayed anonymously to avoid publishing rankings based on automatic scores during the submission period. Until one week after the submission period, teams could select a single primary submission per test set, specify if the primary submission followed a constrained, open or closed system setting, and submit a system description paper abstract. These were mandatory for a system submission to be included in the human evaluation campaign.

## 5 Automatic Evaluation

This year, we received an unusually high number of submitted systems and we were not able to provide manual evaluation for all of them. Therefore, we decided to use automatic metrics to preselect the best performing systems with a method we call AutoRank, which is based on two different metrics:

- MetricX-23-XL (Juraska et al., 2023), a reference-based metric built on top of the mT5 model (Xue, 2020).

---

[16]github.com/AppraiseDev/OCELoT

**(a)** Excerpt of two segments from a larger document. In the first segment, the name *"Kayel"* is omitted which is a major error. In the second segment, there are many minor errors.

**(b)** Example of a video to text translation with several minor errors. The annotator can control the video player.

**Figure 1:** Two screenshots of ESA (Kocmi et al., 2024b) and the annotator instructions. ESA shows multiple segments within a document at once as well as video sources. After marking the individual error spans, the annotator assigns the final segment score from 0 to 100. The tool is implemented in Appraise (Federmann, 2018).

- CometKiwi-DA-XL (Rei et al., 2023), a quality estimation metric built on the XLM-R XL model (Conneau, 2019).

Both metrics are top performing metrics (Freitag et al., 2023), and we intentionally select two distinct metrics (different underlying pretrained systems and architectures) to minimize their bias and potential problems. Although quality estimation is on average slightly worse than reference-based evaluation, it helps us to avoid a potential reference bias when human references are suboptimal (Freitag et al., 2023). Multilingual quality estimation can be fooled when the translation is in the incorrect language, which the reference-based metric will penalize.

To compute MetricX, we used the official implementation[17] and the "google/metricx-23-xl-v2p0" model. MetricX produces scores at the segment level. To obtain scores at the system level, we averaged the segment scores. To compute CometKiwi scores, we used the official implementation[18] with the "Unbabel/wmt23-cometkiwi-da-xl" model, a reference-free model, taking the translation hypothesis and the source segment as inputs. COMET can produce system-level scores so we use them directly.

To merge the two metrics, we first linearly scale the scores of each metric to a range between 1 and

the number of systems for a given language pair. We then average both normalized scores to reach the final automatic ranking, which we refer to as AutoRank. We provide a Jupyter notebook in the WMT24 repository to reproduce the scores.[19]

### 5.1 Selecting Systems for Human Evaluation

When selecting the systems for human evaluation, we prioritize open and constrained systems while penalizing closed systems. We select a subset of 10 to 15 systems per language pair based on AutoRank and following two rules. First, we exclude closed systems that are not among the first third of all systems and we exclude open systems that are not among the top two thirds of all systems. Second, motivated by several very low-performing systems, we also define a hard cutoff point. After this point we do not evaluate any system from any category. The cutoff point is selected as the first gap between two neighboring system's ranks larger than 1.5 of AutoRank. This decision was discussed and published in more detail in Kocmi et al. (2024a).

## 6 Human Evaluation

This year's human evaluation features Error Span Annotation (ESA; Kocmi et al., 2024b) for most languages. For Japanese→Chinese and

---

| Language pairs | Annotators' profile | Tool |
|---|---|---|
| English→Chinese/Japanese/ Hindi/Spanish | Microsoft annotators — bilingual target-language native speakers, professional translators or linguists, experienced in machine translation evaluation. | Appraise ESA |
| Czech→Ukrainian English→Czech | ÚFAL Charles University annotators — linguists, annotators, researchers, and students who were native speakers in the target language and had a very high proficiency in English (for English→Czech) and good knowledge of Czech (for Czech→Ukrainian). | Appraise ESA |
| English→Ukrainian/ Russian/Spanish | Toloka AI paid expert crowd — Bilingual native target-language speakers who were high-performing on the platform. | Appraise ESA |
| English→Icelandic | The Árni Magnússon Institute for Icelandic Studies annotators — bilingual target-language native speakers, paid translators with 10–25 years of experience in Icelandic↔English translation. | Appraise ESA |
| English↔German Japanese→Chinese | Campaign managed by the 2024 metrics shared task. | Google MQM |

**Table 5:** Annotators' profiles and annotation tools for each language pair in the human evaluation. English→Spanish was split between Microsoft and Toloka AI. All annotators were paid a fair wage in their respective countries.

English→German, we rely on the evaluation campaign from the metrics shared task 2024 (Freitag et al., 2024), which uses Multidimensional Quality Metrics (MQM; Lommel et al., 2014).

**Annotation Protocol.** ESA is based on highlighting/marking errors without classifying them into different error types (Kreutzer et al., 2020; Popović, 2020) and represents a compromise between overall scoring (such as direct assessment, DA; Graham et al. 2013) and error classification (such as MQM; Lommel et al. 2014).

The annotators (professional translators but not experts in MQM/ESA-style annotations) were asked to mark each error as well as its severity, "Minor" or "Major", as in Kocmi et al. (2024b); Popović (2020). In addition, the annotators were also asked to assign a score from 0 to 100, similar to DA, to the whole annotation segments (usually a sentence or a paragraph). However, the ESA score should be more robust than DA alone because the annotators are primed by the highlighted errors at the time of the scoring.

The interface is shown in Figure 1 with annotator instructions and other changes from the original implementation by Kocmi et al. (2024b) given in Appendix A. At the start of annotation, each annotator was exposed to an interactive tutorial where they were asked to interact with the system. The length of the context given to the annotators varies depending on the domain, ranging from one to ten sentences, as discussed in Section 6.1. The source for the speech domain is a video which is shown in

| Language pair | Systems | Annotators | | |
|---|---|---|---|---|
| | | Duplication | | Assess./system |
| Cs→Uk | 11 | 1.0 | 14 | 1299 |
| En→Czech | 15 | 1.3 | 20 | 751 |
| En→Spanish | 13 | 1.0 | 14 | 370 |
| En→Hindi | 10 | 1.3 | 15 | 775 |
| En→Icelandic | 10 | 1.0 | 4 | 376 |
| En→Japanese | 12 | 1.5 | 14 | 1212 |
| En→Russian | 13 | 1.0 | 7 | 370 |
| En→Ukrainian | 10 | 1.0 | 8 | 376 |
| En→Chinese | 12 | 1.5 | 12 | 1217 |

**Table 6:** Number systems, annotators, and number of assessments per system in a language pair. Duplication of $d$ means that each segment is annotated by $d$ annotators. All language pairs had 649 segments over 170 documents except for Czech→Ukrainian which had 1954 segments over 302 documents. In total we collected 57k segment-level annotations. English→German and Japanese→Chinese are managed by the metrics shared task 2024.

a native HTML video player.

The output of the ESA annotation is a list of errors and their severity (minor or major) and the final score from 0 to 100 for each segment.

**Human Annotators** Campaigns for different language pairs were managed by various vendors, as described in Table 6. In all cases, professional translators-cum-annotators are used. This is an increasingly strict requirement given the high quality of MT systems, which requires more expert annotators.

### 6.1 Data Preparation

**Document Filtering.** In our setup, all systems for a given language pair are evaluated on the same set of segments. On average, we start with 1092

lines per system, encompassing 184 documents. However, the distribution of document lengths is unbalanced. The majority of the documents (104) consist of just a single line, which is almost exclusively due to video translation segments (103), where each "document" contains strictly one segment. On average, 33 documents per language contain more than 10 segments. We limit these documents to the first 10 segments, motivated by the difficulty of annotating very long documents while maintaining relevant context in mind. After this adjustment, we arrive at an average of 744 lines per system. An overview is shown in Table 6.

**Workload balancing** We use the term "task" as a contained unit of 100 annotation segments. Each annotator is usually assigned to multiple tasks. This 100-segment constraint was kept for historical reasons and will be dropped in future iterations. In order to make it so that each task contains a comparable amount of work, we attempt to balance the number of words in each task to be as constant as possible.

For each task, we show a tutorial at the beginning consisting of 2 documents with 6 segments in total. The tutorial is for German→English translation but does not require any knowledge of German. Finally, we reserve 12 segments for quality control (Section 6.2) in each task. The resulting 82 segments are filled with full documents as much as possible. If that is not possible (i.e., because the next document is too long), a random document is drawn that is either duplicated or incomplete, in order to fill the 100 segments.

**Annotation waves** In order for a segment to be useful in the evaluation, we require that translations by all systems are evaluated. We therefore split (at the document level) the translated data for each language into two "waves", each of which covers a distinct set of source segments. The vendors are instructed to start the second campaign only after the first one is fully complete.

For some language pairs, the vendors finished early. In this case, we prepared an extra two waves, with a different coverage split of the same data, which they annotated afterwards. As a result, some language pairs have multiple annotations per source segment, was shown in Table 6. This is useful to compute inter-annotator agreement but also provides less noisy annotations.

## 6.2 Quality Control

Each task (100 segments) includes 12 quality control segments to ensure the high quality of the annotations. The tasks are created as follows:

1. The task (a maximum of 100 segments) is filled with machine-translated documents to be evaluated.
2. A random document is selected from the task.
3. Segments within the sampled document are perturbed.
4. The perturbed document is shuffled within the task at the document-level.
5. Steps 2-4 are repeated until 12 quality control segments are included in the task.

The segment perturbation is done by randomly sampling a span from the segment and replacing it with random text sampled from the entire corpus in the correct language. Since segment lengths vary and a single perturbation could be lost in a very large paragraph, we apply as many perturbations as there are sentences in the output. See Figure 2 for an example.

**Source**: *Sie haben gestern das Treffen wieder verschoben.*
**Original**: *He postponed the meeting again yesterday.*
**Perturbed**: *He postponed the meeting squirrels are never.*

**Figure 2:** An example of a perturbed translation based on the original system translation. In addition to the original error (the correct pronoun here is *They* and not *He*), we introduce the perturbed part.

After each task is completed, we check whether the perturbed segments received lower scores. Specifically, we compare the distribution of 12 original and 12 perturbed segments with a one-sided Mann-Whitney U test (Mann and Whitney, 1947). If the task fails to pass quality control ($p>0.05$), it is reset and reassigned to another annotator.[20] In the final data, 96% of perturbed segments have lower scores than their original counterparts.

## 6.3 Human Data Analysis

We briefly analyze the data from a broader perspective. The scores given by the annotators are largely concentrated near 100, with a small peak around 0 (see Figure 3). Most languages consistently had very few errors per segment, resulting in higher overall scores (see Table 7). For instance, for the Czech→Ukrainian, an average of 0.2 minor errors

---

[20]Task generation code: github.com/wmt-conference/ ErrorSpanAnnotations/tree/main/preparation/wmt24

and 0.1 major errors per segment means there is approximately one minor error for every 5 segments and one major error for every 10 segments.

The annotation time, which is the primary focus of the analysis in Kocmi et al. (2024b), is similar across most languages with the exception of English→Icelandic. This could be caused either by more meticulous annotators or lower quality of submitted systems, which would require more annotation. The average time per segment is just 22 seconds (see Figure 4).



**Figure 3:** Distribution of final human segment-level scores. The ratings are dominated by the score close to 100.

| Language pair | Minor | Major | Score | Time |
|---|---|---|---|---|
| Czech→Ukrainian | 0.2 | 0.1 | 87.1 | 15.8s |
| English→Czech | 0.6 | 0.2 | 86.2 | 25.3s |
| English→Spanish | 0.7 | 0.4 | 87.1 | 22.0s |
| English→Hindi | 0.5 | 0.2 | 87.3 | 25.7s |
| English→Icelandic | 1.4 | 0.8 | 72.3 | 37.8s |
| English→Japanese | 0.2 | 0.1 | 89.2 | 18.9s |
| English→Russian | 0.5 | 0.3 | 83.4 | 23.0s |
| English→Ukrainian | 0.4 | 0.3 | 84.4 | 21.8s |
| English→Chinese | 0.2 | 0.1 | 87.6 | 16.8s |

**Table 7:** Average number of minor and major errors per segment, average score and annotation time. Despite different annotation crowds, the statistics are balanced.



**Figure 4:** Distribution of annotation times per segment. The vast majority of segments is annotated under one minute.

## 6.4 Domain Difficulty across Languages

In Table 8 we present the maximal obtained score per domain per language. Although absolute scores are not comparable due to different sets of systems

| | Literary | News | Social | Speech | Average |
|---|---|---|---|---|---|
| En.→Czech | 93.1 | 94.9 | 93.3 | 92.1 | 93.3 |
| En.→Spanish | 96.3 | 96.2 | 95.5 | 94.1 | 95.5 |
| En.→Hindi | 95.4 | 93.6 | 91.3 | 88.3 | 92.2 |
| En.→Icelandic | 92.2 | 92.6 | 95.0 | 92.4 | 93.1 |
| En.→Japanese | 92.4 | 93.7 | 91.3 | 92.4 | 92.5 |
| En.→Russian | 94.1 | 93.1 | 92.1 | 86.6 | 91.5 |
| En.→Ukrainian | 93.2 | 93.9 | 94.3 | 85.9 | 91.8 |
| En.→Chinese | 92.0 | 92.5 | 90.7 | 88.4 | 90.9 |
| Average | 93.6 | 93.8 | 93.0 | 90.0 | 92.6 |

**Table 8:** Maximal obtained score per language and per domain across languages evaluated with the same source data (English).

and different groups of annotators, we observe that, across the table, the speech domain obtains the lowest scores for nearly all language pairs suggesting it is the most difficult domain. This is reflected by the fact that the top-performing systems achieve lower scores in the speech domain compared to other domains. This difficulty likely arises from the reliance on ASR text rather than the original audio. This finding is consistent with MQM results from Freitag et al. (2024).

Secondly, we observe that the English→Spanish language pair receives the highest scores, suggesting that either the pair itself or the specific tested domains are relatively easy for top systems, which provide almost flawless translations. These results are consistent with the MQM results from Freitag et al. (2024) where the best system got only -0.12 MQM score, which is close to perfect, while the best German system got -1.58 MQM and the best Japanese-Chinese system an MQM score of -1.22.

Separate scores for each domain, system and language pair are presented in Appendix D.

## 6.5 Clustering of ESA compared to DA+SQM

This year, we revised the human evaluation protocol, ultimately moving from DA+SQM to ESA. In this section, we briefly compare several aspects of both methods. However, due to the absence of a direct head-to-head comparison on the same data and the many changes introduced since last year, this analysis cannot attribute all the improvements solely to the ESA protocol.

ESA produced 59 clusters across 114 systems. This compares to only 37 clusters produces by last year's DA+SQM approach for the same number of systems. In other words, ESA formed a cluster for every 1.9 systems, while DA+SQM created a cluster for every 3.1 systems. This increased clustering efficiency was achieved despite a decrease in

the number of collected samples. With DA+SQM, we collected an average of 1400 annotations per system, whereas ESA required only an average of 750 annotations per system to achieve greater discriminative power.

# 7 Official Ranking Results

We now describe how we compute the final ranking, then discuss the final results and some potential issues with our ranking method. The results are shown on the following two pages in tabular form.

## 7.1 Human Ranking Computation

We calculate three different scores: the human ESA score, rank, and the cluster.

The **human ESA score** is the macro-average of the segment-level ESA scores grouped over the domains. This represents a change compared to previous years, since we used to calculate a simple average over all data. However, with the change towards paragraph-level test sets, the average number of segments per domain is imbalanced and the social domain represents almost half of all segments (see Table 1). To circumvent this imbalance, we use the macro-average as the main human score.

For the statistical analysis and **clustering**, we use the Wilcoxon signed-rank test, a paired non-parametric test (Wilcoxon, 1945), as suggested by Kocmi et al. (2024b). However, given the domain-level imbalanced distribution, we adapted our approach by combining the results from independent domain-level experiments via Stouffer's Z-score method (Stouffer et al., 1949), which combines p-values from individual domain-level Wilcoxon tests. The method produces almost identical clustering as if we had used Wilcoxon over the whole dataset whilst ignoring the imbalance.

**Rank ranges** indicate the number of systems a particular system underperforms or outperforms: the top end of the rank range is $l + 1$, where $l$ is the number of losses, while the bottom is $n - w$, where $n$ is the total number of systems and $w$ is the number of systems against which the system in question significantly wins.

Systems are grouped into ranks that are separated by thick lines, such that systems within the same group do not strictly outperform other systems within the group. In other words, it is not possible to clearly say which system in the cluster is better than the all others. The ranks and clusters are computed with $p < 0.05$.

We say that a system is winning if it ranks in the first cluster, while ignoring the human reference.

The official rankings shown in Section 7.4 are generated on the basis of the ESA scores. Tables with head-to-head comparisons between all systems are included in Appendix E.

## 7.2 Verbosity of LLMs

As pointed out by Briakou et al. (2024), some LLMs produce verbose outputs, including an attempt to explain the translation or a refusal to translate. This creates an issue for both automatic and human evaluation of how to treat such outputs.

During the collection of LLM outputs, we asked the LLM to wrap the translation in a particular type of quotes (```) and post-edited LLM outputs removing all extra details outside of these quotes (keeping the whole answer if no quotes have been found). Therefore LLMs that did not follow the expected output format and produced additional output were not considered in the evaluation.

For future work, we should instruct humans to penalize verbose outputs and strengthen the prompt used for collecting LLM translations.

## 7.3 Human Ranking Discussion

When investigating the official results in Section 7.4, we can make several observations.

The best performing system in the open and constrained systems category is IOL-Research, winning 10 LPs in this category.

The Unbabel-Tower70B system is the best performing participating system winning in 8 LPs. In contrast, this system was ranked the first in all LPs in the automatic evaluation.This highlights that systems can overfit on automatic scores, especially when using Minimum Bayes Risk (MBR; Freitag et al., 2022) with testing metric.

Over all, the best performing system in general seems to be Claude-3.5-Sonnet (wins in 9 LPs); it even outperforms GPT-4 (wins in 5 LPs), which is much more expensive model. Human references are ranked in the first place for 5 language pairs and in the winning cluster for 8 language pairs, suggesting that the reference quality is high and ESA is robust to our setting.

For English→Icelandic, it was almost the case that each system belonged to its own statistically significant cluster. This could be put down to a greater diversity in the quality of systems (also highlighted by more diverse AutoRank scores).

## 7.4 Official Ranking Results Tables

### Results tables legend

The human score is the macro-average of human judgements, grouped by domain. The rank takes into consideration head-to-head wins and losses. AutoRank is calculated from automatic metrics.

Ranking and clustering on human scores is done using Wilcoxon signed rank test for each domain separately and final p-value is combined via Stouffer's Z-score method to align with macro average for human score.

Systems are either constrained (white), open-track (light gray), or closed-track (dark gray).

LLMs that do not officially claim a support a language pair are marked with §.

### Czech→Ukrainian

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-2 | Claude-3.5 § | 93.0 | 1.7 |
| 2-2 | HUMAN-A | 92.7 | - |
| 3-3 | Gemini-1.5-Pro | 92.6 | 2.0 |
| 3-4 | Unbabel-Tower70B | 92.2 | 1.0 |
| 5-5 | IOL-Research | 90.2 | 1.9 |
| 6-7 | CommandR-plus § | 89.7 | 1.9 |
| 6-8 | ONLINE-W | 88.7 | 2.3 |
| 7-9 | GPT-4 § | 88.6 | 2.0 |
| 8-9 | IKUN | 87.1 | 2.3 |
| 10-10 | Aya23 | 86.6 | 2.5 |
| 11-11 | CUNI-Transformer | 85.3 | 3.0 |
| 12-12 | IKUN-C | 82.6 | 3.0 |

### English→Czech

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-2 | HUMAN-A | 92.9 | - |
| 2-2 | Unbabel-Tower70B | 91.6 | 1.0 |
| 2-3 | Claude-3.5 § | 91.2 | 2.1 |
| 4-5 | ONLINE-W | 89.0 | 2.8 |
| 4-6 | CUNI-MH | 88.4 | 2.1 |
| 6-6 | Gemini-1.5-Pro | 88.2 | 2.6 |
| 6-8 | GPT-4 § | 87.7 | 2.6 |
| 8-8 | CommandR-plus § | 86.9 | 2.9 |
| 8-9 | IOL-Research | 86.5 | 2.8 |
| 10-11 | SCIR-MT | 85.4 | 3.2 |
| 10-11 | CUNI-DocTransformer | 84.3 | 4.4 |
| 12-12 | Aya23 | 84.2 | 4.3 |
| 13-13 | CUNI-GA | 82.1 | 2.3 |
| 14-14 | IKUN | 81.7 | 3.9 |
| 15-15 | Llama3-70B § | 77.4 | 4.1 |
| 16-16 | IKUN-C | 75.4 | 4.7 |

### English→German

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-11 | GPT-4 | -1.6 | 1.8 |
| 1-7 | Dubformer | -1.8 | 1.8 |
| 2-10 | ONLINE-B | -1.9 | 1.8 |
| 2-10 | TranssionMT | -1.9 | 1.8 |
| 2-9 | Unbabel-Tower70B | -1.9 | 1.0 |
| 1-9 | HUMAN-B | -2.0 | - |
| 2-12 | Mistral-Large | -2.1 | 2.0 |
| 4-11 | CommandR-plus | -2.3 | 2.0 |
| 8-10 | ONLINE-W | -2.3 | 2.2 |
| 2-12 | Claude-3.5 | -2.4 | 1.9 |
| 3-13 | HUMAN-A | -2.5 | - |
| 10-12 | IOL-Research | -2.5 | 2.3 |
| 5-13 | Gemini-1.5-Pro | -2.8 | 2.2 |
| 14-15 | Aya23 | -3.2 | 2.7 |
| 14-17 | ONLINE-A | -3.5 | 3.0 |
| 15-17 | Llama3-70B § | -4.3 | 2.5 |
| 15-17 | IKUN | -4.3 | 3.0 |
| 18-18 | IKUN-C | -6.1 | 3.8 |
| 19-19 | MSLC | -15.5 | 11.9 |

### English→Spanish

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-1 | HUMAN-A | 95.3 | - |
| 2-2 | Dubformer | 93.4 | 2.0 |
| 3-4 | GPT-4 | 91.9 | 1.9 |
| 4-7 | IOL-Research | 91.4 | 2.3 |
| 5-8 | Mistral-Large | 89.3 | 2.2 |
| 5-9 | Unbabel-Tower70B | 88.9 | 1.0 |
| 3-8 | Claude-3.5 | 88.8 | 2.1 |
| 5-8 | Gemini-1.5-Pro | 88.8 | 2.4 |
| 7-9 | CommandR-plus | 88.3 | 2.1 |
| 9-10 | Llama3-70B § | 87.2 | 2.6 |
| 11-11 | ONLINE-B | 85.6 | 2.7 |
| 12-13 | IKUN | 84.7 | 2.8 |
| 12-13 | IKUN-C | 80.4 | 3.4 |
| 14-14 | MSLC | 63.9 | 7.4 |

### English→Hindi

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-3 | TranssionMT | 91.3 | 1.3 |
| 1-4 | Unbabel-Tower70B | 90.5 | 1.0 |
| 3-3 | Claude-3.5 § | 90.2 | 1.2 |
| 3-4 | ONLINE-B | 90.1 | 1.4 |
| 3-5 | Gemini-1.5-Pro § | 90.0 | 1.6 |
| 6-6 | GPT-4 § | 88.5 | 2.1 |
| 7-8 | HUMAN-A | 88.5 | - |
| 8-8 | IOL-Research | 87.2 | 2.1 |
| 8-9 | Llama3-70B § | 86.7 | 2.1 |
| 10-10 | Aya23 | 84.7 | 3.2 |
| 11-11 | IKUN-C | 70.7 | 5.5 |

### English→Icelandic

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-1 | HUMAN-A | 93.1 | - |
| 2-3 | Dubformer | 84.3 | 2.5 |
| 2-3 | Claude-3.5 § | 81.9 | 2.3 |
| 4-4 | Unbabel-Tower70B | 80.2 | 1.0 |
| 5-5 | AMI | 73.3 | 3.7 |
| 6-6 | IKUN | 71.0 | 3.2 |
| 7-7 | ONLINE-B | 68.0 | 4.2 |
| 8-9 | GPT-4 | 66.3 | 3.4 |
| 8-9 | IKUN-C | 65.2 | 3.7 |
| 10-10 | IOL-Research | 58.0 | 4.3 |
| 11-11 | Llama3-70B § | 41.0 | 6.7 |

### English→Ukrainian

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-2 | Claude-3.5 | 90.5 | 2.0 |
| 1-2 | Unbabel-Tower70B | 89.8 | 1.0 |
| 3-3 | Dubformer | 89.0 | 1.8 |
| 4-6 | HUMAN-A | 87.3 | - |
| 4-6 | Gemini-1.5-Pro | 87.1 | 2.2 |
| 5-8 | ONLINE-W | 86.0 | 2.1 |
| 5-9 | GPT-4 | 84.6 | 2.3 |
| 6-9 | CommandR-plus § | 83.2 | 2.3 |
| 7-9 | IOL-Research | 83.2 | 2.4 |
| 10-10 | IKUN | 78.4 | 2.8 |
| 11-11 | IKUN-C | 67.9 | 3.9 |

### English→Japanese

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-1 | HUMAN-A | 91.8 | - |
| 2-4 | ONLINE-B | 91.1 | 1.4 |
| 3-4 | CommandR-plus | 91.0 | 1.9 |
| 4-4 | GPT-4 | 90.8 | 1.7 |
| 4-5 | Claude-3.5 | 90.8 | 1.5 |
| 4-7 | Gemini-1.5-Pro | 90.0 | 1.7 |
| 7-7 | Unbabel-Tower70B | 89.7 | 1.0 |
| 8-8 | IOL-Research | 89.7 | 2.3 |
| 8-9 | Aya23 | 89.7 | 2.3 |
| 10-10 | NTTSU | 89.4 | 1.9 |
| 11-11 | Team-J | 88.5 | 1.9 |
| 12-12 | Llama3-70B § | 86.8 | 2.6 |
| 13-13 | IKUN-C | 81.7 | 3.9 |

### English→Chinese

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-1 | GPT-4 | 89.6 | 2.0 |
| 2-4 | Unbabel-Tower70B | 89.6 | 1.0 |
| 2-4 | HUMAN-A | 89.4 | - |
| 4-4 | Gemini-1.5-Pro | 89.3 | 1.8 |
| 5-6 | ONLINE-B | 89.3 | 1.7 |
| 6-6 | IOL-Research | 89.0 | 1.8 |
| 6-7 | Claude-3.5 | 88.9 | 1.7 |
| 6-8 | CommandR-plus | 88.3 | 2.2 |
| 9-9 | Llama3-70B § | 86.5 | 2.8 |
| 10-10 | HW-TSC | 86.2 | 2.3 |
| 11-11 | IKUN | 85.3 | 3.1 |
| 12-12 | Aya23 | 85.2 | 3.0 |
| 13-13 | IKUN-C | 82.1 | 3.5 |

### English→Russian

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-1 | HUMAN-A | 89.2 | - |
| 2-3 | Dubformer | 89.1 | 1.9 |
| 3-4 | Claude-3.5 | 88.2 | 2.0 |
| 3-5 | Unbabel-Tower70B | 88.1 | 1.0 |
| 3-7 | Yandex | 87.0 | 1.9 |
| 6-8 | Gemini-1.5-Pro | 85.5 | 2.3 |
| 6-9 | GPT-4 | 85.0 | 2.3 |
| 6-9 | ONLINE-G | 84.6 | 2.2 |
| 5-9 | CommandR-plus § | 84.3 | 2.4 |
| 10-10 | IOL-Research | 82.1 | 2.6 |
| 11-11 | IKUN | 79.2 | 3.2 |
| 12-12 | Aya23 | 78.6 | 3.3 |
| 13-13 | Llama3-70B § | 75.7 | 3.1 |
| 14-14 | IKUN-C | 69.8 | 3.9 |

### Japanese→Chinese

| Rank | System | Human | AutoRank |
|------|--------|-------|----------|
| 1-3 | Claude-3.5 | -1.4 | 1.7 |
| 1-3 | HUMAN-A | -1.5 | - |
| 3-5 | GPT-4 | -1.7 | 2.1 |
| 2-5 | DLUT-GTCOM | -1.7 | 2.0 |
| 4-8 | Unbabel-Tower70B | -1.9 | 1.0 |
| 3-6 | Gemini-1.5-Pro | -2.1 | 1.9 |
| 6-8 | CommandR-plus | -2.2 | 2.8 |
| 6-8 | IOL-Research | -2.4 | 2.2 |
| 9-10 | Llama3-70B § | -3.4 | 3.1 |
| 9-10 | Aya23 | -3.5 | 3.7 |
| 11-12 | Team-J | -4.5 | 2.8 |
| 11-12 | NTTSU | -5.1 | 3.7 |
| 13-13 | ONLINE-B | -5.8 | 5.2 |
| 14-14 | IKUN-C | -7.7 | 5.5 |
| 15-15 | MSLC | -10.7 | 8.9 |

## 8 Test Suites Sub-task: "Help us break LLMs"

The results in the previous tables indicate that the current evaluation methods, despite being more detailed and sophisticated, have difficulties in distinguishing MT output from human translations, or distinguishing the performance among different systems. Additionally, the appearance of LLMs has made it even more clear that generated translations, even those which seem to be fluent and surrounded by seemingly perfect content, can contain serious flaws. The increased interest in this new technology and the use of LLMs for translation, prompted us to set the theme of this year's test suite sub-task as "Help us break LLMs". This was intended as a broader invitation to the NLP community to expose the weaknesses of LLM translations that are hidden behind the apparent overall high quality generation, but also to propose new innovative evaluation methods that may be of high interest for specific use cases. We are thrilled that this year's participation exceeded every precedent, with 11 test suites providing their valuable conclusions, which are presented below.

### 8.1 Setup of the sub-task

Each test suite is a customised extension of the standard test sets, focusing on specific aspects of the MT output. The evaluation of the MT output takes place in a decentralized manner, where test suite providers were invited to submit their customized test sets, following the setup introduced at the Third Conference on Machine Translation (Bojar et al., 2018). Each test suite provider submitted a source-side test set, which was appended by the organisers of the General MT Shared Task to its standard test sets. The corresponding outputs from the systems of the General MT Shared Task were returned to the test suite providers, who were responsible for carrying out the evaluation based on their own individual evaluation concept. The results of each test suite evaluation, together with the relevant analysis, appear in separate description papers, while a summary is given below.

This year's timeline gave the test suite contributors more time. We offered a pre-run in April, when test suite providers were given the opportunity to submit the current version of their corpus in order to receive translation output from online systems, which could help them to carry out the individual (often manual) evaluation in a more timely manner.

### 8.2 Submissions

The test suite sub-task received 11 submissions, out of which 9 completed the entire evaluation cycle. An overview of the test suites can be seen in Table 9. The descriptions of each submission and their main findings are given below.

**Árni Magnússon Institute for Icelandic Studies** (AMI; Ármannsson et al., 2024) The submission of the Árni Magnússon Institute's team to the WMT24 test suite subtask focuses on idiomatic expressions and proper names for the English→Icelandic translation direction. Intuitively and empirically, idioms and proper names are known to be a significant challenge for neural translation models. They create two different test suites. The first evaluates the competency of MT systems in translating common English idiomatic expressions, as well as testing whether systems can distinguish between those expressions and the same phrases when used in a literal context. The second test suite consists of place names that should be translated into their Icelandic exonyms (and correctly inflected) and pairs of Icelandic names that share a surface form between the male and female variants, so that incorrect translations impact meaning as well as readibility. The scores reported are relatively low, especially for idiomatic expressions and place names, and indicate considerable room for improvement.

**Complex Sentence Structure Testset** (CoST; vIIT_HYD; Mukherjee et al., 2024) This test suite presents an evaluation of 16 machine translation systems submitted to the Shared Task for the English-Hindi using our Complex Structures Test suite. Aligning with this year's test suite sub-task theme, "Help us break LLMs", the authors curated a comprehensive test suite encompassing diverse datasets across various categories, including autobiography, poetry, legal, conversation, play, narration, technical, and mixed genres. The evaluation reveals that all the systems struggle significantly with the archaic style of text like legal and technical writings or text with creative twist like conversation and poetry datasets, highlighting their weaknesses in handling complex linguistic structures and stylistic nuances inherent in these text types. This evaluation identifies the strengths and limitations of the models, pointing to specific areas where further research is needed to enhance their performance.[21]

---

[21] github.com/AnanyaCoder/CoST-WMT-24-Test-Suite-Task

| Test suite | Institution | Focus | Language pair | Segments |
|---|---|---|---|---|
| AMI (Ármannsson et al., 2024) | AMI | idiomatic expressions, proper names | en→is | 3,082 |
| COST (Mukherjee et al., 2024) | IIIT_HYD | complex sentence structure | en→hi | 3,908 |
| DFKI (Manakhimova et al., 2024) | DFKI | 110 linguistic phenomena | en→de, en→ru | 54,736 |
| GenderQueer (Friðriksdóttir, 2024) | UI | gender-diverse, queer-inclusive content | en→is | 672 |
| IITP (Bhattacharjee et al., 2024) | IITP | multi-domain dynamics | en→hi | 4,198 |
| Isochrony (Rozanov et al., 2024) | RaskAI, IC | isochrony of translations | en→de, en→es, en→ja, en→ru, en→zh | 10,730 |
| NRCC (Dawkins et al., 2024) | NRCC | speaker-listener gender resolution | en→cs, en→de, en→es, en→is | 53,560 |
| PIA_TQA (Miceli Barone and Sun, 2024) | UEDIN | prompt injection attacks | cs→uk, en→cs, en→de, en→es, en→hi, en→is, en→ja, en→ru, en→uk, en→zh, ja→zh | 250,744 |
| RoCS-MT (Bawden and Sagot, 2023) | Inria | robustness to non-standard user-generated texts | en→cs, en→de, en→es, en→hi, en→is, en→ja, en→ru, en→uk, en→zh | 7883 |

**Table 9:** Overview of the participating test suites.

**DFKI (Manakhimova et al., 2023b)** This test suite offers a fine-grained linguistically motivated analysis of the shared task MT outputs for English–German and English–Russian, based on more than 11,500 manually devised test items, which cover up to 110 phenomena in 14 categories per language direction. Extending their previous test suite submissions (e.g. Avramidis et al., 2020; Macketanz et al., 2021, 2022; Manakhimova et al., 2023a), the submission of this year includes a considerable effort of manual linguistic annotation for the evaluation on 39 MT systems submitted at the Shared Task. Based on the results, LLMs are inferior to NMT in English–German when translating a few linguistic phenomena, though they show quite a competitive performance in English-Russian. Additionally, some LLMs generate very verbose or empty outputs, posing challenges to the evaluation process. Looking more closely at specific phenomena of English-German, LLMs seem to perform worse than the two best performing NMT systems in terms of punctuation, future verb tenses and stripping. For English-Russian, Yandex is weaker in named entities and terminology, Claude in function words, while Unbabel is weaker in verb valency. GPT-4 into Russian performs even worse than several commercial NMT-based systems.

**Indian Institute of Technology Patna** (IITP; domain dynamics; Bhattacharjee et al., 2024) LLMs have demonstrated impressive capabilities in machine translation, leveraging extensive pretraining on vast amounts of data. However, this generalist training often overlooks domain-specific nuances, leading to potential difficulties when translating

specialized texts. This study presents a multi-domain dataset designed to challenge and evaluate the translation abilities of LLMs. The dataset encompasses diverse domains such as judicial, education, literature (specifically religious texts), and noisy user-generated content from online product reviews and forums like Reddit. Each domain consists of approximately 250–300 sentences, carefully curated and randomized in the final compilation. This English-to-Hindi dataset aims to evaluate and expose the limitations of LLM-based translation systems, offering valuable insights into areas requiring further research and development.

**Inria** (RoCS-MT; Bawden and Sagot, 2023), Robust Challenge Set for Machine Translation, is designed to test MT systems' ability to translate user-generated content with non-standard characteristics, such as spelling errors, devowelling, acronymisation, etc. The original English Reddit texts are associated with manual normalisations and translations in five languages (French, German, Czech, Ukrainian and Russian). RoCS-MT was first submitted to the 2023 task, showing that many non-standard phenomena still pose problems for most systems, although more common phenomena are better handled by the larger, closed-source models, presumably due to the large quantity of web-based seen during training. This year's version is largely the same as last year but with some improvements, including modifications to normalisations and to the annotation typology used (all modifications are documented in the GitHub repository).[22] Systems varied greatly in terms of their handling of

---

[22] github.com/rbawden/RoCS-MT

non-standard sentences, with marked differences depending on the type of system. Constrained systems inevitably struggling most, particularly with phenomena affecting the spelling of words (resulting in frequent copying of non-standard source words), a problem also affecting online systems. LLMs exhibited some of the best quality translations, although behaviour varied between translating standard and non-standard input, and additional issues such as refusal to translate and usage notes pose new technical challenges.

**Isochrony Translation**  (Rask AI, Imperial College; Rozanov et al., 2024) MT has come a long way and is readily employed in production systems to serve millions of users daily. With the recent advances in generative AI, a new form of translation is becoming possible – video dubbing. This work motivates the importance of isochronic translation, especially in the context of automatic dubbing, and introduces 'IsoChronoMeter' (ICM). ICM is a simple yet effective metric to measure isochrony of translations in a scalable and resource efficient way without the need for gold data, based on state-of-the-art text-to-speech (TTS) duration predictors. The authors motivate IsoChronoMeter and demonstrate its effectiveness. Using ICM, they demonstrate the short-comings of state-of-the-art translation systems and show the need for new methods. The code has been released.

**National Research Council Canada**  (Speaker-Listener Gender Resolution; gender-res; Dawkins et al., 2024) This test suite assesses the gender resolution tendencies of MT systems in literary dialogue settings. That is, each instance contains dialogue interleaved with additional meta-context. The spoken dialogue refers to either the speaker or listener such that the gender of the referent, if known, must be inferred from the meta-context and informs the correct translation. They find that stereotype factors within the meta-context, such as character descriptions and manner of speaking, affect the gender agreement choices of words within the dialogue. Regression analysis is performed to evaluate the relative influence of these contextual factors compared to structural factors and known stereotype influences (e.g., the internal gender stereotype of an adjective).

**University of Edinburgh Prompt Injection, TruthfulQA**  (PIA; Miceli Barone and Sun, 2024) LLM-based systems typically work by embedding

their input data into prompt templates which contain instructions and/or in-context examples, creating queries which are submitted to a LLM, then parse the LLM response in order to generate the system outputs. Prompt Injection Attacks (PIAs) are a type of subversion of these systems where a malicious user crafts special inputs which interfere with the prompt templates, causing the LLM to respond in ways unintended by the system designer. Recently, Sun and Miceli Barone (2024) proposed a class of PIAs against LLM-based machine translation. Specifically, the task is to translate questions from the TruthfulQA test suite, where an adversarial prompt is prepended to the questions, instructing the system to ignore the translation instruction and answer the questions instead. In this test suite, the authors extend this approach to all the language pairs of the WMT 2024 General Machine Translation task. Moreover, they include additional attack formats in addition to the one originally studied.

**University of Iceland**  (GenderQueer; Friðriksdóttir, 2024) This paper introduces the GenderQueer Test Suite, a novel evaluation set for assessing MT systems' capabilities in handling gender-diverse and queer-inclusive content, focusing on English to Icelandic translation. As MT quality improves, there is an increasing need for specialized evaluation methods that address nuanced aspects of language and identity. The suite evaluates MT systems on various aspects of gender-inclusive translation, including pronoun and adjective agreement, LGBTQIA+ terminology accuracy, and the impact of explicit gender specifications. Its authors evaluated 18 MT systems submitted to the WMT24 English-Icelandic track. Key findings reveal significant performance differences between large language model-based systems and smaller models in handling context for gender agreement. Challenges in translating singular "they" were widespread, while most systems performed well in translating LGBTQIA+ terminology. Accuracy in adjective gender agreement varies, with some models struggling particularly with feminine forms. This evaluation set contributes to the ongoing discussion about inclusive language in MT and natural language processing. By providing a tool for assessing MT systems' handling of gender-diverse content, it aims to enhance the inclusivity of language technology. The methodology and evaluation scripts are made available for adaptation to other languages, promoting further research in this critical area.

## 9 Conclusions

The WMT 2024 General Machine Translation Task covered 11 translation pairs, two of which are non-English: Czech→Ukrainian and Japanese→Chinese. We introduced ESA (Error Span Annotations) as the main human protocol for assessing the translation quality, which enabled more efficient collection of human judgements than MQM while keeping high quality of annotations. In total, 108 human (semi-)professional annotators contributed more than 57,000 judgments.

We received 105 primary submissions from 28 participants, 4 online systems and 8 production large language models, which is a large increase from last year's task. The majority of participants already use LLMs in their systems.

The best performing open system is IOL-Research (wins 10 LPs in it's category), the best performing participating system is Unbabel-Tower70B (wins 8 LPs), and the best performing system in general is Claude-3.5-Sonnet (wins 9 LPs).

While the best performing system based on automatic metrics is Unbabel-Tower70B, it was not the winner across the board in the human evaluation, with the mismatch between the results likely due to metric bias (Kovacs et al., 2024) in MBR. This shows that human evaluation should be used as the final judge of translation quality.

Lastly, we showed promising results in the multimodal evaluation of the speech domain, proving to be a challenging domain for MT systems. On the opposite side, systems were able to produce near-perfect translations in English→Spanish, for the domains that we tested.

## 10 Limitations

We tested the general capabilities of MT systems. However, we have simplified this approach and only used three to five domains. Out of various modalities, we used audio and text.

Although we use human judgements as the gold standard, giving us more reliable signal than automatic metrics, we should mention that human annotations are noisy (Wei and Jia, 2021) and their performance is affected by the quality of other evaluated systems (Mathur et al., 2020). Lastly, different annotators use different ranking strategies, which may have an effect on the system ranking.

Some models may have used Comet or MetricX during their training, for example, using Minimum Bayes Risk. Our automatic evaluation of such models will be biased, giving them artificially higher scores.

Automatic metrics are limited and biased (Karpinska et al., 2022; Moghe et al., 2024), especially in novel domains (Zouhar et al., 2024a), which motivates them being superseded by human evaluation. Another potential problem may have been that test sets we use are paragraph-level; automatic metrics have usually been tested in a sentence-level scenario.

The ESA annotation interface implemented in Appraise is in English only with a tutorial in German→English. This caused difficulties to some of the Czech→Ukrainian annotators we hired, who could not understand English. One such annotator did not pass the initial tutorial and therefore did not participate in the annotation campaign. Next year, we plan to translate the annotation interface to either the source or target language for each translation direction.

## 11 Ethical Considerations

Inappropriate, controversial, and explicit content was filtered out prior to translation, keeping in mind the translators and not exposing them to such content or obliging them to translate it.

Human evaluation using Appraise for the collection of human judgements was fully anonymous. Automatically generated accounts associated with annotation tasks with single-sign-on URLs were distributed randomly among pools of annotators and we do not store any personal information. We do store the mapping between which annotator (with pseudonym) annotated which account. Annotators received standard professional translator's or evaluator's wage with respect to their countries.

Sentences in the Czech→Ukrainian dataset (in Personal, Official and Voice domains) were collected with users' opt-in consent, and any personal data related to people other than well-known people was pseudonymized (using random first names and surnames). Sentences where such pseudonymization would not be enough to preserve reasonable anonymity of the users (e.g., describing events uniquely identifying the persons involved) were not included in the test set.

# Acknowledgments

# References

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina Espa textasciitilde na-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88. Association for Computational Linguistics.

Bjarki Ármannsson, Hinrik Hafsteinsson, Atli Jasonarson, and Steinthor Steingrimsson. 2024. Killing two flies with one stone: An attempt to break llms using english→icelandic idioms and proper names. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, Kelly Marchisio, Max Bartolo, Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. 2024. Aya 23: Open Weight Releases to Further Multilingual Progress.

Eleftherios Avramidis, Annika Grützner-Zahn, Manuel Brack, Patrick Schramowski, Pedro Ortiz Suarez, Malte Ostendorff, Fabio Barth, Shushen Manakhimova, Vivien Macketanz, Georg Rehm, and Kristian Kersting. 2024. Occiglot at WMT24: European open-source large language models evaluated on translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, Aljoscha Burchardt, and Sebastian Möller. 2020. Fine-grained linguistic evaluation for state-of-the-art machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 346–356. Association for Computational Linguistics.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567. Association for Computational Linguistics.

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 conference on machine translation (WMT20). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55. Association for Computational Linguistics.

Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. Findings of the 2019 conference on machine translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61. Association for Computational Linguistics.

Rachel Bawden and Benoît Sagot. 2023. RoCS-MT: Robustness challenge set for machine translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 198–216. Association for Computational Linguistics.

Soham Bhattacharjee, Baban Gain, and Asif Ekbal. 2024. Domain dynamics: Evaluating large language models in english-hindi translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Ondřej Bojar, Christian Buck, Chris Callison-Burch, Christian Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, and Lucia Specia. 2013. Findings of the 2013 Workshop on Statistical Machine Translation. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 1–44. Association for Computational Linguistics.

Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleš Tamchyna. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58. Association for Computational Linguistics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. Findings of the 2017 conference on machine translation (WMT17). In *Proceedings of the Second Conference on Machine Translation*, pages 169–214. Association for Computational Linguistics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198. Association for Computational Linguistics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Barry Haddow, Matthias Huck, Chris Hokamp, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Carolina Scarton, Lucia Specia, and Marco Turchi. 2015. Findings of the 2015 workshop on statistical machine translation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 1–46. Association for Computational Linguistics.

Ondřej Bojar, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, and Christof Monz. 2018. Findings of the 2018 conference on machine translation (WMT18). In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 272–303. Association for Computational Linguistics.

Eleftheria Briakou, Zhongtao Liu, Colin Cherry, and Markus Freitag. 2024. On the implications of verbose llm outputs: A case study in translation evaluation.

Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2007. (meta-) evaluation of machine translation. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 136–158. Association for Computational Linguistics.

Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2008. Further meta-evaluation of machine translation. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 70–106. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, Kay Peterson, Mark Przybocki, and Omar Zaidan. 2010. Findings of the 2010 joint workshop on statistical machine translation and metrics for machine translation. In *Proceedings of the Joint Fifth Workshop on Statistical Machine Translation and MetricsMATR*, pages 17–53. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, and Lucia Specia. 2012. Findings of the 2012 workshop on statistical machine translation. In *Proceedings of the Seventh Workshop on Statistical Machine Translation*, pages 10–51. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, and Josh Schroeder. 2009. Findings of the 2009 Workshop on Statistical Machine Translation. In *Proceedings of the Fourth Workshop on Statistical Machine Translation*, pages 1–28. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, and Omar Zaidan. 2011. Findings of the 2011 workshop on statistical machine translation. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 22–64. Association for Computational Linguistics.

Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. WIT3: Web inventory of transcribed and translated talks. In *Proceedings of the 16th Annual Conference of the European Association for Machine Translation*, pages 261–268. European Association for Machine Translation.

A Conneau. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451. Association for Computational Linguistics.

Hillary Dawkins, Isar Nejadgholi, and Chi-kiu Lo. 2024. WMT24 test suite: Gender resolution in speaker-listener dialogue roles. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Sören Dreano, Derek Molloy, and Noel Murphy. 2024. Cyclegn: a cycle consistent approach for neural machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2020. CCAligned: A massive collection of cross-lingual web-document pairs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5960–5969. Association for Computational Linguistics.

Denis Elshin, Nikolay Karpachev, Boris Gruzdev, Ilya Golovanov, Georgy Ivanov, Alexander Antonov, Nickolay Skachkov, Ekaterina Latypova, Vladimir Layner, Ekaterina Enikeeva, Dmitry Popov, Anton Chekashev, Vladislav Negodin, Vera Frantsuzova, Alexander Chernyshev, and Kirill Denisov. 2024. From general LLM to translation: How we dramatically improve translation quality using human evaluation data for LLM finetuning. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Christian Federmann. 2018. Appraise evaluation framework for machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 86–88. Association for Computational Linguistics.

Markus Freitag, Isaac Caswell, and Scott Roy. 2019. APE at scale and its implications on MT evaluation biases. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 34–44. Association for Computational Linguistics.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are llms breaking mt metrics? results of the wmt24 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628. Association for Computational Linguistics.

Steinunn Rut Friðriksdóttir. 2024. The genderqueer test suite. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. 2012. Building large monolingual dictionaries at the Leipzig corpora collection: From 100 to 200 languages. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 759–765. European Language Resources Association (ELRA).

Thamme Gowda, Zhao Zhang, Chris Mattmann, and Jonathan May. 2021. Many-to-English machine translation tools, data, and pretrained models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 306–316. Association for Computational Linguistics.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In *Proceedings of the 7th Linguistic Annotation Workshop*

*and Interoperability with Discourse*, pages 33–41. Association for Computational Linguistics.

Yvette Graham, Barry Haddow, and Philipp Koehn. 2020. Statistical power and translationese in machine translation evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 72–81. Association for Computational Linguistics.

Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.

Miroslav Hrabal, Josef Jon, Martin Popel, Nam Luu, Danil Semin, and Ondřej Bojar. 2024. CUNI at WMT24 general translation task: Llms, (q)lora, CPO and model merging. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Atli Jasonarson, Hinrik Hafsteinsson, Bjarki Ármannsson, and Steinþór Steingrímsson. 2024. Cogs in a machine, doing what they're meant to do – the AMI submission to the WMT24 general translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767. Association for Computational Linguistics.

Marzena Karpinska, Nishant Raj, Katherine Thai, Yixiao Song, Ankita Gupta, and Mohit Iyyer. 2022. DEMETR: Diagnosing evaluation metrics for translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9540–9561. Association for Computational Linguistics.

Ondrej Klejch, Eleftherios Avramidis, Aljoscha Burchardt, and Martin Popel. 2015. MT-ComparEval: Graphical evaluation interface for machine translation development. *Prague Bull. Math. Linguistics*, 104:63–74.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz,

Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2024a. Preliminary WMT24 Ranking of General MT Systems and LLMs. *arXiv preprint arXiv:2407.19884*.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. GEMBA-MQM: Detecting translation quality error spans with GPT-4. In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775. Association for Computational Linguistics.

Tom Kocmi, Martin Popel, and Ondrej Bojar. 2020. Announcing CzEng 2.0 Parallel Corpus with over 2 Gigawords. *CoRR*, abs/2007.03006.

Tom Kocmi, Vilém Zouhar, Eleftherios Avramidis, Roman Grundkiewicz, Marzena Karpinska, Maja Popović, Mrinmaya Sachan, and Mariya Shmatova. 2024b. Error span annotation: A balanced approach for human evaluation of machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86.

Philipp Koehn and Christof Monz. 2006. Proceedings on the workshop on statistical machine translation. New York, USA. Association for Computational Linguistics.

Minato Kondo, Ryo Fukuda, Xiaotian Wang, Katsuki Chousa, Masato Nishimura, Kosei Buma, Takatomo Kano, and Takehito Utsuro. 2024. NTTSU at WMT2024 general translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Geza Kovacs, Daniel Deutsch, and Markus Freitag. 2024. Mitigating metric bias in minimum bayes risk

decoding. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Julia Kreutzer, Nathaniel Berger, and Stefan Riezler. 2020. Correct me if you can: Learning from error corrections and markings. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 135–144. European Association for Machine Translation.

Keito Kudo, Hiroyuki Deguchi, Makoto Morishita, Ryo Fujii, Takumi Ito, Shintaro Ozaki, Koki Natsumi, Kai Sato, Kazuki Yano, Ryosuke Takahashi, Subaru Kimura, Tomomasa Hara, Yusuke Sakai, and Jun Suzuki. 2024. Document-level translation with LLM reranking: Team-j at WMT 2024 general translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Taku Kudo. 2005. Mecab: Yet another part-of-speech and morphological analyzer. https://taku910.github.io/mecab/. Accessed: 2023-10-02.

Samuel Larkin, Chi-kiu Lo, and Rebecca Knowles. 2024. MSLC24 submissions to the general machine translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Baohang Li, Zekai Ye, yichong huang, Xiaocheng Feng, and Bing Qin. 2024. SCIR-MT's submission for WMT24 general machine translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Baohao Liao, Christian Herold, Shahram Khadivi, and Christof Monz. 2024. IKUN for WMT24 general MT task: Llms are here for multilingual machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Arle Lommel, Aljoscha Burchardt, Maja Popović, Kim Harris, Eleftherios Avramidis, and Hans Uszkoreit. 2014. Using a new analytic measure for the annotation and analysis of MT errors on real data. In *Proceedings of the 17th Annual Conference of the European Association for Machine Translation*, pages 165–172. European Association for Machine Translation.

Samuel Läubli, Sheila Castilho, Graham Neubig, Rico Sennrich, Qinlan Shen, and Antonio Toral. 2020. A Set of Recommendations for Assessing Human–Machine Parity in Language Translation. *Journal of Artificial Intelligence Research (JAIR)*, 67.

Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. 2023. Mega: Moving average equipped gated attention. In *The Eleventh International Conference on Learning Representations*.

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, and Hans Uszkoreit. 2022. A linguistically motivated test suite to semi-automatically evaluate German–English machine translation output. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 936–947. European Language Resources Association.

Vivien Macketanz, Eleftherios Avramidis, Shushen Manakhimova, and Sebastian Möller. 2021. Linguistic evaluation for the 2021 state-of-the-art machine translation systems for German to English and English to German. In *Proceedings of the Sixth Conference on Machine Translation*, pages 1059–1073. Association for Computational Linguistics.

Shushen Manakhimova, Eleftherios Avramidis, Vivien Macketanz, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2023a. Linguistically motivated evaluation of the 2023 state-of-the-art machine translation: Can ChatGPT outperform NMT? In *Proceedings of the Eighth Conference on Machine Translation*, pages 224–245. Association for Computational Linguistics.

Shushen Manakhimova, Eleftherios Avramidis, Vivien Macketanz, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2023b. Linguistically motivated Evaluation of the 2023 State-of-the-art Machine Translation: Can ChatGPT Outperform NMT? In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.

Shushen Manakhimova, Vivien Macketanz, Eleftherios Avramidis, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2024. Investigating the linguistic performance of large language models in machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

H. B. Mann and D. R. Whitney. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The Annals of Mathematical Statistics*, 18(1):50–60.

Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4984–4997. Association for Computational Linguistics.

Antonio Valerio Miceli Barone and Zhifan Sun. 2024. A test suite of prompt injection attacks for LLM-based machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Nikita Moghe, Arnisa Fazla, Chantal Amrhein, Tom Kocmi, Mark Steedman, Alexandra Birch, Rico Sennrich, and Liane Guillou. 2024. Machine Translation

Meta Evaluation through Translation Accuracy Challenge Sets. *Computational Linguistics*, pages 1–60.

Makoto Morishita, Jun Suzuki, and Masaaki Nagata. 2020. JParaCrawl: A large scale web-based English-Japanese parallel corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3603–3609. European Language Resources Association.

Ananya Mukherjee, Saumitra Yadav, and Manish Shrivastava. 2024. Cost of breaking the llms. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Vladimir Aleksandrovich Mynka and Nikolay Mikhaylovskiy. 2024. TSU HITS's submissions to the WMT 2024 general machine translation shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Graham Neubig. 2011. The Kyoto free translation task. http://www.phontron.com/kftt.

Mariana Neves, Cristian Grozea, Philippe Thomas, Roland Roller, Rachel Bawden, Aurélie Névéol, Steffen Castle, Vanessa Bonato, Giorgio Maria Di Nunzio, Federica Vezzani, Maika Vicente Navarro, Lana Yeganova, and Antonio Jimeno Yepes. 2024. Findings of the WMT 2024 Biomedical Translation Shared Task: Test Sets on Abstract Level. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

OpenAI. 2024. GPT-4 Technical Report.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395. Association for Computational Linguistics.

Maja Popović. 2020. Informative manual evaluation of machine translation output. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5059–5069. International Committee on Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191. Association for Computational Linguistics.

Reid Pryzant, Youngjoo Chung, Dan Jurafsky, and Denny Britz. 2018. JESC: Japanese-English subtitle corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. European Language Resources Association (ELRA).

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision.

Ricardo Rei, Nuno M. Guerreiro, Jos textasciitilde A© Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848. Association for Computational Linguistics.

Ricardo Rei, Jose Maria Pombal, Nuno M. Guerreiro, João Alves, Pedro Henrique Martins, Patrick Fernandes, Helena Wu, Tania Vaz, Duarte Alves, Amin Farajian, Sweta Agrawal, Antonio Farinhas, José G. C. de Souza, and André Martins. 2024. Tower v2: Unbabel-IST 2023 submission for the general MT shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Matiss Rikters and Makoto Miwa. 2024. AIST AIRC systems for the WMT 2024 shared tasks. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Nikolai Rozanov, Vikentiy Pankov, Dmitrii Mukhutdinov, and Dima Vypirailenko. 2024. Isochronometer: A simple and effective isochronic translation evaluation metric. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Roberts Rozis and Raivis Skadiņš. 2017. Tilde MODEL - multilingual open data for EU languages. In *Proceedings of the 21st Nordic Conference on Computational Linguistics*, pages 263–265. Association for Computational Linguistics.

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021. WikiMatrix: Mining 135M parallel sentences in 1620 language pairs from Wikipedia. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361. Association for Computational Linguistics.

Danil Semin and Ondřej Bojar. 2024. CUNI-DS submission: A naive transfer learning setup for english-to-russian translation utilizing english-to-czech data. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Samuel A Stouffer, Edward A Suchman, Leland C DeVinney, Shirley A Star, and Robin M Williams Jr. 1949. The american soldier: Adjustment during army life.(studies in social psychology in world war ii), vol. 1.

Roman Sudarikov, Martin Popel, Ondřej Bojar, Aljoscha Burchardt, and Ondřej Klejch. 2016. Using MT-ComparEval. In *Translation Evaluation: From Fragmented Tools and Data Sets to an Integrated Ecosystem*, pages 76–82.

Zhifan Sun and Antonio Valerio Miceli Barone. 2024. Scaling behavior of machine translation with large language models under prompt injection attacks. In *Proceedings of the First edition of the Workshop on the Scaling Behavior of Large Language Models (SCALE-LLM 2024)*, pages 9–23, St. Julian's, Malta. Association for Computational Linguistics.

Shaomu Tan, David Stap, Seth Aycock, Christof Monz, and Di Wu. 2024. Uva-MT's participation in the WMT24 general translation shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Gemini Team. 2024a. Gemini: A family of highly capable multimodal models.

Llama-3 Team. 2024b. The Llama 3 Herd of Models.

Phi-3 Team. 2024c. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218. European Language Resources Association (ELRA).

Antonio Toral, Sheila Castilho, Ke Hu, and Andy Way. 2018. Attaining the unattainable? reassessing claims of human parity in neural machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 113–123. Association for Computational Linguistics.

Christopher Lemmer Webber, Jessica Tallon, Erin Shepherd, Amy Guy, and Evan Prodromou. 2018. ActivityPub, W3C Recommendation. Technical report, W3C.

Johnny Wei and Robin Jia. 2021. The statistical advantage of automatic NLG metrics at the system level. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6840–6854. Association for Computational Linguistics.

Frank Wilcoxon. 1945. Individual comparisons by ranking methods.

Zhanglin Wu, Daimeng Wei, Zongyao Li, Hengchao Shang, Jiaxin GUO, Shaojun Li, Zhiqiang Rao, Yuanchang Luo, Ning Xie, and Hao Yang. 2024. Choose the final translation from NMT and LLM hypotheses using MBR decoding: HW-TSC's submission to the WMT24 general MT shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

L Xue. 2020. mt5: A massively multilingual pretrained text-to-text transformer. *arXiv preprint arXiv:2010.11934*.

Wenbo Zhang. 2024. IOL research machine translation systems for WMT24 general machine translation shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232.

Michał Ziemski, Marcin Junczys-Dowmunt, and Bruno Pouliquen. 2016. The United Nations parallel corpus v1.0. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3530–3534. European Language Resources Association (ELRA).

Hao Zong, Chao Bei, Huan Liu, Conghu Yuan, Wentao Chen, and Degen Huang. 2024. DLUT and GTCOM's neural machine translation systems for WMT24. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Vilém Zouhar, Shuoyang Ding, Anna Currey, Tatyana Badeka, Jenyuan Wang, and Brian Thompson. 2024a. Fine-tuned machine translation metrics struggle in unseen domains. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 488–500, Bangkok, Thailand. Association for Computational Linguistics.

Vilém Zouhar, Věra Kloudová, Martin Popel, and Ondřej Bojar. 2024b. Evaluating optimal reference translations. *Natural Language Processing*, page 1–24.

## A  Error Span Annotation Miscellaneous

### A.1  Annotation Guidelines

**Higlighting errors:**  Highlight the text fragment where you have identified a translation error (drag or click start & end). Click repeatedly on the highlighted fragment to increase its severity level or to remove the selection.
• **Minor Severity:** Style/grammar/lexical choice could be better/more natural.
• **Major Severity:** Seriously changed meaning, difficult to read, decreases usability.
If something is missing from the text, mark it as an error on the **[MISSING]** word. The highlights do not have to have character-level precision. It's sufficient if you highlight the word or rough area where the error appears. Each error should have a separate highlight.

**Score:**  After highlighting all errors, please set the overall segment translation scores. The quality levels associated with numerical scores on the slider:
• 0%: No meaning preserved: Nearly all information is lost in the translation.
• 33%: Some meaning preserved: Some of the meaning is preserved but significant parts are missing. The narrative is hard to follow due to errors. Grammar may be poor.
• 66%: Most meaning preserved and few grammar mistakes: The translation retains most of the meaning. It may have some grammar mistakes or minor inconsistencies.
• 100%: Perfect meaning and grammar: The meaning and grammar of the translation is completely consistent with the source.

### A.2  Changes to Interface

Since the original study of Kocmi et al. (2024b), we used an updated version of the interface. Apart from minor quality of life changes, a noticeable change is the addition of a pop-up bubble that shows the exact score of the segment (see Figure 5). While it appears as a minor change, it might change the annotator behavior that prefer for example certain numbers, as annotators did in translation evaluation study of Zouhar et al. (2024b).



**Figure 5:** Interacting with the score slider shows the exact score to the annotator in the updated ESA interface.

| Dataset | Segments | Tokens | | Characters | |
|---|---|---|---|---|---|
| | | Source | Target | Source | Target |
| **Czech→Ukrainian** | Segs | Czech | Ukrainian | Czech | Ukrainian |
| OPUS | 9.8M | 103.0M | 102.9M | 752.0M | 1.3B |
| Facebook-wikimatrix-1 | 849.0k | 10.4M | 10.1M | 76.0M | 127.3M |
| ELRC | 130.0k | 2.5M | 2.6M | 19.6M | 35.3M |
| (Total) | 10.8M | 115.9M | 115.6M | 847.6M | 1.4B |
| **English→Czech** | Segs | English | Czech | English | Czech |
| ParaCrawl-paracrawl-9 | 50.6M | 692.1M | 626.3M | 4.3B | 4.7B |
| Facebook-wikimatrix-1 | 2.1M | 33.6M | 29.7M | 206.8M | 216.6M |
| Tilde | 2.1M | 42.3M | 38.3M | 276.5M | 303.7M |
| Statmt-europarl-10 | 644.4k | 15.6M | 13.0M | 94.3M | 98.1M |
| Statmt-wikititles-3 | 410.9k | 1.0M | 965.6k | 7.5M | 7.6M |
| Statmt-news_commentary-18.1 | 265.4k | 5.7M | 5.2M | 36.2M | 39.8M |
| Statmt-commoncrawl_wmt13-1 | 161.8k | 3.3M | 2.9M | 20.7M | 20.7M |
| (Total) | 56.3M | 793.7M | 716.3M | 5.0B | 5.4B |
| **English→German** | Segs | English | German | English | German |
| ParaCrawl-paracrawl-9 | 278.3M | 4.3B | 4.0B | 26.4B | 29.5B |
| Facebook-wikimatrix-1 | 6.2M | 100.5M | 97.0M | 623.7M | 701.2M |
| Tilde | 5.2M | 107.4M | 102.7M | 698.6M | 822.1M |
| Statmt-commoncrawl_wmt13-1 | 2.4M | 51.4M | 47.0M | 314.2M | 340.5M |
| Statmt-europarl-10 | 1.8M | 45.5M | 42.4M | 272.9M | 312.1M |
| Statmt-wikititles-3 | 1.5M | 3.6M | 3.1M | 26.5M | 25.5M |
| Statmt-news_commentary-18.1 | 437.5k | 9.6M | 9.8M | 61.2M | 74.3M |
| (Total) | 295.9M | 4.6B | 4.3B | 28.4B | 31.7B |
| **English→Hindi** | Segs | English | Hindi | English | Hindi |
| AllenAi-nllb-1 | 33.2M | 327.0M | 311.6M | 1.8B | 3.8B |
| OPUS | 12.1M | 147.6M | 165.7M | 919.3M | 2.2B |
| AI4Bharath-samananthar-0.2 | 8.5M | 135.8M | 152.3M | 819.0M | 2.0B |
| Statmt-ccaligned-1 | 8.2M | 114.5M | 129.8M | 724.3M | 1.7B |
| Anuvaad | 3.0M | 58.5M | 61.6M | 359.5M | 836.2M |
| IITB-hien_train-1.5 | 1.6M | 19.8M | 21.4M | 114.7M | 283.6M |
| Facebook-wikimatrix-1 | 696.1k | 12.0M | 13.5M | 74.0M | 182.4M |
| Statmt-pmindia-1 | 56.8k | 1.1M | 1.2M | 6.7M | 16.6M |
| JoshuaDec-indian_training-1 | 37.7k | 562.6k | 659.1k | 3.4M | 8.9M |
| Neulab-tedtalks_train-1 | 18.8k | 372.6k | 491.2k | 1.9M | 4.4M |
| Statmt-news_commentary-18.1 | 4.9k | 149.7k | 167.7k | 963.6k | 2.3M |
| ELRC | 245 | 4.9k | 6.3k | 31.6k | 85.7k |
| (Total) | 67.3M | 817.3M | 858.4M | 4.9B | 11.1B |
| **English→Icelandic** | Segs | English | Icelandic | English | Icelandic |
| OPUS | 16.4M | 174.9M | 166.5M | 1.0B | 1.1B |
| ParaCrawl-paracrawl-9 | 3.0M | 45.1M | 42.7M | 266.1M | 292.2M |
| ParIce-eea_train-20.05 | 1.7M | 26.7M | 24.2M | 170.4M | 179.5M |
| Statmt-ccaligned-1 | 1.2M | 18.6M | 17.8M | 115.6M | 124.4M |
| Tilde | 420.7k | 6.3M | 6.1M | 41.7M | 45.3M |
| ParIce-ema_train-20.05 | 399.1k | 6.1M | 5.9M | 40.4M | 43.9M |
| Facebook-wikimatrix-1 | 313.9k | 5.7M | 4.8M | 34.5M | 34.0M |
| Statmt-wikititles-3 | 50.2k | 99.0k | 88.4k | 722.2k | 763.3k |
| EU | 4.7k | 54.4k | 52.3k | 369.0k | 398.5k |
| (Total) | 23.4M | 283.7M | 268.2M | 1.7B | 1.8B |
| **English→Russian** | Segs | English | Russian | English | Russian |
| Statmt-backtrans_ruen-wmt20 | 39.4M | 746.5M | 596.3M | 4.5B | 7.8B |
| OPUS | 25.2M | 563.8M | 520.7M | 3.7B | 7.3B |
| ParaCrawl-paracrawl-1_bonus | 5.4M | 101.3M | 80.4M | 632.5M | 1.1B |
| Facebook-wikimatrix-1 | 5.2M | 86.8M | 76.5M | 537.7M | 1.0B |
| Statmt-wikititles-3 | 1.2M | 3.1M | 2.9M | 22.8M | 39.3M |
| Statmt-yandex-wmt22 | 1.0M | 21.3M | 18.7M | 131.0M | 250.8M |
| Statmt-commoncrawl_wmt13-1 | 878.4k | 18.8M | 17.4M | 116.2M | 214.6M |
| Statmt-news_commentary-18.1 | 377.7k | 8.7M | 8.1M | 55.7M | 112.1M |
| Tilde | 34.3k | 752.7k | 702.8k | 4.8M | 10.0M |
| (Total) | 78.6M | 1.6B | 1.3B | 9.7B | 17.7B |

**Table 10:** Statistics for parallel training data provided for General/News Translation Task. Suffixes, k, M, and B, are short for thousands, millions, and billions, respectively.

| Dataset | Segments | Tokens | | Characters | |
|---|---|---|---|---|---|
| | | Source | Target | Source | Target |
| **English→Spanish** | Segs | English | Spanish | English | Spanish |
| ParaCrawl-paracrawl-9 | 269.4M | 4.4B | 4.8B | 26.7B | 30.0B |
| OPUS | 223.4M | 4.1B | 4.6B | 26.3B | 30.0B |
| Statmt-ccaligned-1 | 98.4M | 1.2B | 1.3B | 7.7B | 8.6B |
| LinguaTools-wikititles-2014 | 16.6M | 41.3M | 46.0M | 304.8M | 335.2M |
| Facebook-wikimatrix-1 | 6.5M | 120.1M | 137.4M | 742.9M | 854.5M |
| Tilde | 3.8M | 80.0M | 92.9M | 521.0M | 603.4M |
| EU | 3.7M | 70.7M | 80.6M | 457.1M | 519.6M |
| Statmt-europarl-7 | 2.0M | 49.1M | 51.6M | 294.5M | 324.6M |
| Statmt-commoncrawl_wmt13-1 | 1.8M | 40.8M | 43.5M | 248.8M | 272.8M |
| Statmt-news_commentary-18.1 | 500.2k | 11.1M | 13.1M | 71.1M | 83.5M |
| Neulab-tedtalks_train-1 | 196.0k | 4.1M | 3.9M | 20.4M | 20.6M |
| (Total) | 626.2M | 10.2B | 11.2B | 63.4B | 71.6B |
| **English→Ukrainian** | Segs | English | Ukrainian | English | Ukrainian |
| ParaCrawl-paracrawl-1_bonus | 13.4M | 505.8M | 487.5M | 3.3B | 6.0B |
| Statmt-ccaligned-1 | 8.5M | 119.4M | 104.1M | 755.4M | 1.3B |
| Facebook-wikimatrix-1 | 2.6M | 41.5M | 35.6M | 257.6M | 447.3M |
| ELRC | 129.9k | 3.0M | 2.6M | 19.6M | 35.7M |
| Tilde | 1.6k | 36.1k | 34.2k | 238.0k | 477.9k |
| (Total) | 24.6M | 669.8M | 629.8M | 4.3B | 7.8B |
| **English→Japanese** | Segs | English | | English | Japanese |
| KECL-paracrawl-3 | 25.7M | 599.0M | | 3.7B | 4.6B |
| Facebook-wikimatrix-1 | 3.9M | 61.6M | | 379.1M | 455.0M |
| StanfordNLP-jesc_train-1 | 2.8M | 19.3M | | 104.0M | 119.6M |
| Statmt-wikititles-3 | 757.0k | 1.9M | | 14.0M | 18.7M |
| Phontron-kftt_train-1 | 440.3k | 9.7M | | 59.9M | 49.1M |
| Statmt-ted-wmt20 | 241.7k | 4.0M | | 23.0M | 27.3M |
| Statmt-news_commentary-18.1 | 1.9k | 40.3k | | 253.2k | 318.5k |
| (Total) | 33.9M | 695.7M | | 4.3B | 5.2B |
| **English→Chinese** | Segs | English | | English | Chinese |
| Statmt-backtrans_enzh-wmt20 | 19.8M | 364.2M | | 2.2B | 2.0B |
| OPUS | 17.5M | 417.3M | | 2.7B | 2.1B |
| ParaCrawl-paracrawl-1_bonus | 14.2M | 217.6M | | 1.3B | 1.2B |
| Facebook-wikimatrix-1 | 2.6M | 49.9M | | 311.1M | 277.8M |
| Statmt-wikititles-3 | 922.0k | 2.4M | | 17.8M | 16.3M |
| Statmt-news_commentary-18.1 | 442.9k | 9.8M | | 62.7M | 55.2M |
| (Total) | 55.3M | 1.1B | | 6.6B | 5.6B |
| **Japanese→Chinese** | Segs | | | Japanese | Chinese |
| OPUS | 19.6M | | | 1.4B | 1.1B |
| KECL-paracrawl-2wmt24 | 4.6M | | | 1.0B | 705.0M |
| LinguaTools-wikititles-2014 | 1.7M | | | 35.2M | 27.5M |
| Facebook-wikimatrix-1 | 1.3M | | | 145.1M | 113.6M |
| KECL-paracrawl-2 | 83.9k | | | 18.9M | 14.1M |
| Neulab-tedtalks_train-1 | 5.2k | | | 490.9k | 376.0k |
| Statmt-news_commentary-18.1 | 1.6k | | | 272.8k | 197.3k |
| (Total) | 27.2M | | | 2.6B | 1.9B |

**Table 11:** Training dataset statistics (continued from Table 10 on previous page).

## B  System Submission Summaries

This section lists all the submissions to the translation task and provides the authors' descriptions of their submission.

### B.1  AIST-AIRC (Rikters and Miwa, 2024)

At WMT 2024 AIST AIRC participated in the General Machine Translation shared task and the Biomedical Translation task (Neves et al., 2024). We trained constrained track models for translation between English, German, and Japanese. Before training the final models, we first filtered the parallel data, then performed iterative back-translation as well as parallel data distillation. We experimented with training baseline Transformer models, Mega models, and fine-tuning open-source T5 and Gemma model checkpoints using the filtered parallel data. Our primary submissions contain translations from ensembles of two Mega model checkpoints and our contrastive submissions are generated by our fine-tuned T5 model checkpoints.

### B.2  AMI (Jasonarson et al., 2024)

This paper presents the submission of the Arni Magnusson Institute's team to the WMT24 General translation task. We work on the English→Icelandic translation direction. Our system comprises four translation models and a grammar correction model. For training our systems we carefully curate our datasets, aggressively filtering out sentence pairs that may detrimentally affect the quality of our systems output. Some of our data are collected from human translations and some are synthetically generated. A part of the synthetic data is generated using an LLM, and we find that it increases the translation capability of our system significantly.

### B.3  CUNI-DS (Semin and Bojar, 2024)

We present a naive transfer learning approach for English-to-Russian translation, leveraging English-to-Czech data within the constrained track of WMT24. Utilizing the Mistral-7B-0.1 model in its 4-bit quantized variant, we employ QLoRA adapter training. The approach involves two phases: first, training the adapters on the English-to-Czech CzEng 2.0 dataset, followed by fine-tuning the adapters further for English-to-Russian translation with additional corpora. The training spans a total of 48 hours. Evaluation is performed using WMT22 and WMT23 datasets, including the paragraph-level version of the latter.
Phase 1: Training on English-to-Czech Data
    Dataset: CzEng 2.0, with examples packed into chunks of sequence length 2048.
    Parameters: Warmup Steps: 20, Learning Rate: 2e-5, Weight Decay: 1e-2, Cumulative Batch Size: 32
    Instructions: Alpaca-like instructions
    Duration: 24 hours on a single A100 GPU, using the Unsloth library.
Phase 2: Fine-Tuning for English-to-Russian
    Data: Yandex Corpus and News Commentary v18.1, with the latter divided into chunks of 10 sentences.
    Regimen: Training with parameters similar to Phase 1.
    Duration: An additional 24 hours, totaling 48 hours of training.

### B.4  CUNI-{Transformer, DocTransformer, GA, MH, NL} (Hrabal et al., 2024)

This paper presents the contributions of Charles University teams to the WMT24 General Translation task (English to Czech, German and Russian, and Czech to Ukrainian), and the WMT24 Translation into Low-Resource Languages of Spain task.

Our most elaborate submission, CUNI-MH for English→Czech, is the result of fine-tuning Mistral 7B v0.1 for translation using a three-stage process: Supervised fine-tuning using QLoRA, Contrastive Preference Optimization, and merging of model checkpoints. We also describe the CUNI-GA, CUNI-Transformer and CUNI-DocTransformer submissions, which are based on our systems from the previous year.

Our en2ru system CUNI-DS uses a similar first stage as CUNI-MH (QLoRA for English→Czech) and follows with transferring to en2ru.

For en2de (CUNI-NL), we experimented with a LLM-based speech translation system, to translate without the speech input.

For the Translation into Low-Resource Languages of Spain task, we performed QLoRA fine-tuning of a large LLM on a small amount of synthetic (backtranslated) data.

### B.5  CycleL and CycleL2 (Dreano et al., 2024)

CycleGN is a fully self-supervised Neural Machine Translation framework relying on the Transformer architecture that does not require parallel data. Its approach is similar to a Discriminator-less CycleGAN, hence the "non-adversarial" name, specifically tailored for non-parallel text datasets. The foundational concept of our research posits that in an ideal scenario, retro-translations of generated translations should revert to the original source sentences. Consequently, a pair of models can be trained using a Cycle Consistency Loss (CCL) only, with one model translating in one direction and the second model in the opposite direction.

In the context of this research, two sub-categories of non-parallel datasets are introduced. A "permuted" dataset is defined as a parallel dataset wherein the sentences of one language have been systematically rearranged. Consequently, this results in a non-parallel corpus where it is guaranteed that each sentence has a corresponding translation located at an unspecified index within the dataset. A "non-intersecting" dataset is a non-parallel dataset for which it is guaranteed that no sentence has an exact translation.

Masked Language Modeling (MLM) is a pre-training strategy implemented in BERT, where a specified proportion of the input tokens are substituted with a unique $mask$ token. The objective of the neural network under this paradigm is to accurately reconstruct the original sentence from this degraded input.

In inference mode, Transformers are able to generate sentences without labels. Thus, the first step is to generate pseudo-labels in inference, that are then used as labels during training. However, the models consistently converge towards a trivial solution in which the input, the generated pseudo-labels and the output are identical, achieving an optimal outcome on the CCL function, registering a value of zero. CycleGN demonstrates how MLM pre-training can be leveraged to move away from this trivial path and perform actual text translation.

As a contribution to the WMT24 challenge, this study explores the efficacy of the CycleGN architectural framework in learning translation tasks across eleven language pairs under the permuted condition and four under the non-intersecting condition.

Moreover, two additional language pairs from the previous WMT edition were trained and the evaluations demonstrate the robust adaptability of CycleGN in learning translation tasks.

### B.6  DLUT-GTCOM (Zong et al., 2024)

This paper presents the submission from Global Tone Communication Co., Ltd. and Dalian University of Technology for the WMT24 shared general Machine Translation (MT) task at the Conference on Empirical Methods in Natural Language Processing (EMNLP). Our participation encompasses two language pairs: English to Japanese and Japanese to Chinese. The systems are developed without particular constraints or requirements, facilitating extensive research in machine translation. We emphasize back-translation, utilize multilingual translation models, and apply fine-tuning strategies to improve performance. Additionally, we integrate both human-generated and machine-generated data to fine-tune our models, leading to enhanced translation accuracy. The automatic evaluation results indicate that our system ranks first in terms of BLEU score for the Japanese to Chinese translation.

### B.7  HW-TSC (Wu et al., 2024)

This paper presents the submission of Huawei Translate Services Center (HW-TSC) to the WMT24 general machine translation (MT) shared task, where we participate in the English to Chinese (en→zh) language pair. Similar to previous years' work, we use training strategies such as regularized dropout, bidirectional training, data diversification, forward translation, back translation, alternated training, curriculum learning, and transductive ensemble learning to train the neural machine translation (NMT) model based on the deep Transformer-big architecture. The difference is that we also use continue pre-training, supervised fine-tuning, and contrastive preference optimization to train the large language model (LLM) based MT

model. By using Minimum Bayesian risk (MBR) decoding to select the final translation from multiple hypotheses for NMT and LLM-based MT models, our submission receives competitive results in the final evaluation.

### B.8  IKUN and IKUN-C ([Liao et al., 2024](#))

This paper introduces two multilingual systems, IKUN and IKUN-C, developed for the general machine translation task in WMT24. IKUN and IKUN-C represent an open system and a constrained system, respectively, built on Llama-3-8b and Mistral-7B-v0.3. Both systems are designed to handle all 11 language directions using a single model. According to automatic evaluation metrics, IKUN-C achieved 6 first-place and 3 second-place finishes among all constrained systems, while IKUN secured 1 first-place and 2 second-place finishes across both open and constrained systems. These encouraging results suggest that large language models (LLMs) are nearing the level of proficiency required for effective multilingual machine translation. The systems are based on a two-stage approach: first, continuous pre-training on monolingual data in 10 languages, followed by fine-tuning on high-quality parallel data for 11 language directions. The primary difference between IKUN and IKUN-C lies in their monolingual pre-training strategy. IKUN-C is pre-trained using constrained monolingual data, whereas IKUN leverages monolingual data from the OSCAR dataset. In the second phase, both systems are fine-tuned on parallel data sourced from NTREX, Flores, and WMT16-23 for all 11 language pairs.

### B.9  IOL-Research ([Zhang, 2024](#))

This paper illustrates the submission system of the IOL Research team for the WMT24 General Machine Translation shared task. We submitted translations for all translation directions in the general machine translation task. According to the official track categorization, our system qualifies as an open system due to the utilization of open-source resources in developing our machine translation model. With the growing prevalence of large language models (LLMs) as a conventional approach for managing diverse NLP tasks, we have developed our machine translation system by leveraging the capabilities of LLMs. Overall, We first performed continued pretraining using the open-source LLMs with tens of billions of parameters to enhance the model's multilingual capabilities. Subsequently, we employed open-source Large Language Models, equipped with hundreds of billions of parameters, to generate synthetic data. This data was then blended with a modest quantity of additional open-source data for precise supervised fine-tuning. In the final stage, we also used ensemble learning to improve translation quality.

### B.10  MSLC ([Larkin et al., 2024](#))

The MSLC (Metric Score Landscape Challenge) submissions for English–German, English–Spanish, and Japanese–Chinese are constrained systems built using Transformer models for the purpose of better evaluating metric performance in the WMT24 Metrics Task. They are intended to be representative of the performance of systems that can be built relatively simple using constrained data and with minimal modifications to the translation training pipeline.

### B.11  NTTSU ([Kondo et al., 2024](#))

The NTTSU team's submission leverages several large language models developed through a training procedure that includes continual pre-training and supervised fine-tuning. For paragraph-level translation, we generated synthetic paragraph-aligned data and utilized this data for training.

In the task of translating Japanese to Chinese, we particularly focused on the speech domain translation. Specifically, we built Whisper models for Japanese automatic speech recognition (ASR). We used YODAS dataset for Whisper training. Since this data contained many noisy data pairs, we combined the Whisper outputs using ROVER for polishing the transcriptions. Furthermore, to enhance the robustness of the translation model against errors in the transcriptions, we performed data augmentation by forward translation from audio, using both ASR and base translation models.

To select the best translation from multiple hypotheses of the models, we applied Minimum Bayes Risk decoding + reranking, incorporating scores such as COMET-QE, COMET, and cosine similarity by LaBSE.

### B.12    Occiglot ([Avramidis et al., 2024](#))

This document describes the submission of the very first version of the Occiglot open-source large language model to the General MT Shared Task of the 9th Conference of Machine Translation (WMT24). Occiglot is an open-source, community-based LLM based on Mistral-7B, which went through language-specific continual pre-training and subsequent instruction tuning, including instructions relevant to machine translation. We examine the automatic metric scores for translating the WMT24 test set and provide a detailed linguistically-motivated analysis.

Despite Occiglot performing worse than many of the other system submissions, we observe that it performs better than Mistral7B, which has been based upon, which indicates the positive effect of the language specific continual-pretraining and instruction tuning.

We see the submission of this very early version of the model as a motivation to unite community forces and pursue future LLM research on the translation task.

### B.13    SCIR-MT ([Li et al., 2024](#))

This paper introduces the submission of SCIR research center of Harbin Institute of Technology participating in the WMT24 machine translation evaluation task of constrained track for English to Czech. Our approach involved a rigorous process of cleaning and deduplicating both monolingual and bilingual data, followed by a three-stage model training recipe. During the testing phase, we used the beam serach decoding method to generate a large number of candidate translations. Furthermore, we employed COMET-MBR decoding to identify optimal translations.

### B.14    Team-J ([Kudo et al., 2024](#))

We participated in the constrained track for English-Japanese and Japanese-Chinese translations at the WMT 2024 General Machine Translation Task. Our approach was to generate a large number of sentence-level translation candidates and select the most probable translation using minimum Bayes risk (MBR) decoding and document-level large language model (LLM) re-ranking. We first generated hundreds of translation candidates from multiple translation models and retained the top 30 candidates using MBR decoding. In addition, we continually pre-trained LLMs on the target language corpora to leverage document-level information. We utilized LLMs to select the most probable sentence sequentially in context from the beginning of the document.

### B.15    TranssionMT

Hyper-SNMT represents a cutting-edge approach in the field of machine translation. Hyper-SNMT is based on embedding sentences in a hyperbolic space, where distances naturally reflect language hierarchy and dependencies. This novel embedding space enables the model to achieve more accurate translations, especially for languages with complex grammatical structures and rich morphology. Both speed and accuracy are significantly improved compared to existing models. This submission is highlighting the portential of Hyper-SNMT to revolutionize the field of neural machine translation.

### B.16    TSU-HITs ([Mynka and Mikhaylovskiy, 2024](#))

This paper describes the TSU HITS team's submission system for the WMT'24 general translation task. We focused on exploring the capabilities of discrete diffusion models for the English-to-Russian, German, Czech, Spanish translation tasks in the constrained track. Our submission system consists of a set of discrete diffusion models for each language pair. The main advance is using a separate length regression model to determine the length of the output sequence more precisely.

### B.17    Unbabel-Tower70B ([Rei et al., 2024](#))

In this work, we present Tower v2, an improved iteration of the state-of-the-art open-weight Tower models, and the backbone of our submission to the WMT24 General Translation shared task. Tower v2 introduces key improvements including expanded language coverage, enhanced data quality, and increased model capacity up to 70B parameters. Our final submission combines these advancements with quality-aware decoding strategies, selecting translations based on multiple translation quality signals. The resulting

system demonstrates significant improvement over previous versions, outperforming closed commercial systems like GPT-4o, Claude 3.5, and DeepL even at a smaller 7B scale.

### B.18  UvA-MT (Tan et al., 2024)

Fine-tuning Large Language Models (FT-LLMs) with parallel data has emerged as a promising paradigm in recent machine translation research. In this paper, we explore the effectiveness of FT-LLMs and compare them to traditional encoder-decoder Neural Machine Translation (NMT) systems under the WMT24 general MT shared task across three high-resource directions: English to Chinese, English to Japanese, and Japanese to Chinese. We implement several techniques, including Quality Estimation (QE) data filtering, supervised fine-tuning, and post-editing that integrate NMT systems with LLMs. We demonstrate that fine-tuning LLaMA2 on a high-quality but relatively small bitext dataset (100K) yields COMET results comparable to much smaller encoder-decoder NMT systems trained on over 22 million bitexts. However, this approach largely underperforms on surface-level metrics like BLEU and ChrF. We further control the data quality using the COMET-based quality estimation method. Our experiments show that 1) filtering low COMET scores largely improves encoder-decoder systems, but 2) no clear gains are observed for LLMs when further refining the fine-tuning set. Finally, we show that combining NMT systems with LLMs via post-editing generally yields the best performance in our experiments.

### B.19  Yandex (Elshin et al., 2024)

In this paper, we present the methodology employed by the NLP team at Yandex LLC for participating in the WMT 2024 General MT Translation track, focusing on English-to-Russian translation. Our approach involves training a YandexGPT LLM-based model for translation tasks using a multi-stage process to ensure high-quality and contextually accurate translations.

Initially, we utilize a pre-trained model, trained on a large corpus of high-quality monolingual texts in various languages, crawled from various open sources, not limited to English and Russian. This extensive pre-training allows the model to capture a broad spectrum of linguistic nuances and structures. Following this, the model is fine-tuned on a substantial parallel corpus of high-quality texts collected from diverse open sources, including websites, books, and subtitles. These texts are meticulously aligned at both the sentence and paragraph levels to enhance the model's contextual understanding and translation accuracy.

In the subsequent stage, we employ p-tuning on an internal high-quality corpus of paragraph-aligned data. This step ensures that the model is finely adjusted to handle complex paragraph-level translations with greater fluency and coherence.

Next, we apply the Contrastive Pretraining Objective (CPO) method, as described in the paper CPO, using a human-annotated translation corpus. This stage focuses on refining the model's performance based on metrics evaluated at the paragraph level, emphasizing both the accuracy of the translation and the fluency of the resulting texts. The CPO method helps the model to better distinguish between subtle contextual differences, thereby improving translation quality.

In the final stage, we address the importance of preserving the content structure in translations, which is crucial for the General MT test set. To achieve this, we introduce a synthetic corpus based on web pages and video subtitles, and use it during HE markup finetune training. This encourages the model to maintain the original text's tag structure. This step ensures that the translated output retains the structural integrity of the source web pages, providing a seamless user experience.

Our multi-stage approach, combining extensive pre-training, targeted fine-tuning, advanced p-tuning, and structure-preserving techniques, ensures that our model delivers high-quality, fluent, and structurally consistent translations suitable for practical applications and competitive benchmarks.

## C   Translator Brief

In this project we wish to translate data from several domains for use in evaluation of Machine Translation (MT). The translations produced by you will be compared against the translations produced by a variety of different MT systems. They will be released to the research community to provide a benchmark, or "gold-standard" measure for translation quality. The translation therefore needs to be a high-quality rendering of the source text into the target language, as if it was originally written directly in the target language. However, there are some constraints imposed by the intended usage:

- All translations must be "from scratch", without post-editing from machine translation or usage of CAT tools. Using post-editing would bias the evaluation, so we need to avoid it. We can detect post-editing and will reject translations that are post-edited.

- Translation should preserve the paragraph boundaries but may change number of sentences per paragraph. The source texts contain one paragraph per line and the translations should be the same.

- Translators should avoid inserting parenthetical explanations into the translated text and obviously avoid losing any pieces of information from the source text. We will check the translations for quality and will reject translations that contain errors.

- If the original data contain errors, typos, or other problems, do not change the source sentences, instead try to prepare correct translation as if the error wouldn't be in the source.

- The data contain several domains, each folder containing one domain source.

The source files will be delivered as text files (sometimes known as "notepad" files), with one paragraph per line. We need the translations to be returned in the same format. The translation file needs to have the same name as the original file.

**Speech Domain**   The texts are the transcriptions of audio, edited by native speakers. Each file represents one segment of audio (you are also provided with correspondent audio in WAW format). Phrases said by different speakers are located on different lines. Audios correspond to different domains, they differ in formality, style, topics and number of speakers. The idea is to translate using the most similar language in the target language, matching as best as possible the characteristics of the source text.

**Social domain**   The texts are from the social network Mastodon (similar to Twitter). Each file represents a thread or part of a thread from one or several users. Different posts within a thread are presented on different lines in the file, although individual posts can also span several lines. The sentences have been selected so that they do not contain offensive or sensitive content (hate speech, taking-drugs, suicide, politically sensitive topics, etc.). However, profanities were kept as they were taken to be illustrative of the sociolect of online language. If however, you do not feel comfortable with translating something, please leave the whole line blank and let us know that you have not translated it. The texts are particular in that they may contain spelling errors, slang, acronyms, marks of expressivity, etc. The idea is to translate using the most natural language in the target language, matching as best as possible the style and familiarity of the source text.

- Spelling mistakes should not be preserved in their translations, i.e. the translation should be spelt correctly

- Marks of expressivity (e.g. asterisks *wow*, capitals letters WOW) should be conserved as best as possible. However, we recommend not to attempt to reproduce repeated characters (e.g. woooow) in translation, as the choice as to which character to repeat is often arbitrary.

- There will be abbreviations and acronyms (e.g. btw -> by the way, fwiw -> for what it's worse). These do not need to be translated using abbreviation or acronyms unless an abbreviation/acronym is the best translation choice in the target language.

- Users have been pseudo-anonymised (e.g. @user1, @user2). These should be left as they are, i.e. not translated.

- Platform-specific elements such as hashtags should be translated as hashtags, but the content should be translated as appropriate into the target language.

- Punctuation can be added if it necessary to avoid comprehension difficulties. Otherwise we recommend following the punctuation of the source text.

A file entitled README-social-domain-translation-notes.pdf has been distributed with the texts to translate. This file should not be translated. It contains some notes to provide additional context on the topic and terms used in some of the texts.

# D   Official Ranking Results (extends Section 7.4)

**Results tables legend**

The human score is the macro-average of human judgements, grouped by domain. The rank takes into consideration head-to-head wins and losses. AutoRank is calculated from automatic metrics.

Ranking and clustering on human scores is done using Wilcoxon signed rank test for each domain separately and final p-value is combined via Stouffer's Z-score method to align with macro average for human score.

Systems are either constrained (white), open-track (light gray), or closed-track (dark gray).

LLMs that do not officially claim a support a language pair are marked with §.

Human scores for individual domains are marked by an up arrow ↑ if their difference from the system domain score is larger than the standard deviation over all domains for given system (row) and down arrow ↓ indicates that the domain score is worse than the overall score.

Underlined domain scores indicate that the domain score is better than the domain score of system above it (of a better ranked system).

|  |  |  |  | Czech→Ukrainian |  |  |  |  |  |  |
| Rank | System | Human | AutoRank | CometKiwi | MetricX | education | news | official | personal | voice |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 1-2 | Claude-3.5 § | 93.0 | 1.7 | -0.7 | 1.0 | ↓ 90.4 | 91.7 | ↑ 95.3 | ↑ 95.4 | 92.2 |
| 2-2 | HUMAN-A | 92.7 | - | - | - | 92.6 | 93.0 | 92.0 | ↑ 94.9 | ↓ 91.1 |
| 3-3 | Gemini-1.5-Pro | 92.6 | 2.0 | -0.7 | 1.2 | ↓ 88.6 | 94.7 | 94.5 | 93.6 | 91.9 |
| 3-4 | Unbabel-Tower70B | 92.2 | 1.0 | -0.7 | 0.9 | ↓ 86.8 | 93.5 | 94.8 | 94.1 | 91.8 |
| 5-5 | IOL-Research | 90.2 | 1.9 | -0.7 | 1.3 | ↓ 80.8 | 89.9 | 92.7 | 94.6 | 93.0 |
| 6-7 | CommandR-plus § | 89.7 | 1.9 | -0.7 | 1.3 | ↓ 83.4 | 89.6 | ↑ 93.8 | 92.1 | 89.4 |
| 6-8 | ONLINE-W | 88.7 | 2.3 | -0.7 | 1.4 | ↓ 84.4 | 89.4 | 87.9 | ↑ 91.3 | 90.4 |
| 7-9 | GPT-4 § | 88.6 | 2.0 | -0.7 | 1.4 | ↓ 83.2 | 87.9 | 89.0 | ↑ 92.4 | 90.3 |
| 8-9 | IKUN | 87.1 | 2.3 | -0.7 | 1.6 | ↓ 77.6 | 86.8 | 89.7 | 91.2 | 90.3 |
| 10-10 | Aya23 | 86.6 | 2.5 | -0.7 | 1.9 | ↓ 77.4 | 91.1 | 88.5 | 87.6 | 88.3 |
| 11-11 | CUNI-Transformer | 85.3 | 3.0 | -0.6 | 2.0 | ↓ 83.2 | 85.2 | 84.8 | ↑ 88.0 | 85.3 |
| 12-12 | IKUN-C | 82.6 | 3.0 | -0.6 | 2.4 | 79.6 | ↓ 70.0 | 87.2 | 88.4 | 87.8 |
| | Mistral-Large § | - | 2.3 | - | - | - | - | - | - | - |
| | TranssionMT | - | 2.6 | - | - | - | - | - | - | - |
| | ONLINE-B | - | 2.6 | - | - | - | - | - | - | - |
| | ONLINE-A | - | 2.6 | - | - | - | - | - | - | - |
| | Llama3-70B § | - | 2.6 | - | - | - | - | - | - | - |
| | ONLINE-G | - | 2.8 | - | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 9.1 | - | - | - | - | - | - | - |
| | BJFU-LPT | - | 11.5 | - | - | - | - | - | - | - |
| | CycleL | - | 21.0 | - | - | - | - | - | - | - |

## English→Czech

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|---|---|---|---|---|---|---|---|---|---|
| 1-2 | HUMAN-A | 92.9 | - | - | - | 93.1 | ↑ 94.5 | 92.0 | 92.1 |
| 2-2 | Unbabel-Tower70B | 91.6 | 1.0 | -0.7 | 1.8 | 91.7 | 94.1 | 93.3 | ↓ 87.5 |
| 2-3 | Claude-3.5 § | 91.2 | 2.1 | -0.7 | 2.4 | 91.2 | ↑ 94.9 | 91.6 | ↓ 87.2 |
| 4-5 | ONLINE-W | 89.0 | 2.8 | -0.7 | 2.8 | 91.0 | ↑ 92.1 | 88.2 | ↓ 84.9 |
| 4-6 | CUNI-MH | 88.4 | 2.1 | -0.7 | 2.3 | 89.7 | ↑ 91.9 | 88.0 | ↓ 84.1 |
| 6-6 | Gemini-1.5-Pro | 88.2 | 2.6 | -0.7 | 2.8 | 88.6 | 89.3 | ↓ 85.2 | 89.6 |
| 6-8 | GPT-4 § | 87.7 | 2.6 | -0.7 | 2.9 | ↓ 85.2 | 89.5 | ↑ 90.1 | 86.1 |
| 8-8 | CommandR-plus § | 86.9 | 2.9 | -0.7 | 2.9 | ↓ 85.2 | 87.5 | ↑ 88.6 | 86.2 |
| 8-9 | IOL-Research | 86.5 | 2.8 | -0.7 | 3.0 | 84.7 | ↑ 90.4 | 86.3 | 84.5 |
| 10-11 | SCIR-MT | 85.4 | 3.2 | -0.7 | 3.3 | 85.0 | ↑ 92.4 | 82.2 | 82.1 |
| 10-11 | CUNI-DocTransformer | 84.3 | 4.4 | -0.6 | 4.0 | 83.1 | ↑ 90.7 | 80.9 | 82.4 |
| 12-12 | Aya23 | 84.2 | 4.3 | -0.6 | 4.0 | 81.6 | ↑ 89.9 | 84.9 | ↓ 80.3 |
| 13-13 | CUNI-GA | 82.1 | 2.3 | -0.7 | 3.7 | 82.8 | ↑ 88.5 | 81.7 | ↓ 75.3 |
| 14-14 | IKUN | 81.7 | 3.9 | -0.6 | 3.7 | 80.2 | ↑ 87.0 | 82.2 | ↓ 77.5 |
| 15-15 | Llama3-70B § | 77.4 | 4.1 | -0.6 | 4.0 | ↓ 65.4 | 83.0 | 82.4 | 78.8 |
| 16-16 | IKUN-C | 75.4 | 4.7 | -0.6 | 4.3 | ↓ 70.5 | 77.7 | 77.5 | 75.7 |
| | TranssionMT | - | 3.5 | - | - | - | - | - | - |
| | ONLINE-A | - | 3.6 | - | - | - | - | - | - |
| | Mistral-Large § | - | 3.7 | - | - | - | - | - | - |
| | ONLINE-B | - | 4.0 | - | - | - | - | - | - |
| | CUNI-Transformer | - | 4.7 | - | - | - | - | - | - |
| | ONLINE-G | - | 5.7 | - | - | - | - | - | - |
| | NVIDIA-NeMo | - | 7.6 | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 15.0 | - | - | - | - | - | - |
| | TSU-HITs | - | 19.5 | - | - | - | - | - | - |
| | CycleL2 | - | 24.2 | - | - | - | - | - | - |
| | CycleL | - | 27.0 | - | - | - | - | - | - |

## English→German

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|---|---|---|---|---|---|---|---|---|---|
| 1-11 | GPT-4 | -1.6 | 1.8 | -0.7 | 1.4 | -0.7 | -1.4 | -0.9 | ↓ -3.6 |
| 1-7 | Dubformer | -1.8 | 1.8 | -0.7 | 1.2 | -1.2 | -1.3 | -0.6 | ↓ -4.2 |
| 2-10 | ONLINE-B | -1.9 | 1.8 | -0.7 | 1.4 | -1.3 | -1.5 | -1.2 | ↓ -3.6 |
| 2-10 | TranssionMT | -1.9 | 1.8 | -0.7 | 1.4 | -1.3 | -1.2 | -1.2 | ↓ -3.9 |
| 2-9 | Unbabel-Tower70B | -1.9 | 1.0 | -0.7 | 1.1 | -1.4 | -2.0 | ↑ -0.8 | ↓ -3.5 |
| 1-9 | HUMAN-B | -2.0 | - | - | - | -0.8 | -1.4 | -0.8 | ↓ -4.9 |
| 2-12 | Mistral-Large | -2.1 | 2.0 | -0.7 | 1.5 | -1.5 | -1.9 | -1.1 | ↓ -3.9 |
| 4-11 | CommandR-plus | -2.3 | 2.0 | -0.7 | 1.4 | -1.7 | -2.4 | ↑ -1.1 | ↓ -3.9 |
| 8-10 | ONLINE-W | -2.3 | 2.2 | -0.7 | 1.5 | -2.1 | -1.3 | -1.7 | ↓ -4.1 |
| 2-12 | Claude-3.5 | -2.4 | 1.9 | -0.7 | 1.4 | -1.1 | -1.0 | -1.2 | ↓ -6.0 |
| 3-13 | HUMAN-A | -2.5 | - | - | - | -2.0 | -1.8 | -1.0 | ↓ -5.0 |
| 10-12 | IOL-Research | -2.5 | 2.3 | -0.7 | 1.6 | -2.0 | -1.7 | -1.6 | ↓ -4.9 |
| 5-13 | Gemini-1.5-Pro | -2.8 | 2.2 | -0.7 | 1.5 | ↓ -5.0 | ↑ -1.3 | -1.9 | -2.9 |
| 14-15 | Aya23 | -3.2 | 2.7 | -0.7 | 1.8 | -2.3 | -2.7 | -2.2 | ↓ -5.7 |
| 14-17 | ONLINE-A | -3.5 | 3.0 | -0.7 | 1.8 | -2.8 | -1.9 | -2.3 | ↓ -6.9 |
| 15-17 | Llama3-70B § | -4.3 | 2.5 | -0.7 | 1.7 | -4.8 | -2.9 | ↑ -2.3 | ↓ -7.1 |
| 15-17 | IKUN | -4.3 | 3.0 | -0.7 | 1.8 | -3.5 | -4.3 | ↑ -2.4 | ↓ -7.1 |
| 18-18 | IKUN-C | -6.1 | 3.8 | -0.6 | 2.0 | -7.6 | -3.4 | -3.3 | ↓ -9.9 |
| 19-19 | MSLC | -15.5 | 11.9 | -0.4 | 4.4 | -15.3 | -11.5 | ↑ -8.2 | ↓ -26.8 |
| | Phi-3-Medium § | - | 3.4 | - | - | - | - | - | - |
| | ONLINE-G | - | 3.5 | - | - | - | - | - | - |
| | CUNI-NL | - | 4.2 | - | - | - | - | - | - |
| | AIST-AIRC | - | 7.2 | - | - | - | - | - | - |
| | NVIDIA-NeMo | - | 7.4 | - | - | - | - | - | - |
| | Occiglot | - | 8.2 | - | - | - | - | - | - |
| | TSU-HITs | - | 13.3 | - | - | - | - | - | - |
| | CycleL2 | - | 27.0 | - | - | - | - | - | - |
| | CycleL | - | 27.0 | - | - | - | - | - | - |

**English→Spanish**

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|------|--------|-------|----------|-----------|---------|----------|------|--------|--------|
| 1-1 | HUMAN-A | 95.3 | - | - | - | 95.2 | ↑ 96.2 | 95.5 | ↓ 94.1 |
| 2-2 | Dubformer | 93.4 | 2.0 | -0.7 | 2.2 | 95.3 | 94.5 | 94.4 | ↓ 89.4 |
| 3-4 | GPT-4 | 91.9 | 1.9 | -0.7 | 2.5 | 93.5 | 94.0 | 93.2 | ↓ 87.0 |
| 4-7 | IOL-Research | 91.4 | 2.3 | -0.7 | 2.8 | ↑ 96.3 | 92.5 | 90.9 | ↓ 86.0 |
| 5-8 | Mistral-Large | 89.3 | 2.2 | -0.7 | 2.7 | 90.5 | 90.4 | 91.0 | ↓ 85.2 |
| 5-9 | Unbabel-Tower70B | 88.9 | 1.0 | -0.7 | 1.9 | 86.2 | ↑ 93.7 | 91.1 | ↓ 84.6 |
| 3-8 | Claude-3.5 | 88.8 | 2.1 | -0.7 | 2.6 | 91.5 | 92.8 | 90.4 | ↓ 80.5 |
| 5-8 | Gemini-1.5-Pro | 88.8 | 2.4 | -0.7 | 2.8 | 89.6 | ↑ 92.3 | 87.0 | ↓ 86.2 |
| 7-9 | CommandR-plus | 88.3 | 2.1 | -0.7 | 2.6 | 88.2 | 89.3 | ↑ 90.8 | ↓ 84.8 |
| 9-10 | Llama3-70B § | 87.2 | 2.6 | -0.7 | 3.0 | ↑ 89.4 | 87.1 | 87.9 | ↓ 84.2 |
| 11-11 | ONLINE-B | 85.6 | 2.7 | -0.7 | 3.1 | 87.4 | 88.6 | 86.8 | ↓ 79.4 |
| 12-13 | IKUN | 84.7 | 2.8 | -0.7 | 3.3 | 85.4 | ↑ 92.4 | 82.8 | ↓ 78.3 |
| 12-13 | IKUN-C | 80.4 | 3.4 | -0.7 | 3.5 | 83.3 | ↑ 85.6 | 79.0 | ↓ 73.6 |
| 14-14 | MSLC | 63.9 | 7.4 | -0.5 | 6.4 | 59.3 | ↑ 78.8 | 55.9 | 61.7 |
| | ONLINE-W | - | 2.7 | - | - | - | - | - | - |
| | TranssionMT | - | 2.8 | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 3.0 | - | - | - | - | - | - |
| | ONLINE-A | - | 3.0 | - | - | - | - | - | - |
| | Aya23 | - | 3.1 | - | - | - | - | - | - |
| | ONLINE-G | - | 3.2 | - | - | - | - | - | - |
| | NVIDIA-NeMo | - | 4.5 | - | - | - | - | - | - |
| | Occiglot | - | 5.9 | - | - | - | - | - | - |
| | TSU-HITs | - | 16.3 | - | - | - | - | - | - |
| | CycleL | - | 24.0 | - | - | - | - | - | - |

**English→Hindi**

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|------|--------|-------|----------|-----------|---------|----------|------|--------|--------|
| 1-3 | TranssionMT | 91.3 | 1.3 | -0.6 | 3.3 | ↑ 94.0 | 93.0 | 89.8 | ↓ 88.2 |
| 1-4 | Unbabel-Tower70B | 90.5 | 1.0 | -0.7 | 3.1 | 90.9 | ↑ 92.7 | 90.7 | ↓ 87.7 |
| 3-3 | Claude-3.5 § | 90.2 | 1.2 | -0.6 | 3.3 | 95.4 | 93.6 | 91.0 | ↓ 81.1 |
| 3-4 | ONLINE-B | 90.1 | 1.4 | -0.6 | 3.3 | 91.8 | 90.4 | 91.3 | ↓ 86.9 |
| 3-5 | Gemini-1.5-Pro § | 90.0 | 1.6 | -0.6 | 3.6 | 90.3 | ↑ 91.9 | 89.4 | ↓ 88.3 |
| 6-6 | GPT-4 § | 88.5 | 2.1 | -0.6 | 4.5 | 89.9 | 90.4 | 89.2 | ↓ 84.4 |
| 7-8 | HUMAN-A | 88.5 | - | - | - | 88.8 | ↓ 88.1 | ↑ 88.9 | 88.2 |
| 8-8 | IOL-Research | 87.2 | 2.1 | -0.6 | 4.3 | 87.2 | ↑ 88.9 | 87.7 | ↓ 84.9 |
| 8-9 | Llama3-70B § | 86.7 | 2.1 | -0.6 | 4.6 | 86.4 | 87.1 | ↓ 86.1 | 87.1 |
| 10-10 | Aya23 | 84.7 | 3.2 | -0.6 | 5.4 | 83.3 | ↑ 86.9 | ↓ 83.1 | 85.7 |
| 11-11 | IKUN-C | 70.7 | 5.5 | -0.5 | 7.1 | 71.2 | ↓ 59.2 | ↑ 80.2 | 72.4 |
| | CommandR-plus § | - | 2.3 | - | - | - | - | - | - |
| | ONLINE-A | - | 3.5 | - | - | - | - | - | - |
| | ONLINE-G | - | 4.2 | - | - | - | - | - | - |
| | Mistral-Large § | - | 5.0 | - | - | - | - | - | - |
| | NVIDIA-NeMo | - | 5.8 | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 7.4 | - | - | - | - | - | - |
| | IKUN | - | 7.7 | - | - | - | - | - | - |
| | ONLINE-empty | - | 15.3 | - | - | - | - | - | - |
| | CycleL | - | 20.0 | - | - | - | - | - | - |

**English→Icelandic**

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|------|--------|-------|----------|-----------|---------|----------|------|--------|--------|
| 1-1 | HUMAN-A | 93.1 | - | - | - | 92.2 | 92.6 | ↑95.0 | 92.4 |
| 2-3 | Dubformer | 84.3 | 2.5 | -0.7 | 3.4 | 84.1 | 83.1 | ↑87.5 | 82.5 |
| 2-3 | Claude-3.5 § | 81.9 | 2.3 | -0.7 | 3.6 | 80.2 | 83.9 | ↑87.2 | ↓76.4 |
| 4-4 | Unbabel-Tower70B | 80.2 | 1.0 | -0.7 | 2.5 | ↓76.6 | 80.6 | ↑84.3 | 79.2 |
| 5-5 | AMI | 73.3 | 3.7 | -0.7 | 4.9 | ↑75.2 | 72.8 | 74.1 | ↓71.1 |
| 6-6 | IKUN | 71.0 | 3.2 | -0.7 | 4.3 | ↓66.8 | ↑74.7 | 73.6 | 69.1 |
| 7-7 | ONLINE-B | 68.0 | 4.2 | -0.7 | 5.5 | 70.5 | ↓59.4 | ↑74.0 | 67.9 |
| 8-9 | GPT-4 | 66.3 | 3.4 | -0.7 | 4.7 | 66.5 | 65.5 | ↑69.5 | ↓63.9 |
| 8-9 | IKUN-C | 65.2 | 3.7 | -0.7 | 4.9 | ↓59.6 | 68.2 | ↑69.3 | 63.8 |
| 10-10 | IOL-Research | 58.0 | 4.3 | -0.7 | 5.7 | ↓49.4 | 59.6 | 61.4 | 61.4 |
| 11-11 | Llama3-70B § | 41.0 | 6.7 | -0.6 | 8.0 | 39.8 | 40.0 | ↑44.0 | 40.3 |
| | TranssionMT | - | 4.2 | - | - | - | - | - | - |
| | ONLINE-A | - | 5.5 | - | - | - | - | - | - |
| | ONLINE-G | - | 6.9 | - | - | - | - | - | - |
| | CommandR-plus § | - | 9.8 | - | - | - | - | - | - |
| | Mistral-Large § | - | 10.4 | - | - | - | - | - | - |
| | Aya23 § | - | 15.2 | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 16.2 | - | - | - | - | - | - |
| | ONLINE-empty | - | 18.1 | - | - | - | - | - | - |
| | TSU-HITs | - | 19.2 | - | - | - | - | - | - |
| | CycleL | - | 21.0 | - | - | - | - | - | - |

**English→Japanese**

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|------|--------|-------|----------|-----------|---------|----------|------|--------|--------|
| 1-1 | HUMAN-A | 91.8 | - | - | - | 92.4 | 93.0 | ↓89.5 | 92.4 |
| 2-4 | ONLINE-B | 91.1 | 1.4 | -0.8 | 2.4 | 91.7 | ↑92.6 | 91.1 | ↓88.9 |
| 3-4 | CommandR-plus | 91.0 | 1.9 | -0.7 | 2.7 | 92.2 | ↑93.7 | 89.5 | ↓88.5 |
| 4-4 | GPT-4 | 90.8 | 1.7 | -0.7 | 2.7 | ↑91.9 | 91.3 | ↓89.9 | 90.1 |
| 4-5 | Claude-3.5 | 90.8 | 1.5 | -0.7 | 2.3 | 91.4 | ↑92.8 | 91.3 | ↓87.6 |
| 4-7 | Gemini-1.5-Pro | 90.0 | 1.7 | -0.7 | 2.5 | 91.1 | ↑92.2 | ↓88.1 | 88.7 |
| 7-7 | Unbabel-Tower70B | 89.7 | 1.0 | -0.8 | 2.0 | ↓88.2 | ↑91.6 | 89.8 | 89.2 |
| 8-8 | IOL-Research | 89.7 | 2.3 | -0.7 | 3.1 | 91.0 | 90.6 | 90.3 | ↓86.9 |
| 8-9 | Aya23 | 89.7 | 2.3 | -0.7 | 3.1 | 90.1 | ↑92.1 | 88.4 | ↓87.9 |
| 10-10 | NTTSU | 89.4 | 1.9 | -0.7 | 2.6 | 90.0 | ↑93.2 | 88.4 | ↓86.2 |
| 11-11 | Team-J | 88.5 | 1.9 | -0.7 | 2.9 | ↓85.0 | 90.1 | ↑91.3 | 87.5 |
| 12-12 | Llama3-70B § | 86.8 | 2.6 | -0.7 | 3.5 | 89.3 | ↑89.8 | 85.2 | ↓82.7 |
| 13-13 | IKUN-C | 81.7 | 3.9 | -0.7 | 4.3 | ↓77.5 | ↑88.5 | 81.2 | 79.8 |
| | DLUT-GTCOM | - | 2.6 | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 2.8 | - | - | - | - | - | - |
| | ONLINE-W | - | 2.9 | - | - | - | - | - | - |
| | Mistral-Large § | - | 2.9 | - | - | - | - | - | - |
| | ONLINE-A | - | 3.0 | - | - | - | - | - | - |
| | IKUN | - | 3.1 | - | - | - | - | - | - |
| | ONLINE-G | - | 6.4 | - | - | - | - | - | - |
| | AIST-AIRC | - | 6.6 | - | - | - | - | - | - |
| | UvA-MT | - | 6.7 | - | - | - | - | - | - |
| | NVIDIA-NeMo | - | 6.9 | - | - | - | - | - | - |
| | CycleL | - | 24.0 | - | - | - | - | - | - |

**English→Russian**

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|------|--------|-------|----------|-----------|---------|----------|------|--------|--------|
| 1-1 | HUMAN-A | 89.2 | - | - | - | ↑94.0 | 88.3 | 87.7 | 86.6 |
| 2-3 | Dubformer | 89.1 | 1.9 | -0.7 | 2.8 | 90.7 | 88.5 | ↑92.1 | ↓84.9 |
| 3-4 | Claude-3.5 | 88.2 | 2.0 | -0.7 | 3.0 | ↑94.1 | 93.1 | 85.7 | ↓80.0 |
| 3-5 | Unbabel-Tower70B | 88.1 | 1.0 | -0.7 | 2.4 | 87.5 | 91.2 | 90.6 | ↓83.2 |
| 3-7 | Yandex | 87.0 | 1.9 | -0.7 | 2.9 | 89.6 | ↑91.8 | 84.5 | ↓82.0 |
| 6-8 | Gemini-1.5-Pro | 85.5 | 2.3 | -0.7 | 3.2 | ↑90.7 | 84.9 | 83.4 | 82.9 |
| 6-9 | GPT-4 | 85.0 | 2.3 | -0.7 | 3.4 | ↑89.3 | 85.4 | 84.6 | ↓80.7 |
| 6-9 | ONLINE-G | 84.6 | 2.2 | -0.7 | 3.3 | 88.3 | 88.8 | 84.6 | ↓76.6 |
| 5-9 | CommandR-plus § | 84.3 | 2.4 | -0.7 | 3.4 | 86.7 | 84.5 | 85.7 | ↓80.5 |
| 10-10 | IOL-Research | 82.1 | 2.6 | -0.7 | 3.7 | 84.8 | 86.4 | 84.2 | ↓73.1 |
| 11-11 | IKUN | 79.2 | 3.2 | -0.7 | 4.1 | 80.2 | ↑87.2 | 78.5 | ↓70.9 |
| 12-12 | Aya23 | 78.6 | 3.3 | -0.7 | 4.2 | 77.8 | ↑82.9 | 78.5 | ↓75.3 |
| 13-13 | Llama3-70B § | 75.7 | 3.1 | -0.7 | 4.1 | 77.0 | ↑80.1 | 76.3 | ↓69.5 |
| 14-14 | IKUN-C | 69.8 | 3.9 | -0.6 | 4.7 | 65.1 | ↑78.3 | 72.9 | ↓62.6 |
| | ONLINE-W | - | 2.6 | - | - | - | - | - | - |
| | Mistral-Large § | - | 2.7 | - | - | - | - | - | - |
| | ONLINE-B | - | 3.1 | - | - | - | - | - | - |
| | TranssionMT | - | 3.1 | - | - | - | - | - | - |
| | ONLINE-A | - | 3.4 | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 3.9 | - | - | - | - | - | - |
| | CUNI-DS | - | 5.9 | - | - | - | - | - | - |
| | NVIDIA-NeMo | - | 7.2 | - | - | - | - | - | - |
| | TSU-HITs | - | 10.8 | - | - | - | - | - | - |
| | CycleL | - | 24.3 | - | - | - | - | - | - |
| | CycleL2 | - | 25.0 | - | - | - | - | - | - |

**English→Ukrainian**

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|------|--------|-------|----------|-----------|---------|----------|------|--------|--------|
| 1-2 | Claude-3.5 | 90.5 | 2.0 | -0.7 | 3.0 | 93.2 | 93.9 | 92.2 | ↓82.7 |
| 1-2 | Unbabel-Tower70B | 89.8 | 1.0 | -0.7 | 2.2 | 92.5 | 92.8 | 91.1 | ↓82.9 |
| 3-3 | Dubformer | 89.0 | 1.8 | -0.7 | 2.7 | ↓84.4 | 91.3 | ↑94.3 | 85.9 |
| 4-6 | HUMAN-A | 87.3 | - | - | - | 89.6 | ↑91.5 | ↓83.8 | 84.1 |
| 4-6 | Gemini-1.5-Pro | 87.1 | 2.2 | -0.7 | 3.0 | ↑90.1 | 88.8 | 85.3 | ↓84.4 |
| 5-8 | ONLINE-W | 86.0 | 2.1 | -0.7 | 2.8 | 86.7 | ↑88.9 | 86.8 | ↓81.8 |
| 5-9 | GPT-4 | 84.6 | 2.3 | -0.7 | 3.3 | 81.2 | ↑90.3 | 84.5 | 82.4 |
| 6-9 | CommandR-plus § | 83.2 | 2.3 | -0.7 | 3.2 | 79.6 | ↑89.1 | 83.6 | 80.4 |
| 7-9 | IOL-Research | 83.2 | 2.4 | -0.7 | 3.4 | 80.6 | ↑90.2 | 83.1 | ↓78.8 |
| 10-10 | IKUN | 78.4 | 2.8 | -0.7 | 3.7 | 83.2 | ↑88.2 | 72.7 | ↓69.7 |
| 11-11 | IKUN-C | 67.9 | 3.9 | -0.6 | 4.7 | ↓65.2 | 69.0 | 68.3 | 69.2 |
| | ONLINE-G | - | 2.3 | - | - | - | - | - | - |
| | Mistral-Large § | - | 2.4 | - | - | - | - | - | - |
| | ONLINE-B | - | 3.1 | - | - | - | - | - | - |
| | TranssionMT | - | 3.1 | - | - | - | - | - | - |
| | Llama3-70B § | - | 3.2 | - | - | - | - | - | - |
| | Aya23 | - | 3.3 | - | - | - | - | - | - |
| | ONLINE-A | - | 3.3 | - | - | - | - | - | - |
| | NVIDIA-NeMo | - | 6.2 | - | - | - | - | - | - |
| | Phi-3-Medium § | - | 11.1 | - | - | - | - | - | - |
| | CycleL | - | 21.0 | - | - | - | - | - | - |

### English→Chinese

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | social | speech |
|---|---|---|---|---|---|---|---|---|---|
| 1-1 | GPT-4 | 89.6 | 2.0 | -0.7 | 3.3 | 88.7 | ↑91.2 | 90.3 | ↓88.4 |
| 2-4 | Unbabel-Tower70B | 89.6 | 1.0 | -0.7 | 2.3 | 90.0 | ↑92.3 | 90.2 | ↓85.8 |
| 2-4 | HUMAN-A | 89.4 | - | - | - | 89.9 | 90.1 | 90.7 | ↓86.8 |
| 4-4 | Gemini-1.5-Pro | 89.3 | 1.8 | -0.7 | 3.1 | 92.0 | ↑92.5 | ↓85.2 | 87.5 |
| 5-6 | ONLINE-B | 89.3 | 1.7 | -0.7 | 2.9 | ↑91.9 | 89.7 | 90.3 | ↓85.0 |
| 6-6 | IOL-Research | 89.0 | 1.8 | -0.7 | 3.1 | 91.0 | 90.8 | 88.3 | ↓86.1 |
| 6-7 | Claude-3.5 | 88.9 | 1.7 | -0.7 | 3.0 | 92.0 | 90.8 | 89.5 | ↓83.4 |
| 6-8 | CommandR-plus | 88.3 | 2.2 | -0.7 | 3.3 | 85.9 | ↑90.8 | 90.4 | 85.9 |
| 9-9 | Llama3-70B § | 86.5 | 2.8 | -0.7 | 3.9 | 87.5 | 86.8 | 87.0 | ↓84.6 |
| 10-10 | HW-TSC | 86.2 | 2.3 | -0.7 | 3.4 | 87.1 | ↑91.5 | 84.9 | ↓81.4 |
| 11-11 | IKUN | 85.3 | 3.1 | -0.6 | 4.0 | 88.6 | ↑89.1 | 82.1 | ↓81.5 |
| 12-12 | Aya23 | 85.2 | 3.0 | -0.7 | 4.1 | 85.4 | ↑88.3 | 85.5 | ↓81.7 |
| 13-13 | IKUN-C | 82.1 | 3.5 | -0.6 | 4.2 | 81.0 | ↑85.9 | 83.1 | ↓78.6 |
|  | ONLINE-W | - | 2.2 | - | - | - | - | - | - |
|  | Mistral-Large § | - | 2.8 | - | - | - | - | - | - |
|  | Phi-3-Medium § | - | 3.1 | - | - | - | - | - | - |
|  | ONLINE-A | - | 3.3 | - | - | - | - | - | - |
|  | UvA-MT | - | 4.3 | - | - | - | - | - | - |
|  | ONLINE-G | - | 4.8 | - | - | - | - | - | - |
|  | NVIDIA-NeMo | - | 7.3 | - | - | - | - | - | - |
|  | CycleL | - | 20.1 | - | - | - | - | - | - |
|  | CycleL2 | - | 22.0 | - | - | - | - | - | - |

### Japanese→Chinese

| Rank | System | Human | AutoRank | CometKiwi | MetricX | literary | news | speech |
|---|---|---|---|---|---|---|---|---|
| 1-3 | Claude-3.5 | -1.4 | 1.7 | -0.6 | 3.5 | -0.5 | -0.8 | ↓-3.0 |
| 1-3 | HUMAN-A | -1.5 | - | - | - | -0.7 | -0.8 | ↓-3.2 |
| 3-5 | GPT-4 | -1.7 | 2.1 | -0.6 | 3.8 | -1.0 | -0.8 | ↓-3.2 |
| 2-5 | DLUT-GTCOM | -1.7 | 2.0 | -0.6 | 3.3 | -0.5 | -1.1 | ↓-3.7 |
| 4-8 | Unbabel-Tower70B | -1.9 | 1.0 | -0.6 | 3.2 | -1.0 | -1.2 | ↓-3.5 |
| 3-6 | Gemini-1.5-Pro | -2.1 | 1.9 | -0.6 | 3.5 | -1.6 | -0.8 | ↓-3.8 |
| 6-8 | CommandR-plus | -2.2 | 2.8 | -0.6 | 4.1 | -0.7 | -1.3 | ↓-4.6 |
| 6-8 | IOL-Research | -2.4 | 2.2 | -0.6 | 3.9 | -1.4 | -1.1 | ↓-4.8 |
| 9-10 | Llama3-70B § | -3.4 | 3.1 | -0.6 | 4.7 | -2.0 | -2.2 | ↓-6.2 |
| 9-10 | Aya23 | -3.5 | 3.7 | -0.6 | 5.0 | -2.1 | -1.9 | ↓-6.4 |
| 11-12 | Team-J | -4.5 | 2.8 | -0.6 | 4.0 | -3.1 | -2.0 | ↓-8.5 |
| 11-12 | NTTSU | -5.1 | 3.7 | -0.6 | 5.3 | -2.8 | -2.1 | ↓-10.5 |
| 13-13 | ONLINE-B | -5.8 | 5.2 | -0.5 | 5.5 | -4.2 | -3.7 | ↓-9.5 |
| 14-14 | IKUN-C | -7.7 | 5.5 | -0.5 | 6.2 | -5.1 | -3.4 | ↓-14.4 |
| 15-15 | MSLC | -10.7 | 8.9 | -0.5 | 8.8 | -9.1 | ↑-4.0 | ↓-19.0 |
|  | Mistral-Large § | - | 3.5 | - | - | - | - | - |
|  | Phi-3-Medium § | - | 4.0 | - | - | - | - | - |
|  | IKUN | - | 4.4 | - | - | - | - | - |
|  | UvA-MT | - | 5.2 | - | - | - | - | - |
|  | ONLINE-W | - | 5.3 | - | - | - | - | - |
|  | ONLINE-A | - | 6.8 | - | - | - | - | - |
|  | ONLINE-G | - | 10.3 | - | - | - | - | - |
|  | CycleL | - | 23.0 | - | - | - | - | - |

# E    Head to head comparisons

Following tables show differences in average human scores for each language pair. The number in each of cell shows the difference in average human scores for the systems in the column and row.

Because there are many systems and data conditions, the significance of each pairwise comparison needs to be quantified. We apply Wilcoxon signed-rank test to measure the likelihood that such differences could occur simply by chance. In the following tables ⋆ indicates statistical significance at $p < 0.05$, † indicates statistical significance at $p < 0.01$, and ‡ indicates statistical significance at $p < 0.001$.

Each table contains final rows showing the macro-average score achieved by that system and the rank range. Gray lines separate clusters based on non-overlapping rank ranges.

**Head to head comparison for Czech→Ukrainian systems**

| | Claude-3.5 | refA | Gemini-1.5-Pro | Unbabel-Tower70B | IOL-Research | CommandR-plus | ONLINE-W | GPT-4 | IKUN | Aya23 | CUNI-Transformer | IKUN-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Claude-3.5 | – | 0.3† | 0.4‡ | 0.8 | 2.8‡ | 3.3‡ | 4.3‡ | 4.4‡ | 5.9‡ | 6.4‡ | 7.7‡ | 10.4‡ |
| refA | – | – | 0.1‡ | 0.5† | 2.5‡ | 3.0‡ | 4.0‡ | 4.2‡ | 5.6‡ | 6.1‡ | 7.4‡ | 10.1‡ |
| Gemini-1.5-Pro | – | – | – | 0.4★ | 2.4‡ | 3.0‡ | 4.0‡ | 4.1‡ | 5.5‡ | 6.1‡ | 7.3‡ | 10.1‡ |
| Unbabel-Tower70B | – | – | – | – | 2.0‡ | 2.5‡ | 3.5‡ | 3.6‡ | 5.1‡ | 5.6‡ | 6.9‡ | 9.6‡ |
| IOL-Research | – | – | – | – | – | 0.5† | 1.5‡ | 1.6‡ | 3.1‡ | 3.6‡ | 4.9‡ | 7.6‡ |
| CommandR-plus | – | – | – | – | – | – | 1.0 | 1.1★ | 2.5‡ | 3.1‡ | 4.4‡ | 7.1‡ |
| ONLINE-W | – | – | – | – | – | – | – | 0.1 | 1.6‡ | 2.1‡ | 3.4‡ | 6.1‡ |
| GPT-4 | – | – | – | – | – | – | – | – | 1.4 | 2.0‡ | 3.3‡ | 6.0‡ |
| IKUN | – | – | – | – | – | – | – | – | – | 0.6‡ | 1.8‡ | 4.5‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | – | 1.3‡ | 4.0† |
| CUNI-Transformer | – | – | – | – | – | – | – | – | – | – | – | 2.7‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 93.0 | 92.7 | 92.6 | 92.2 | 90.2 | 89.7 | 88.7 | 88.6 | 87.1 | 86.6 | 85.3 | 82.6 |
| Ranks | 1-2 | 2-2 | 3-3 | 3-4 | 5-5 | 6-7 | 6-8 | 7-9 | 8-9 | 10-10 | 11-11 | 12-12 |

**Head to head comparison for English→Czech systems**

| | refA | Unbabel-Tower70B | Claude-3.5 | ONLINE-W | CUNI-MH | Gemini-1.5-Pro | GPT-4 | CommandR-plus | IOL-Research | SCIR-MT | CUNI-DocTransformer | Aya23 | CUNI-GA | IKUN | Llama3-70B | IKUN-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| refA | – | 1.3‡ | 1.7 | 3.9‡ | 4.5‡ | 4.7‡ | 5.2‡ | 6.0‡ | 6.4‡ | 7.5‡ | 8.7‡ | 8.8‡ | 10.8‡ | 11.2‡ | 15.5‡ | 17.5‡ |
| Unbabel-Tower70B | – | – | 0.4† | 2.6‡ | 3.2‡ | 3.4‡ | 3.9‡ | 4.7‡ | 5.2‡ | 6.2‡ | 7.4‡ | 7.5‡ | 9.5‡ | 9.9‡ | 14.2‡ | 16.3‡ |
| Claude-3.5 | – | – | – | 2.2‡ | 2.8‡ | 3.0‡ | 3.5‡ | 4.3‡ | 4.7‡ | 5.8‡ | 6.9‡ | 7.0‡ | 9.1‡ | 9.5‡ | 13.8‡ | 15.8‡ |
| ONLINE-W | – | – | – | – | 0.6 | 0.8‡ | 1.3‡ | 2.1‡ | 2.6‡ | 3.6‡ | 4.8‡ | 4.9‡ | 6.9‡ | 7.3‡ | 11.6‡ | 13.7‡ |
| CUNI-MH | – | – | – | – | – | 0.2‡ | 0.7 | 1.5‡ | 2.0‡ | 3.0‡ | 4.2‡ | 4.3‡ | 6.3‡ | 6.7‡ | 11.0‡ | 13.1‡ |
| Gemini-1.5-Pro | – | – | – | – | – | – | 0.4‡ | 1.3‡ | 1.7‡ | 2.8‡ | 3.9‡ | 4.0‡ | 6.1‡ | 6.5‡ | 10.8‡ | 12.8‡ |
| GPT-4 | – | – | – | – | – | – | – | 0.9★ | 1.3‡ | 2.3‡ | 3.5‡ | 3.6‡ | 5.7‡ | 6.0‡ | 10.3‡ | 12.4‡ |
| CommandR-plus | – | – | – | – | – | – | – | – | 0.4‡ | 1.5‡ | 2.6‡ | 2.7‡ | 4.8‡ | 5.2‡ | 9.5‡ | 11.5‡ |
| IOL-Research | – | – | – | – | – | – | – | – | – | 1.0‡ | 2.2‡ | 2.3★ | 4.4‡ | 4.7‡ | 9.1‡ | 11.1‡ |
| SCIR-MT | – | – | – | – | – | – | – | – | – | – | 1.2 | 1.3‡ | 3.3‡ | 3.7‡ | 8.0‡ | 10.1‡ |
| CUNI-DocTransformer | – | – | – | – | – | – | – | – | – | – | – | 0.1‡ | 2.2‡ | 2.5‡ | 6.9‡ | 8.9‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | – | – | – | 2.1‡ | 2.4‡ | 6.8‡ | 8.8‡ |
| CUNI-GA | – | – | – | – | – | – | – | – | – | – | – | – | – | 0.4† | 4.7‡ | 6.7‡ |
| IKUN | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 4.3‡ | 6.4‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 2.0‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 92.9 | 91.6 | 91.2 | 89.0 | 88.4 | 88.2 | 87.7 | 86.9 | 86.5 | 85.4 | 84.3 | 84.2 | 82.1 | 81.7 | 77.4 | 75.4 |
| Ranks | 1-2 | 2-2 | 2-3 | 4-5 | 4-6 | 6-6 | 6-8 | 8-8 | 8-9 | 10-11 | 10-11 | 12-12 | 13-13 | 14-14 | 15-15 | 16-16 |

**Head to head comparison for English→German systems**

| | GPT-4 | Dubformer | ONLINE-B | TranssionMT | Unbabel-Tower70B | refB | Mistral-Large | CommandR-plus | ONLINE-W | Claude-3.5 | refA | IOL-Research | Gemini-1.5-Pro | Aya23 | ONLINE-A | Llama3-70B | IKUN | IKUN-C | MSLC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPT-4 | – | 0.2 | 0.3 | 0.3 | 0.3 | 0.3 | 0.5 | 0.6 | 0.7 | 0.7† | 0.8 | 0.9‡ | 1.1 | 1.6‡ | 1.8‡ | 2.6‡ | 2.7‡ | 4.4‡ | 13.8‡ |
| Dubformer | – | – | 0.1★ | 0.1★ | 0.1 | 0.2 | 0.3 | 0.5† | 0.5★ | 0.5 | 0.6 | 0.7‡ | 1.0★ | 1.4‡ | 1.7‡ | 2.4‡ | 2.5‡ | 4.2‡ | 13.7‡ |
| ONLINE-B | – | – | – | 0.0 | 0.0 | 0.1 | 0.2 | 0.4★ | 0.4★ | 0.5 | 0.6‡ | 0.6‡ | 0.9 | 1.3‡ | 1.6‡ | 2.4‡ | 2.4‡ | 4.2‡ | 13.6‡ |
| TranssionMT | – | – | – | – | 0.0† | 0.1 | 0.2 | 0.4 | 0.4‡ | 0.4 | 0.6★ | 0.6‡ | 0.9‡ | 1.3‡ | 1.6‡ | 2.3‡ | 2.4‡ | 4.2‡ | 13.6‡ |
| Unbabel-Tower70B | – | – | – | – | – | 0.1 | 0.2 | 0.2 | 0.4 | 0.4‡ | 0.4 | 0.6‡ | 0.9 | 1.3‡ | 1.6‡ | 2.3‡ | 2.3‡ | 4.2‡ | 13.6‡ |
| refB | – | – | – | – | – | – | 0.1★ | 0.2 | 0.3‡ | 0.4 | 0.5 | 0.5 | 0.8 | 1.2‡ | 1.5‡ | 2.3‡ | 2.3‡ | 4.1‡ | 13.5‡ |
| Mistral-Large | – | – | – | – | – | – | – | 0.2 | 0.2 | 0.2 | 0.4 | 0.4★ | 0.7 | 1.1‡ | 1.4‡ | 2.1‡ | 2.1‡ | 3.9‡ | 13.4‡ |
| CommandR-plus | – | – | – | – | – | – | – | – | 0.0★ | 0.1 | 0.2 | 0.3★ | 0.5 | 0.9‡ | 1.2‡ | 2.0‡ | 2.0‡ | 3.8‡ | 13.2‡ |
| ONLINE-W | – | – | – | – | – | – | – | – | – | 0.0† | 0.2 | 0.2 | 0.5‡ | 0.9‡ | 1.2‡ | 1.9‡ | 2.0‡ | 3.7‡ | 13.1‡ |
| Claude-3.5 | – | – | – | – | – | – | – | – | – | – | 0.1 | 0.2‡ | 0.4 | 0.9‡ | 1.1‡ | 1.9‡ | 2.0‡ | 3.7‡ | 13.1‡ |
| refA | – | – | – | – | – | – | – | – | – | – | – | 0.1 | 0.3 | 0.8‡ | 1.0‡ | 1.8‡ | 1.9‡ | 3.6‡ | 13.0‡ |
| IOL-Research | – | – | – | – | – | – | – | – | – | – | – | – | 0.2‡ | 0.7★ | 0.9★ | 1.7‡ | 1.8‡ | 3.5‡ | 12.9‡ |
| Gemini-1.5-Pro | – | – | – | – | – | – | – | – | – | – | – | – | – | 0.4‡ | 0.7‡ | 1.5‡ | 1.5‡ | 3.3‡ | 12.7‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 0.3 | 1.0★ | 1.1† | 2.8‡ | 12.3‡ |
| ONLINE-A | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 0.8 | 0.8 | 2.6‡ | 12.0‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 0.1 | 1.8‡ | 11.2‡ |
| IKUN | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 1.7‡ | 11.1‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 9.4‡ |
| MSLC | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | -1.6 | -1.8 | -1.9 | -1.9 | -1.9 | -2.0 | -2.1 | -2.3 | -2.3 | -2.4 | -2.5 | -2.5 | -2.8 | -3.2 | -3.5 | -4.3 | -4.3 | -6.1 | -15.5 |
| Ranks | 1-11 | 1-7 | 2-10 | 2-10 | 2-9 | 1-9 | 2-12 | 4-11 | 8-10 | 2-12 | 3-13 | 10-12 | 5-13 | 14-15 | 14-17 | 15-17 | 15-17 | 18-18 | 19-19 |

**Head to head comparison for English→Spanish systems**

| | refA | Dubformer | GPT-4 | IOL-Research | Mistral-Large | Unbabel-Tower70B | Claude-3.5 | Gemini-1.5-Pro | CommandR-plus | Llama3-70B | ONLINE-B | IKUN | IKUN-C | MSLC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| refA | – | 1.9† | 3.3‡ | 3.8‡ | 6.0‡ | 6.3‡ | 6.5‡ | 6.5‡ | 7.0‡ | 8.1‡ | 9.7‡ | 10.5‡ | 14.9‡ | 31.3‡ |
| Dubformer | – | – | 1.5★ | 2.0† | 4.1‡ | 4.5‡ | 4.6‡ | 4.6‡ | 5.1‡ | 6.2‡ | 7.8‡ | 8.7‡ | 13.0‡ | 29.5‡ |
| GPT-4 | – | – | – | 0.5‡ | 2.7‡ | 3.0† | 3.1 | 3.2★ | 3.6‡ | 4.8‡ | 6.4‡ | 7.2‡ | 11.6‡ | 28.0‡ |
| IOL-Research | – | – | – | – | 2.2† | 2.5 | 2.7 | 2.7 | 4.3‡ | 5.9‡ | 6.7‡ | | 11.1‡ | 27.5‡ |
| Mistral-Large | – | – | – | – | – | 0.4★ | 0.5 | 0.5★ | 1.0 | 2.1 | 3.7† | 4.5‡ | 8.9‡ | 25.4‡ |
| Unbabel-Tower70B | – | – | – | – | – | – | 0.1 | 0.1 | 0.6 | 1.7† | 3.3‡ | 4.2‡ | 8.5‡ | 25.0‡ |
| Claude-3.5 | – | – | – | – | – | – | – | 0.0 | 0.5★ | 1.6† | 3.2‡ | 4.0‡ | 8.4‡ | 24.9‡ |
| Gemini-1.5-Pro | – | – | – | – | – | – | – | – | 0.5★ | 1.6 | 3.2‡ | 4.0‡ | 8.4‡ | 24.9‡ |
| CommandR-plus | – | – | – | – | – | – | – | – | – | 1.1‡ | 2.7† | 3.5‡ | 7.9‡ | 24.4‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | – | 1.6★ | 2.4‡ | 6.8‡ | 23.2‡ |
| ONLINE-B | – | – | – | – | – | – | – | – | – | – | – | 0.8★ | 5.2‡ | 21.7‡ |
| IKUN | – | – | – | – | – | – | – | – | – | – | – | – | 4.4 | 20.8‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – | – | 16.4‡ |
| MSLC | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 95.3 | 93.4 | 91.9 | 91.4 | 89.3 | 88.9 | 88.8 | 88.8 | 88.3 | 87.2 | 85.6 | 84.7 | 80.4 | 63.9 |
| Ranks | 1-1 | 2-2 | 3-4 | 4-7 | 5-8 | 5-9 | 3-8 | 5-8 | 7-9 | 9-10 | 11-11 | 12-13 | 12-13 | 14-14 |

**Head to head comparison for English→Hindi systems**

| | TranssionMT | Unbabel-Tower70B | Claude-3.5 | ONLINE-B | Gemini-1.5-Pro | GPT-4 | refA | IOL-Research | Llama3-70B | Aya23 | IKUN-C |
|---|---|---|---|---|---|---|---|---|---|---|---|
| TranssionMT | – | 0.7 | 1.0‡ | 1.1‡ | 1.3 | 2.7‡ | 2.8‡ | 4.1‡ | 4.6‡ | 6.5‡ | 20.5‡ |
| Unbabel-Tower70B | – | – | 0.3‡ | 0.4 | 0.5 | 2.0† | 2.0‡ | 3.4‡ | 3.8‡ | 5.8‡ | 19.8‡ |
| Claude-3.5 | – | – | – | 0.1‡ | 0.3‡ | 1.7‡ | 1.8‡ | 3.1‡ | 3.6‡ | 5.5‡ | 19.5‡ |
| ONLINE-B | – | – | – | – | 0.1† | 1.6‡ | 1.6‡ | 3.0‡ | 3.4‡ | 5.4‡ | 19.4‡ |
| Gemini-1.5-Pro | – | – | – | – | – | 1.5★ | 1.5‡ | 2.8‡ | 3.3‡ | 5.2‡ | 19.3‡ |
| GPT-4 | – | – | – | – | – | – | 0.0‡ | 1.3† | 1.8‡ | 3.8‡ | 17.8‡ |
| refA | – | – | – | – | – | – | – | 1.3★ | 1.8 | 3.7‡ | 17.7‡ |
| IOL-Research | – | – | – | – | – | – | – | – | 0.5† | 2.4‡ | 16.4‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | 1.9† | 16.0‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | – | 14.0‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 91.3 | 90.5 | 90.2 | 90.1 | 90.0 | 88.5 | 88.5 | 87.2 | 86.7 | 84.7 | 70.7 |
| Ranks | 1-3 | 1-4 | 3-3 | 3-4 | 3-5 | 6-6 | 7-8 | 8-8 | 8-9 | 10-10 | 11-11 |

**Head to head comparison for English→Icelandic systems**

| | refA | Dubformer | Claude-3.5 | Unbabel-Tower70B | AMI | IKUN | ONLINE-B | GPT-4 | IKUN-C | IOL-Research | Llama3-70B |
|---|---|---|---|---|---|---|---|---|---|---|---|
| refA | – | 8.8‡ | 11.1‡ | 12.9‡ | 19.8‡ | 22.0‡ | 25.1‡ | 26.7‡ | 27.9‡ | 35.1‡ | 52.0‡ |
| Dubformer | – | – | 2.3 | 4.1‡ | 11.0‡ | 13.2‡ | 16.3‡ | 17.9‡ | 19.1‡ | 26.3‡ | 43.3‡ |
| Claude-3.5 | – | – | – | 1.8‡ | 8.7‡ | 10.9‡ | 14.0‡ | 15.6‡ | 16.7‡ | 24.0‡ | 40.9‡ |
| Unbabel-Tower70B | – | – | – | – | 6.9‡ | 9.1‡ | 12.2‡ | 13.8‡ | 15.0‡ | 22.2‡ | 39.1‡ |
| AMI | – | – | – | – | – | 2.2★ | 5.3‡ | 6.9‡ | 8.1‡ | 15.3‡ | 32.3‡ |
| IKUN | – | – | – | – | – | – | 3.1† | 4.7‡ | 5.9‡ | 13.1‡ | 30.0‡ |
| ONLINE-B | – | – | – | – | – | – | – | 1.6‡ | 2.8‡ | 10.0‡ | 26.9‡ |
| GPT-4 | – | – | – | – | – | – | – | – | 1.2 | 8.4‡ | 25.3‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | 7.2‡ | 24.2‡ |
| IOL-Research | – | – | – | – | – | – | – | – | – | – | 16.9‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 93.1 | 84.3 | 81.9 | 80.2 | 73.3 | 71.0 | 68.0 | 66.3 | 65.2 | 58.0 | 41.0 |
| Ranks | 1-1 | 2-3 | 2-3 | 4-4 | 5-5 | 6-6 | 7-7 | 8-9 | 8-9 | 10-10 | 11-11 |

**Head to head comparison for English→Japanese systems**

| | refA | ONLINE-B | CommandR-plus | GPT-4 | Claude-3.5 | Gemini-1.5-Pro | Unbabel-Tower70B | IOL-Research | Aya23 | NTTSU | Team-J | Llama3-70B | IKUN-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| refA | – | 0.7‡ | 0.9★ | 1.0‡ | 1.0‡ | 1.8† | 2.1‡ | 2.1‡ | 2.2‡ | 2.4‡ | 3.3‡ | 5.1‡ | 10.1‡ |
| ONLINE-B | – | – | 0.1† | 0.3‡ | 0.3 | 1.1 | 1.4‡ | 1.4‡ | 1.4★ | 1.7‡ | 2.6‡ | 4.3‡ | 9.3‡ |
| CommandR-plus | – | – | – | 0.1‡ | 0.2† | 0.9 | 1.2‡ | 1.3‡ | 1.3‡ | 1.5‡ | 2.5‡ | 4.2‡ | 9.2‡ |
| GPT-4 | – | – | – | – | 0.0‡ | 0.8† | 1.1‡ | 1.1‡ | 1.2‡ | 1.4‡ | 2.3‡ | 4.0‡ | 9.1‡ |
| Claude-3.5 | – | – | – | – | – | 0.8★ | 1.1‡ | 1.1‡ | 1.1‡ | 1.4‡ | 2.3‡ | 4.0‡ | 9.1‡ |
| Gemini-1.5-Pro | – | – | – | – | – | – | 0.3‡ | 0.3‡ | 0.4 | 0.6‡ | 1.5‡ | 3.2‡ | 8.3‡ |
| Unbabel-Tower70B | – | – | – | – | – | – | – | 0.0‡ | 0.1‡ | 0.3‡ | 1.2‡ | 2.9‡ | 8.0‡ |
| IOL-Research | – | – | – | – | – | – | – | – | 0.1† | 0.3‡ | 1.2‡ | 2.9‡ | 8.0‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | 0.2★ | 1.2‡ | 2.9‡ | 7.9‡ |
| NTTSU | – | – | – | – | – | – | – | – | – | – | 0.9‡ | 2.7‡ | 7.7‡ |
| Team-J | – | – | – | – | – | – | – | – | – | – | – | 1.7‡ | 6.8‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | – | – | – | 5.0‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 91.8 | 91.1 | 91.0 | 90.8 | 90.8 | 90.0 | 89.7 | 89.7 | 89.7 | 89.4 | 88.5 | 86.8 | 81.7 |
| Ranks | 1-1 | 2-4 | 3-4 | 4-4 | 4-5 | 4-7 | 7-7 | 8-8 | 8-9 | 10-10 | 11-11 | 12-12 | 13-13 |

## Head to head comparison for English→Russian systems

| | refA | Dubformer | Claude-3.5 | Unbabel-Tower70B | Yandex | Gemini-1.5-Pro | GPT-4 | ONLINE-G | CommandR-plus | IOL-Research | IKUN | Aya23 | Llama3-70B | IKUN-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| refA | – | 0.1‡ | 0.9★ | 1.0‡ | 2.2† | 3.7★ | 4.1‡ | 4.6† | 4.8‡ | 7.1‡ | 10.0‡ | 10.5‡ | 13.4‡ | 19.4‡ |
| Dubformer | – | – | 0.8‡ | 0.9 | 2.1† | 3.6‡ | 4.0‡ | 4.5‡ | 4.7‡ | 6.9‡ | 9.8‡ | 10.4‡ | 13.3‡ | 19.3‡ |
| Claude-3.5 | – | – | – | 0.1‡ | 1.3 | 2.8† | 3.2‡ | 3.7† | 3.9† | 6.1‡ | 9.0‡ | 9.6‡ | 12.5‡ | 18.5‡ |
| Unbabel-Tower70B | – | – | – | – | 1.1 | 2.7‡ | 3.1‡ | 3.6‡ | 3.8‡ | 6.0‡ | 8.9‡ | 9.5‡ | 12.4‡ | 18.4‡ |
| Yandex | – | – | – | – | – | 1.5★ | 1.9★ | 2.4 | 2.6 | 4.9‡ | 7.8‡ | 8.3‡ | 11.2‡ | 17.2‡ |
| Gemini-1.5-Pro | – | – | – | – | – | – | 0.4 | 0.9† | 1.1 | 3.3‡ | 6.2‡ | 6.8‡ | 9.7‡ | 15.7‡ |
| GPT-4 | – | – | – | – | – | – | – | 0.5 | 0.7 | 2.9† | 5.8‡ | 6.4‡ | 9.3‡ | 15.3‡ |
| ONLINE-G | – | – | – | – | – | – | – | – | 0.2 | 2.5★ | 5.3‡ | 5.9‡ | 8.8‡ | 14.8‡ |
| CommandR-plus | – | – | – | – | – | – | – | – | – | 2.2† | 5.1‡ | 5.7‡ | 8.6‡ | 14.6‡ |
| IOL-Research | – | – | – | – | – | – | – | – | – | – | 2.9† | 3.5‡ | 6.4‡ | 12.4‡ |
| IKUN | – | – | – | – | – | – | – | – | – | – | – | 0.6★ | 3.5★ | 9.5‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | – | – | – | 2.9★ | 8.9‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | – | – | – | – | 6.0‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 89.2 | 89.1 | 88.2 | 88.1 | 87.0 | 85.5 | 85.0 | 84.6 | 84.3 | 82.1 | 79.2 | 78.6 | 75.7 | 69.8 |
| Ranks | 1-1 | 2-3 | 3-4 | 3-5 | 3-7 | 6-8 | 6-9 | 6-9 | 5-9 | 10-10 | 11-11 | 12-12 | 13-13 | 14-14 |

## Head to head comparison for English→Ukrainian systems

| | Claude-3.5 | Unbabel-Tower70B | Dubformer | refA | Gemini-1.5-Pro | ONLINE-W | GPT-4 | CommandR-plus | IOL-Research | IKUN | IKUN-C |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Claude-3.5 | – | 0.6 | 1.5† | 3.2‡ | 3.3† | 4.4‡ | 5.9‡ | 7.3‡ | 7.3‡ | 12.1‡ | 22.6‡ |
| Unbabel-Tower70B | – | – | 0.8‡ | 2.6† | 2.7★ | 3.8† | 5.2‡ | 6.7‡ | 6.7‡ | 11.4‡ | 21.9‡ |
| Dubformer | – | – | – | 1.7‡ | 1.8‡ | 2.9† | 4.4‡ | 5.8‡ | 5.8‡ | 10.6‡ | 21.1‡ |
| refA | – | – | – | – | 0.1 | 1.2★ | 2.7 | 4.1★ | 4.1‡ | 8.8‡ | 19.4‡ |
| Gemini-1.5-Pro | – | – | – | – | – | 1.1 | 2.5★ | 4.0★ | 4.0‡ | 8.7‡ | 19.2‡ |
| ONLINE-W | – | – | – | – | – | – | 1.4 | 2.9 | 2.9★ | 7.6‡ | 18.1‡ |
| GPT-4 | – | – | – | – | – | – | – | 1.4 | 1.4 | 6.2‡ | 16.7‡ |
| CommandR-plus | – | – | – | – | – | – | – | – | 0.0 | 4.8‡ | 15.3‡ |
| IOL-Research | – | – | – | – | – | – | – | – | – | 4.7‡ | 15.3‡ |
| IKUN | – | – | – | – | – | – | – | – | – | – | 10.5‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 90.5 | 89.8 | 89.0 | 87.3 | 87.1 | 86.0 | 84.6 | 83.2 | 83.2 | 78.4 | 67.9 |
| Ranks | 1-2 | 1-2 | 3-3 | 4-6 | 4-6 | 5-8 | 5-9 | 6-9 | 7-9 | 10-10 | 11-11 |

## Head to head comparison for English→Chinese systems

| | GPT-4 | Unbabel-Tower70B | refA | Gemini-1.5-Pro | ONLINE-B | IOL-Research | Claude-3.5 | CommandR-plus | Llama3-70B | HW-TSC | IKUN | Aya23 | IKUN-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPT-4 | – | 0.0‡ | 0.3★ | 0.3‡ | 0.4‡ | 0.6‡ | 0.7† | 1.4★ | 3.2‡ | 3.4‡ | 4.3‡ | 4.4‡ | 7.5‡ |
| Unbabel-Tower70B | – | – | 0.2 | 0.3‡ | 0.3‡ | 0.6‡ | 0.7‡ | 1.3 | 3.1‡ | 3.4‡ | 4.3‡ | 4.4‡ | 7.5‡ |
| refA | – | – | – | 0.1‡ | 0.1‡ | 0.3‡ | 0.5★ | 1.1 | 2.9‡ | 3.1‡ | 4.1‡ | 4.2‡ | 7.2‡ |
| Gemini-1.5-Pro | – | – | – | – | 0.1‡ | 0.3† | 0.4‡ | 1.1‡ | 2.9‡ | 3.1‡ | 4.0‡ | 4.1‡ | 7.2‡ |
| ONLINE-B | – | – | – | – | – | 0.2★ | 0.3 | 1.0† | 2.8‡ | 3.0‡ | 3.9‡ | 4.0‡ | 7.1‡ |
| IOL-Research | – | – | – | – | – | – | 0.1★ | 0.8‡ | 2.6‡ | 2.8‡ | 3.7‡ | 3.8‡ | 6.9‡ |
| Claude-3.5 | – | – | – | – | – | – | – | 0.6† | 2.5‡ | 2.7‡ | 3.6‡ | 3.7‡ | 6.8‡ |
| CommandR-plus | – | – | – | – | – | – | – | – | 1.8‡ | 2.0‡ | 2.9‡ | 3.0‡ | 6.1‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | 0.2‡ | 1.1‡ | 1.2‡ | 4.3‡ |
| HW-TSC | – | – | – | – | – | – | – | – | – | – | 0.9‡ | 1.0‡ | 4.1‡ |
| IKUN | – | – | – | – | – | – | – | – | – | – | – | 0.1‡ | 3.2‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | – | – | – | 3.1‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | 89.6 | 89.6 | 89.4 | 89.3 | 89.3 | 89.0 | 88.9 | 88.3 | 86.5 | 86.2 | 85.3 | 85.2 | 82.1 |
| Ranks | 1-1 | 2-4 | 2-4 | 4-4 | 5-6 | 6-6 | 6-7 | 6-8 | 9-9 | 10-10 | 11-11 | 12-12 | 13-13 |

## Head to head comparison for Japanese→Chinese systems

| | Claude-3.5 | refA | GPT-4 | DLUT-GTCOM | Unbabel-Tower70B | Gemini-1.5-Pro | CommandR-plus | IOL-Research | Llama3-70B | Aya23 | Team-J | NTTSU | ONLINE-B | IKUN-C | MSLC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Claude-3.5 | – | 0.1 | 0.3★ | 0.3 | 0.5‡ | 0.7★ | 0.8‡ | 1.0‡ | 2.0‡ | 2.0‡ | 3.1‡ | 3.7‡ | 4.4‡ | 6.2‡ | 9.3‡ |
| refA | – | – | 0.1★ | 0.2 | 0.4‡ | 0.5★ | 0.7‡ | 0.9‡ | 1.9‡ | 1.9‡ | 3.0‡ | 3.6‡ | 4.3‡ | 6.1‡ | 9.2‡ |
| GPT-4 | – | – | – | 0.1★ | 0.2 | 0.4 | 0.5‡ | 0.8† | 1.7‡ | 1.8‡ | 2.9‡ | 3.4‡ | 4.1‡ | 6.0‡ | 9.0‡ |
| DLUT-GTCOM | – | – | – | – | 0.2★ | 0.3 | 0.5‡ | 0.7‡ | 1.7‡ | 1.7‡ | 2.8‡ | 3.4‡ | 4.1‡ | 5.9‡ | 9.0‡ |
| Unbabel-Tower70B | – | – | – | – | – | 0.2 | 0.3 | 0.5 | 1.5‡ | 1.6‡ | 2.6‡ | 3.2‡ | 3.9‡ | 5.7‡ | 8.8‡ |
| Gemini-1.5-Pro | – | – | – | – | – | – | 0.1‡ | 0.4† | 1.4‡ | 1.4‡ | 2.5‡ | 3.1‡ | 3.8‡ | 5.6‡ | 8.7‡ |
| CommandR-plus | – | – | – | – | – | – | – | 0.3 | 1.2‡ | 1.3‡ | 2.3‡ | 2.9‡ | 3.6‡ | 5.5‡ | 8.5‡ |
| IOL-Research | – | – | – | – | – | – | – | – | 1.0‡ | 1.0‡ | 2.1‡ | 2.7‡ | 3.4‡ | 5.2‡ | 8.3‡ |
| Llama3-70B | – | – | – | – | – | – | – | – | – | 0.0 | 1.1† | 1.7‡ | 2.4‡ | 4.2‡ | 7.3‡ |
| Aya23 | – | – | – | – | – | – | – | – | – | – | 1.1‡ | 1.7★ | 2.4‡ | 4.2‡ | 7.3‡ |
| Team-J | – | – | – | – | – | – | – | – | – | – | – | 0.6 | 1.3‡ | 3.1‡ | 6.2‡ |
| NTTSU | – | – | – | – | – | – | – | – | – | – | – | – | 0.7‡ | 2.5‡ | 5.6‡ |
| ONLINE-B | – | – | – | – | – | – | – | – | – | – | – | – | – | 1.8† | 4.9‡ |
| IKUN-C | – | – | – | – | – | – | – | – | – | – | – | – | – | – | 3.1‡ |
| MSLC | – | – | – | – | – | – | – | – | – | – | – | – | – | – | – |
| Scores | -1.4 | -1.5 | -1.7 | -1.7 | -1.9 | -2.1 | -2.2 | -2.4 | -3.4 | -3.5 | -4.5 | -5.1 | -5.8 | -7.7 | -10.7 |
| Ranks | 1-3 | 1-3 | 3-5 | 2-5 | 4-8 | 3-6 | 6-8 | 6-8 | 9-10 | 9-10 | 11-12 | 11-12 | 13-13 | 14-14 | 15-15 |

# Are LLMs Breaking MT Metrics?
# Results of the WMT24 Metrics Shared Task

**Markus Freitag**[(1)], **Nitika Mathur**[(2)], **Daniel Deutsch**[(1)], **Chi-kiu Lo** 羅致翹[(3)],
**Eleftherios Avramidis**[(4)], **Ricardo Rei**[(5)], **Brian Thompson**[(6)], **Frédéric Blain**[(7)], **Tom Kocmi**[(8)],
**Jiayi Wang**[(9)], **David I. Adelani**[(10,11)], **Marianna Buchicchio**[(5)], **Chrysoula Zerva**[(12,13)], **Alon Lavie**[(14)]

[(1)]Google Research [(2)]Oracle [(3)]National Research Council Canada
[(4)]German Research Center for Artificial Intelligence (DFKI) [(5)]Unbabel [(6)]Amazon [(7)]Tilburg University
[(8)]Microsoft [(9)]University College London [(10)]McGill University [(11)]Mila - Quebec AI Institute
[(12)]Instituto Superior Técnico [(13)]Instituto de Telecomunicações [(14)]Phrase
`wmt-metrics@googlegroups.com`

## Abstract

The WMT24 Metrics Shared Task evaluated the performance of automatic metrics for machine translation (MT), with a major focus on LLM-based translations that were generated as part of the WMT24 General MT Shared Task. As LLMs become increasingly popular in MT, it is crucial to determine whether existing evaluation metrics can accurately assess the output of these systems.

To provide a robust benchmark for this evaluation, human assessments were collected using Multidimensional Quality Metrics (MQM), continuing the practice from recent years. Furthermore, building on the success of the previous year, a challenge set subtask was included, requiring participants to design contrastive test suites that specifically target a metric's ability to identify and penalize different types of translation errors.

Finally, the meta-evaluation procedure was refined to better reflect real-world usage of MT metrics, focusing on pairwise accuracy at both the system- and segment-levels.

We present an extensive analysis on how well metrics perform on three language pairs: English→Spanish (Latin America), Japanese→Chinese, and English→German. The results strongly confirm the results reported last year, that fine-tuned neural metrics continue to perform well, even when used to evaluate LLM-based translation systems.

| metric | | avg corr |
|---|---|---|
| MetaMetrics-MT | **1** | **0.725** |
| MetricX-24-Hybrid | **1** | **0.721** |
| XCOMET | **1** | **0.719** |
| MetricX-24-Hybrid-QE* | 2 | 0.714 |
| gemba_esa* | 2 | 0.711 |
| XCOMET-QE* | 3 | 0.695 |
| COMET-22 | 3 | 0.688 |
| BLEURT-20 | 3 | 0.686 |
| MetaMetrics-MT-QE* | 3 | 0.684 |
| bright-qe* | 4 | 0.681 |
| BLCOM_1 | 4 | 0.664 |
| sentinel-cand-mqm* | 5 | 0.650 |
| PrismRefMedium | 5 | 0.646 |
| PrismRefSmall | 5 | 0.642 |
| CometKiwi* | 5 | 0.640 |
| damonmonli | 5 | 0.635 |
| YiSi-1 | 6 | 0.630 |
| BERTScore | 7 | 0.617 |
| MEE4 | 7 | 0.609 |
| chrF | 8 | 0.608 |
| chrfS | 8 | 0.606 |
| spBLEU | 9 | 0.593 |
| BLEU | 9 | 0.589 |
| XLsimMqm* | 10 | 0.515 |
| sentinel-src-mqm* | 10 | 0.513 |
| sentinel-ref-mqm | 10 | 0.513 |

Table 1: Official ranking of primary submissions to the WMT24 Metric Task. The final score is the weighted average correlation over 6 different tasks. Starred metrics are reference-free, and underlined metrics are baselines. See Table 14 for the pairwise comparisons from which the ranks were derived.

## 1 Introduction

The Metrics Shared Task[1] has been a key component of WMT since 2008, serving as a way to validate the use of automatic MT evaluation metrics and drive the development of new metrics. We evaluate reference-based automatic metrics that score MT output by comparing the translations with a

reference translation generated by human translators, who are instructed to translate "from scratch" without post-editing from MT. In addition, we also invited submissions of reference-free metrics (quality estimation metrics or QE metrics) that compare MT outputs directly with the source segments. All metrics are evaluated based on their agreement with human ratings when scoring MT systems and human translations at the system and sentence level.

---

[1]`https://www2.statmt.org/wmt24/metrics-task.html`

The final ranking of this year's submitted primary metrics is shown in Table 1. Below are some of the key details and changes implemented for this year's Metrics Shared Task:

- **Language Pairs**: For this year, we focus on three language pairs, all on the paragraph-level: (i) English→German (en→de), English→Spanish (Latin America) (en→es), and Japanese→Chinese (ja→zh).

- **Human Evaluation**: Like last year, we collected our own human quality ratings for our three language pairs leveraging professional translators performing MQM annotations (Lommel et al., 2014; Freitag et al., 2021). We released and uploaded[2] all MQM annotations, and we recommend using Marot[3] for looking into this data.

- **Meta Evaluation**: This year, we designed the meta-evaluation to evaluate metrics on how they are used in practice, by focusing on pairwise accuracy at the system- and segment-levels and removing Pearson correlation. At the system-level, we use a new statistic called soft pairwise accuracy (Thompson et al., 2024), and, like last year, we use pairwise accuracy with tie calibration (Deutsch et al., 2023) at the segment-level.

- **Challenge Sets Subtask**: The submission format of the challenge sets changed to provide for more flexibility on how the participants could challenge the metrics. In contrast to previous years, when the challenge items were evaluated in a rigid pairwise manner on whether the metric scores can distinguish between a good and a bad translation, this year's participants could provide single translations and then employ an evaluation concept of their own. This year's subtask features 4 submissions that test the ability of the metrics to evaluate MT outputs on African languages, the biomedical domain, on more than a hundred linguistically-motivated phenomena, as well as on low- to mid-quality outputs and specific challenges (empty strings, wrong/mixed language output and language variants).

- **Understand Magnitude of Score Difference**: Similar to last year, we include two analyses to understand the meaning of the score differences

that metrics present with respect to the statistical significance of MT system rankings according to human annotations and metric scores. These analyses provide additional assistance for MT researchers to build an intuition on the relationship between the magnitude of metric score differences and the reliability of the improved translation quality.

- **MTME**: Similar to last year, all the data has been uploaded to MTME[4], and all results in this paper are calculated with this analysis tool. We encourage every metric developer to use MTME to calculate contrastive scores to enhance consistency and comparability going forward.

Our main findings are:

- Two metametrics (which are both ensemble metrics), MetricX-24-Hybrid and XCOMET, are the winners of the WMT24 Metrics Shared Task (Table 1);

- Fine-tuned neural metrics continue to be strong in performance and are effective quality estimators, even for LLM-based translations;

- Results from the challenge sets independently suggest that it is important for metric researchers to test the performance of metrics in diverse collections of linguistic phenomena, languages and domains, including low-resource languages, mixed languages and irregular outputs, and on a wide range of translation quality, in order to minimize anomalous and unexpected behaviours of metrics (Section 9).

The rest of the paper is organized as follows: Section 2 describes the test data. Section 3 presents an overview of the conducted expert-based human evaluation. Section 4 describes the metrics evaluated this year (baselines and participants). Section 5 describes the conducted meta-evaluation. Section 6 reports our main results. Section 7 interprets and evaluates metrics' scores beyond correlations. Section 8 summarizes our results for the WMT24 General MT Shared Task language-pairs based on their new ESA human evaluation methodology (Kocmi et al., 2024c). Section 9 presents a description of the submitted challenge sets along with their findings. Finally, Section 10 summarizes our most important conclusions.

---

[2]https://github.com/google/wmt-mqm-human-evaluation
[3]https://github.com/google-research/google-research/tree/master/marot

[4]https://github.com/google-research/mt-metrics-eval

48

## 2 Translation Systems

Similar to previous years' editions, the source, reference texts, and MT system outputs for the metrics task are mainly derived from the WMT24 General MT Shared Task (Kocmi et al., 2024a). The domains cover news, literary, speech, and social. We do not provide any sentence splitting, thus many segments contain multiple sentences. Each language pair contains a comparable number of sentences from each domain, resulting in reasonably balanced test sets. Data statistics can be seen in Table 2. The language pairs en→de and en→es have the same source segments; ja→zh consists of segments from only 3 different domains.

| | news | literary | speech | social |
|---|---|---|---|---|
| | #tokens | | | |
| en→{de,es} | 9,268 | 9,601 | 9,611 | 9,829 |
| ja→zh | 14,896 | 14,541 | 11,025 | |
| | #docs (#segments/doc) | | | |
| en→{de,es} | 17 (8.8) | 8 (25.8) | 111 (1.0) | 34 (15.6) |
| ja→zh | 45 (6.0) | 15 (21.1) | 136 (1.0) | |
| | #sents (#sents/doc) | | | |
| en→{de,es} | 333 (19.6) | 607 (75.9) | 685 (6.2) | 759 (22.3) |
| ja→zh | 634 (14.1) | 875 (58.3) | 332 (2.4) | |

Table 2: Test set statistics split by domain. Statistics are calculated on the source side.

The reference translations provided for the test sets are produced by professional translators.

For more details regarding the test sets, we refer the reader to the WMT24 General MT Shared Task findings paper (Kocmi et al., 2024a). All data has been released and can be downloaded[5].

## 3 MQM Human Evaluation

Automatic metrics are commonly evaluated by measuring correlations with corresponding human ratings. The quality of these human ratings is critical, and recent findings (Freitag et al., 2021) have shown that crowdsourced human ratings are not sufficiently reliable for evaluating high quality MT outputs. Furthermore, an evaluation schema based on MQM (Lommel et al., 2014), which requires explicit error annotation is more effective than an evaluation schema that only asks raters for a single scalar value per translation. Similar to last year, we decided to conduct our own MQM-based

human evaluation on a subset of translation system submissions and language pairs which we believe are most interesting for evaluating current metrics. Instead of evaluating all MT system submissions, we restrict our human evaluation to the top scoring submissions, as determined based on baseline automatic scores. MQM is a general framework that provides a hierarchy of translation errors which can be tailored to specific applications. Google and Unbabel sponsored the human evaluation for this year's metrics task for a subset of language pairs using either professional translators (English→German, Japanese→Chinese) or trusted and trained raters (English→Spanish). The error annotation typology and guidelines used by Google's and Unbabel's annotators differ slightly and are described in the following two sections.

### 3.1 English→German & Japanese→Chinese

Annotations for en→de and ja→zh were sponsored and executed by Google, using 18 professional translators (10 for en→de, 8 for ja→zh) having access to the full document context. Each segment gets annotated by a single rater. Instead of assigning a scalar value to each translation, annotators were instructed to label error spans within each segment in a document, paying particular attention to document context. Each error was highlighted in the text, and labelled with an error category and a severity. Segments that are too badly garbled to permit reliable identification of individual errors are assigned a special *Non-translation* error. Error severities are assigned independent of category, and consist of *Major*, *Minor*, and *Neutral* levels, corresponding respectively to actual translation or grammatical errors, smaller imperfections and purely subjective opinions about the translation. Since we are ultimately interested in scoring segments, we adopt the weighting scheme shown in Table 3.

| Severity | Category | Weight |
|---|---|---|
| Major | Non-translation<br>all others | 25<br>5 |
| Minor | Fluency/Punctuation<br>all others | 0.1<br>1 |
| Neutral | all | 0 |

Table 3: Google's MQM error weighting.

Recent research demonstrated that rater assignment is crucial for reliable human evaluation and we adopted the suggested Pseudo-Side-by-Side

(pSxS) rater assignment as suggested in (Riley et al., 2024). For more details, exact annotator instructions and a list of error categories, we refer the reader to Freitag et al. (2021) as the exact same setup was used for the previous three metrics tasks.

## 3.2 English→Spanish (Latin America)

The annotations for the en→es (Latin America)[6] language pair were sourced from Unbabel, who engaged four professional native language annotators possessing extensive translation experience. Much like Google's approach, these annotators were provided with the full document context, comprising up to ten segments. Their task was to identify and classify errors by highlighting them, following Unbabel's MQM 3.0 typology[7].

The annotators were instructed to classify the errors based on severity, with Unbabel's classification encompassing not only "Minor" and "Major" error severities (analogous to Google's criteria) but also a "Critical" error severity. However, to ensure consistency in our evaluation process, we opted to align with the Google methodology outlined previously. Specifically, we treated all annotated "Critical" errors as "Major" errors, and we applied a weighting scheme for punctuation errors, as detailed in Table 3.

## 3.3 Human Evaluation Results

Due to the fact that we ran our own human evaluation, we were only able to evaluate a subset of the test segments. In Table 4, you can see the number of segments and documents for each language pair and test set that we used for human evaluation. In all cases, the MQM score for a segment is the sum of the scores for the errors in that segment, and the MQM score for a test set is the average of the MQM scores of the segments that were annotated.

The results of the MQM human evaluation can be seen in Table 5. It's important to note a non-intentional, but important difference in our human evaluation setting for the speech domain between the three language pairs. For English→German and English→Spanish, we asked human annotators to compare translations against the ASR output, which inadvertently disadvantaged participants who used audio input, including those providing human translations, as these translations rely on an

error-free input. This is evident in the higher MQM scores for the speech domain for both language pairs for human translations and the dubformer system (which also utilizes audio input). However, for Japanese→Chinese, the human annotators compared against the cleaned human transcription. This mismatch was not intentional and we will discuss the impact on the correlation numbers in Section 6.

## 4 Baselines and Submissions

We computed scores for several baseline metrics in order to compare submissions against previous well-studied metrics. We will start by describing those baselines, and then we will describe the submissions from participating teams. An overview of the evaluated metrics can be seen in Table 6.

### 4.1 Baselines

**SacreBLEU baselines** We use the following metrics from SacreBLEU (Post, 2018) as baselines:

- **BLEU (Papineni et al., 2002)** is based on the precision of $n$-grams between the MT output and its reference, weighted by a brevity penalty. Using SacreBLEU we obtained sentence-BLEU values using the `sentence_bleu` Python function and for corpus-level BLEU we used `corpus_bleu` (both with default arguments[8]).

- **SPBLEU (NLLB Team et al., 2022)** are BLEU scores computed with subword tokenization by the standardized FLORES-200 Sentencepiece models. We used the command line SacreBLEU to compute the sentence level SPBLEU[9] and we averaged the segment-level scores to obtain a corpus-level score.

- **CHRF (Popović, 2015)** uses character $n$-grams instead of word $n$-grams to compare the MT output with the reference. For CHRF we used the SacreBLEU `sentence_chrf` function (with default arguments[10]) for segment-level scores and we average those scores to obtain a corpus-level score.

---

[6]Since the testset is for Spanish from Mexico rather than Spanish from Spain, the conducted annotations were collected taking that variant in consideration.

[7]see Unbabel Annotation Guidelines - Typology 3.0

[8]lnrefs.1|case.mixed|lang.LANGPAIR|tok.13a|smooth.exp| version.2.3.0. For to-zh and to-ja language pairs, we use tok.zh and tok.ja-mecab

[9]nrefs:1|case:mixed|eff:yes|tok:flores200|smooth:exp| version:2.3.0

[10]chrF2||lang.LANGPAIR|nchars.6|space.false|version.2.3.0

| language | news | social | speech | literary |
|---|---|---|---|---|
| en→de | 90/149 (17/17) | 258/531 (34/34) | 111/111 (1/1) | 27/206 (8/8) |
| en→es | 124/149 (14/17) | 281/531 (20/34) | 107/111 (1/1) | 110/206 (5/8) |
| ja→zh | 255/269 (45/45) | n/a | 136/136 (1/1) | 168/316 (15/15) |

Table 4: Numbers of MQM-annotated segments per domain (number of docs in brackets).

**BERTSCORE (Zhang et al., 2020)** leverages contextual embeddings from pre-trained transformers to create soft-alignments between words in candidate and reference sentences using cosine similarity. Based on the alignment matrix, BERTSCORE returns a precision, recall and F1 score. We used F1 without TF-IDF weighting.

**BLEURT (Sellam et al., 2020)** is a learned metric fine-tuned on Direct Assessments (DA). Unlike COMET, BLEURT encodes the translation and the reference together and utilizes the `[CLS]` token as an embedding to represent the pair. We employed the BLEURT20 checkpoint (Pu et al., 2021), which was trained on top of RemBERT using DA data from previous shared tasks spanning from 2015 to 2019, along with additional synthetic data created from Wikipedia articles.

**COMET-22 (Rei et al., 2022a)** is a learned metric fine-tuned using DA from previous WMT Translation shared tasks. This metric relies on sentence embeddings from the source, translation, and reference to produce a final score. We utilized the default model `wmt22-comet-da` provided in version 2.0.2 of the `Unbabel/COMET` framework. This model employs XLM-R large as its backbone model and is trained on data from the 2017 to 2019 WMT shared tasks, in combination with the MLQE-PE corpus (Fomicheva et al., 2022).

**COMETKIWI (Rei et al., 2022b)** is a reference-free learned metric that functions similarly to BLEURT, but instead of encoding the translation along with its reference, it uses the source. We utilized the `wmt22-cometkiwi-da` model, which was a top-performing reference-free metric from the WMT22 shared task. This reference-free metric is fine-tuned on the same data as `wmt22-comet-da` using the version 2.0.2 of the `Unbabel/COMET` framework.

**PRISMREFSMALL AND PRISMREFMEDIUM (Thompson and Post, 2020a,b)** are both reference-based PRISM that uses a multilingual MT model in zero-shot paraphrase model to score the candidate translation conditioned on the reference, and

the reference conditioned on the candidate translation, and averages the two scores. As LLMs have become quite capable multi-lingual MT models, we opted to use Llama3.1 (Llama Team, 2024) as the underlying MT model this year. PRISMREFSMALL corresponds to Llama3.1 8B and PRISMREFMEDIUM corresponds to Llama3.1 70B. The long context window of LLMs allows us to compute scores for entire documents, while still averaging scores for each sentence to produce sentence-level scores (Vernikos et al., 2022). We chunked longer documents into sub-documents of up to 10 sentences, and added a penalty for producing no output.

**YISI-1 (Lo, 2019)** is an MT evaluation metric that measures the semantic similarity between a machine translation and human references by aggregating the IDF-weighted lexical semantic similarities based on the contextual embeddings extracted from pre-trained language models (e.g. RoBERTa, CamemBERT, XLM-RoBERTa, etc.).

### 4.2 Metric Submissions

The rest of this section summarizes the participating metrics.

**BLCOM_1 and BLCOM** Unfortunately, we have no information about these submission.

**BRIGHT-QE** is a referenceless metric, which uses the XLM-XL encoder to perform multi-stage fine-tuning according to the XCOMET framework. In the first stage of training, we used DA 2017 2022 corpus, and gradually reduced the weight of REF-based loss with the idea of curriculum learning, trying to reduce the model's dependence on reference and better align the semantics of the translation and source text; in the second stage, we used batch softmax to normalize scores, and introduced KL divergence loss to learn to modify the minor rank error that MSE loss cannot solve, so as to obtain better Pearson correlation; finally, we further fine-tuned on high-quality MQM corpus to achieve better consistency with human expert MQM.

| English→German ↓ | | | | | |
|---|---|---|---|---|---|
| System | all | news | social | speech | literary |
| Dubformer | 1.58 | 1.29 | 0.60 | 4.22 | 1.15 |
| GPT-4 | 1.58 | 1.39 | 0.88 | 3.60 | 0.69 |
| Unbabel-Tower70B | 1.65 | 1.99 | 0.78 | 3.46 | 1.41 |
| ONLINE-B | 1.81 | 1.48 | 1.22 | 3.59 | 1.30 |
| TranssionMT | 1.81 | 1.24 | 1.18 | 3.87 | 1.33 |
| refB | 1.84 | 1.38 | 0.80 | 4.92 | 0.81 |
| Mistral-Large | 1.93 | 1.95 | 1.12 | 3.91 | 1.46 |
| CommandR-plus | 2.01 | 2.40 | 1.07 | 3.95 | 1.74 |
| refA | 2.12 | 1.84 | 1.01 | 4.96 | 2.04 |
| Gemini-1.5-Pro | 2.20 | 1.29 | 1.93 | 2.90 | 4.97 |
| ONLINE-W | 2.22 | 1.32 | 1.75 | 4.09 | 2.12 |
| Claude-3.5 | 2.28 | 1.00 | 1.23 | 6.04 | 1.13 |
| IOL_Research | 2.39 | 1.66 | 1.61 | 4.91 | 2.01 |
| Aya23 | 3.09 | 2.69 | 2.20 | 5.71 | 2.26 |
| ONLINE-A | 3.30 | 1.93 | 2.29 | 6.88 | 2.85 |
| Llama3-70B | 3.62 | 2.91 | 2.28 | 7.08 | 4.76 |
| IKUN | 3.86 | 4.35 | 2.36 | 7.09 | 3.48 |
| IKUN-C | 5.07 | 3.39 | 3.34 | 9.87 | 7.63 |
| MSLC | 13.46 | 11.54 | 8.24 | 26.80 | 15.29 |

| English→Spanish ↓ | | | | | |
|---|---|---|---|---|---|
| System | all | news | social | speech | literary |
| GPT-4 | 0.12 | 0.03 | 0.14 | 0.24 | 0.03 |
| Unbabel-Tower70B | 0.20 | 0.21 | 0.04 | 0.68 | 0.14 |
| Claude-3.5 | 0.26 | 0.06 | 0.21 | 0.60 | 0.29 |
| Mistral-Large | 0.26 | 0.16 | 0.28 | 0.50 | 0.12 |
| Gemini-1.5-Pro | 0.39 | 0.18 | 0.56 | 0.54 | 0.06 |
| Dubformer | 0.43 | 0.29 | 0.07 | 2.00 | 0.01 |
| Llama3-70B | 0.52 | 0.10 | 0.28 | 2.17 | 0.02 |
| refA | 0.55 | 0.20 | 0.12 | 2.42 | 0.20 |
| IOL_Research | 0.57 | 0.44 | 0.33 | 1.39 | 0.56 |
| CommandR-plus | 0.62 | 0.50 | 0.34 | 0.52 | 1.55 |
| ONLINE-W | 0.64 | 0.17 | 0.27 | 2.36 | 0.46 |
| IKUN | 0.94 | 0.86 | 0.74 | 1.01 | 1.46 |
| ONLINE-B | 1.08 | 1.01 | 0.59 | 1.76 | 1.77 |
| Aya23 | 1.52 | 1.52 | 1.09 | 2.03 | 2.12 |
| MSLC | 6.80 | 4.09 | 4.63 | 10.99 | 11.36 |

| Japanese→Chinese ↓ | | | | |
|---|---|---|---|---|
| System | all | news | speech | literary |
| Claude-3.5 | 1.22 | 0.76 | 2.96 | 0.76 |
| refA | 1.32 | 0.77 | 3.15 | 0.77 |
| GPT-4 | 1.45 | 0.82 | 3.25 | 0.82 |
| DLUT_GTCOM | 1.52 | 1.06 | 3.66 | 1.06 |
| Unbabel-Tower70B | 1.69 | 1.16 | 3.53 | 1.16 |
| Gemini-1.5-Pro | 1.78 | 0.84 | 3.80 | 0.84 |
| CommandR-plus | 1.91 | 1.28 | 4.61 | 1.28 |
| IOL_Research | 2.10 | 1.14 | 4.82 | 1.14 |
| Aya23 | 3.03 | 1.86 | 6.44 | 1.86 |
| Llama3-70B | 3.07 | 2.16 | 6.16 | 2.16 |
| Team-J | 3.91 | 2.02 | 8.46 | 2.02 |
| NTTSU | 4.34 | 2.11 | 10.51 | 2.11 |
| ONLINE-B | 5.27 | 3.72 | 9.52 | 3.72 |
| IKUN-C | 6.60 | 3.45 | 14.41 | 3.45 |
| MSLC | 9.19 | 4.01 | 19.04 | 4.01 |

Table 5: MQM human evaluations for generalMT2024. Lower average error counts represent higher MT quality. Systems above any solid line are significantly better than those below, based on all domains with p < 0.05.

**CHRFS (Mukherjee and Shrivastava, 2024)** is an unsupervised reference-based metric, a semantic version of CHRF++ that integrates sentence embeddings to evaluate translation quality more comprehensively. By combining traditional character and word n-gram analysis with semantic information derived from embeddings, CHRFS captures both syntactic accuracy and sentence-level semantics.

**DAMONMONLI and MONMONLI** is a proof-of-concept of multiple ideas. A multi-lingual NLI model is used to extract embeddings for (mt, src) and (mt, ref) pairs, based on findings of Chen and Eger (2023). A multi-task learning approach is employed where different human annotations from WMT22 and WMT23 are used as different tasks. For each task, it uses a separate regression head that learns a monotonic function of the metric's score(Runje and Shankaranarayana, 2023). The main metric "DAMONMONLI" also includes a domain adversarial loss (Ganin and Lempitsky, 2015) to make metric representations robust against shifts in MT systems and language pairs.

**GEMBA-ESA (Kocmi and Federmann, 2023)** is an extension of previous work on an LLM-based metric, with an updated prompt to reflect the new human evaluation protocol ESA (Kocmi et al., 2024c) used at WMT General MT task. It contains a two-step approach where in the first step, MQM error spans are collected and in a second step, the final score is assigned.

**MEE4 (Mukherjee and Shrivastava, 2023a)** is an unsupervised, reference-based metric (an improved version of MEE) focusing on computing contextual and syntactic equivalences, along with lexical, morphological, and semantic similarity. The goal is to comprehensively evaluate the fluency and adequacy of MT outputs while also considering the surrounding context. Fluency is determined by analysing syntactic correlations, while context is evaluated by comparing sentence similarities using sentence embeddings. The ultimate score is derived from a weighted amalgamation of three distinct similarity measures: a) Syntactic similarity, which is established using a modified BLEU score. b) Lexical, morphological, and semantic similarity, quantified through explicit unigram matching. c) Contextual similarity, gauged by sentence similarity scores obtained from the Language-Agnostic BERT model.

**METAMETRICS-MT** (Anugraha et al., 2024; Winata et al., 2024) is a machine translation

| | metric | broad category | supervised | ref. free | citation | availability (https://github.com/) |
|---|---|---|---|---|---|---|
| **baselines** | BLEU | lexical overlap | | | Papineni et al. (2002) | mjpost/sacrebleu |
| | SPBLEU | lexical overlap | | | NLLB Team et al. (2022) | mjpost/sacrebleu |
| | CHRF | lexical overlap | | | Popović (2015) | mjpost/sacrebleu |
| | BERTSCORE | embedding similarity | | | Zhang et al. (2020) | Tiiiger/bert_score |
| | BLEURT-20 | fine-tuned metric | ✓ | | Sellam et al. (2020) | google-research/bleurt |
| | COMET-22 | fine-tuned metric | ✓ | | Rei et al. (2022a) | Unbabel/COMET |
| | COMETKIWI | fine-tuned metric | ✓ | ✓ | Rei et al. (2022b) | Unbabel/COMET |
| | PRISMREFSMALL | MT-model metric | | | Thompson and Post (2020a,b) | thompsonb/prism |
| | PRISMREFMEDIUM | MT-model metric | | | Thompson and Post (2020a,b) | thompsonb/prism |
| | YISI-1 | embedding similarity | | | Lo (2019) | chikiulo/yisi |
| **primary submissions** | BLCOM_1 | N/A | N/A | N/A | N/A | (not available) |
| | BRIGHT-QE | fine-tuned metric | ✓ | ✓ | N/A | https://bright.pcl.ac.cn/en/ |
| | CHRFS | lexical and embedding similarity | | | (Mukherjee and Shrivastava, 2024) | AnanyaCoder/chrF-S |
| | COMETKIWI-XXL | fine-tuned metric | ✓ | ✓ | Rei et al. (2023) | Unbabel/COMET |
| | DAMONMONLI | finetuned metric | ✓ | ✓ | N/A | (not available) |
| | GEMBA-ESA | LLM prompt-based metric | | ✓ | Kocmi and Federmann (2023) | MicrosoftTranslator/GEMBA |
| | MEE4 | lexical & embedding similarity | | | Mukherjee and Shrivastava (2023b) | AnanyaCoder/WMT22Submission |
| | METAMETRICS-MT | ensemble metric | ✓ | | Anugraha et al. (2024) | meta-metrics/metametrics |
| | METAMETRICS-MT-QE | ensemble metric | ✓ | ✓ | Anugraha et al. (2024) | gentaiscool/meta-metrics |
| | METRICX-24-HYBRID | fine-tuned metric | ✓ | | Juraska et al. (2024) | google-research/metricx |
| | METRICX-24-HYBRID-QE | fine-tuned metric | ✓ | ✓ | Juraska et al. (2024) | google-research/metricx |
| | SENTINEL-CAND-MQM | fine-tuned metric | ✓ | ✓ | Perrella et al. (2024) | SapienzaNLP/guardians-mt-eval |
| | SENTINEL-REF-MQM | fine-tuned metric | ✓ | | Perrella et al. (2024) | SapienzaNLP/guardians-mt-eval |
| | SENTINEL-SRC-MQM | fine-tuned metric | ✓ | ✓ | Perrella et al. (2024) | SapienzaNLP/guardians-mt-eval |
| | XCOMET | fine-tuned metric | ✓ | | Guerreiro et al. (2023) | Unbabel/COMET |
| | XCOMET-QE | fine-tuned metric | ✓ | ✓ | Guerreiro et al. (2023) | Unbabel/COMET |
| | XLSIMMQM | fine-tuned metric | ✓ | | Mukherjee and Shrivastava (2023b) | AnanyaCoder/XLsim |

Table 6: Baseline metrics and primary submissions for the metrics task. Supervised metrics are trained on MT evaluation data such as DA or MQM scores.

(MT) metric developed from our METAMET-RICS (Winata et al., 2024), specifically designed to better align with human preferences using Bayesian optimization with Gaussian Processes (GP). By systematically integrating multiple existing metrics, we create a sparse allocation that only includes metrics enhancing the overall correlation score. We optimize this metric by maximizing Kendall scores from the WMT shared task (MQM) 2020-2022. METAMETRICS-MT achieves state-of-the-art performance for reference-based metrics, while its reference-free variant, METAMETRICS-MT-QE, demonstrates competitive correlation with human scores in the WMT24 metric shared task. By strategically assigning weights to combined metrics, METAMETRICS-MT aims to be as competitive as, if not superior to, any individual metric. To address missing values when reference data is unavailable, we propose a hybrid variant, METAMETRICS-MT-HYBRID, which utilizes both metrics to compensate for the absence of reference data in the reference-based setting.

**METRICX-24 (Juraska et al., 2024)** is a learned regression-based metric that builds on top of its predecessor from 2023. Similar to METRICX-23, it is based on the mT5-XXL pretrained language model, which is fine-tuned in two stages on DA and MQM scores from WMT 2015-22, and it implements three major design improvements. First, the training data in both stages is augmented with synthetic examples to make the metric more robust to several common failure modes, such as fluent but unrelated translation, or undertranslation. Second, a small proportion of DA data is mixed in during the second stage of fine-tuning in order to preserve the performance on non-MQM language pairs. Finally, the model's training is done on a mixture of examples that include the source only, the reference only, or both, which allows the model to operate in both a QE and a reference-based mode (and the latter either with or without the source included). Hence, both METRICX-24-HYBRID and METRICX-24-HYBRID-QE submission are in fact the exact same model, only with the references excluded from the input in the latter case.

**SENTINEL-CAND-MQM, SENTINEL-REF-MQM and SENTINEL-SRC-MQM (Perrella et al., 2024)** are designed explicitly to scrutinize the accuracy, robustness, and fairness of the meta-evaluation process. The three sentinel metrics are trained only on the candidate, reference and source sentence re-

spectively on DA and MQM data from WMT 2017 to 2022.

**XCOMET AND XCOMET-QE (Guerreiro et al., 2023)** models are trained using both a sentence-level signal and span-level supervision coming from MQM data from previous years, along with some synthetic data that mimics hallucinations. We ensemble XCOMET-XXL and XCOMET-XL to give a single unified score.

**XLSIMMQM (Mukherjee and Shrivastava, 2023b)** is an enhanced version of XLSIM, a supervised reference-based evaluation metric, which we have transformed into a reference-free model to improve its applicability across multiple language pairs. Unlike the original XLSIM, which was limited to the English-German language pair, XLSIMMQM is trained on a filtered comprehensive dataset curated from WMT-MQM (2020-22), ensuring broader applicability and robustness. The filtered datasets (train, dev and test) contains uniform distribution across good, medium and poor-quality sentences; this careful balancing of the dataset leads to a better, reliable and robust metric.

## 5 Meta Evaluation

The goal of metric meta-evaluation is to quantify how well automatic metrics agree with human ratings of translation quality. There are a multitude of ways to approach this problem, as evidenced by the variety of solutions proposed by previous years' editions of the shared task. For instance—to name just a few possible design decisions—the agreement can be measured at the system or segment level; the agreement function can be Pearson, Spearman, Kendall, pairwise agreement, or $L_2$ loss; the agreement can be computed per domain or on the full dataset. None of these approaches are necessarily right or wrong, but rather each method evaluates a different property of the metric.

Because there is no one way to evaluate a metric, the past two iterations of the Metrics Shared Task defined a variety of "tasks" (or different configurations of meta-evaluations) that evaluated some aspect of a metric, then calculated an overall quality score by averaging the individual task scores. Implicitly, this approach defines a "high-quality" metric as one that performs well across the tasks on average. In 2022, there were 201 tasks that varied along dimensions such as language pair, domain, correlation granularity, correlation statistic, etc. In

2023, the number of tasks was reduced to 10, measuring only pairwise accuracy and Pearson at both the system and segment levels.

For this year's meta-evaluation, we follow the same approach of averaging performance across tasks, but focus the tasks to better align with how evaluation metrics are used in practice. The two main use cases that we targeted were using metrics to rank a set of MT systems and using a metric to rank a set of translations for the same source segment. The former setting is widely used by academics and practitioners in industry to determine whether one model produces better translations than another, and the latter setting has applications in Minimum Bayes Risk Decoding and Quality Estimation Reranking either directly as decoding method (Fernandes et al., 2022; Freitag et al., 2022) or to further fine-tune models (Finkelstein and Freitag, 2024; Finkelstein et al., 2024). The latter one is getting more popular and can introduce metric biases (Kovacs et al., 2024) that is an emerging challenge for metrics. As such, we defined one task to quantify how well metrics work for each of these two use cases separately for all three language pairs, resulting in a total of six tasks.

At the system-level, we use the recently proposed metric called soft pairwise accuracy, or SPA (Thompson et al., 2024). One of the drawbacks of standard pairwise accuracy (or the very related Kendall's $\tau$) that has been used in previous years' shared tasks is that it does not account for the uncertainty of the system ranking. For example, if the human ranking of two systems is almost arbitrary (e.g, a statistical tie) but the metric ranking is quite certain, standard pairwise accuracy will either reward or penalize the metric nearly randomly. The reverse case—a certain human ranking and uncertain metric ranking—also nearly arbitrarily rewards or penalizes metrics. If both rankings are uncertain, the metric will again be rewarded nearly randomly, and the penalty for an incorrect ranking is equal to when the metric was very certain but also wrong.

SPA addresses this problem by using $p$-values as a proxy for certainty, calculating $p$-values between two systems using both the metric and human scores, then taking 1.0 minus the absolute difference between the two $p$-values as the metric's score for that pair. This rewards metrics that result in the same statistical conclusion as the human scores. Now, statistical ties do not randomly reward or penalize metrics, but instead the score is proportional to whether or not the metric and human have

| language | ref used | scored ref |
|----------|----------|------------|
| en→de | B | A |
| en→es | A | – |
| ja→zh | A | – |

Table 7: Use of reference translations.

| task | lang | level | correlation | wt |
|------|------|-------|-------------|-----|
| 1 | en→de | system | SPA | 1 |
| 2 | en→de | segment | $acc^*_{eq}$ | 1 |
| 3 | en→es | system | SPA | 1 |
| 4 | en→es | segment | $acc^*_{eq}$ | 1 |
| 5 | ja→zh | system | SPA | 1 |
| 6 | ja→zh | segment | $acc^*_{eq}$ | 1 |

Table 8: For each language pair, soft pairwise accuracy (SPA) was used at the system-level and $acc^*_{eq}$ at the segment-level. Each task was given equal weight in the overall average. See §5 for explanations of SPA and $acc^*_{eq}$.

the same level of certainty in the ranking.

At the segment-level, we follow last year's meta-evaluation and meta-evaluate metrics using "group-by-item" segment-level accuracy with tie calibration (Deutsch et al., 2023) denoted $acc^*_{eq}$.

The six tasks (shown in Table 8) receive equal weighting in the overall average, which is the final score for the metric.

**Removing Pearson's Correlation:** Notably, the meta-evaluation this year only focuses on evaluating rankings and does not include any correlation that evaluates the absolute value of the scores predicted by metrics, like Pearson's correlation. This decision was made because using metrics to rank systems or translations is much more common in practice than using a metric to approximate the absolute quality score as derived by humans, which is more similar to a Pearson correlation.

**Limitations:** Like previous years, we acknowledge that this approach is not perfect. One problem is that we need to combine correlations and accuracies that may have different dynamic ranges, which could result in certain tasks carrying more weight than others in the overall ranking. However, to simplify the implementation, we assigned equal weight to all tasks, which worked well in last year's evaluation.

### 5.1 Rank Assignment

For each task, we assign ranks to metrics based on their significance clusters in the same way that we

did last year, detailed below.

We compare all pairs of metrics and determine whether the difference in their correlation scores is significant, according to the PERM-BOTH hypothesis test of Deutsch et al. (2021). We use 1000 resampling runs and set $p = 0.05$. As advocated by Wei et al. (2022), we divide the sample into blocks of 100, compute significance after each block (cumulative over all blocks sampled so far), and stop early if the p-value is $< 0.02$ or $> 0.50$.

The $\mathrm{acc}^*_{eq}$ statistic creates a problem for significance testing because it optimizes a latent tie threshold for each metric on each test set (just one threshold for all item-wise score vectors). Since the permutation test for comparing two metrics creates two new vectors by randomly swapping elements of the original vectors on each draw, this necessitates the very expensive step of finding two new tie thresholds for each draw. To reduce the expense, we used the following approximate procedure. First find an optimal threshold for each input metric on the current test set, then create all pairs of item-wise scores and assign a correct/incorrect status to each pair by examining whether the metric's ranking matches the human ranking. Then perform the permutation test on these pairwise status vectors rather than the original score vectors. This approximation has more degrees of freedom than the original test, and can sample pairs that would never result from swapping the original score vectors, but our experiments showed that it is a reasonable proxy for the correct procedure.

To compute overall p-values based on weighted average scores of two metrics across all tasks, we cache the results of the draws for the per-task significance tests. In all cases, these are vectors of $K$ pairs of correlation or accuracy statistics. Where $K < 1000$ due to early stopping, we duplicate elements to get 1000 examples. Then for $i$ in 1..1000 we compare the weighted average of the pairs from the $i$th draw across all tasks, and record the results to produce an overall p-value.

**Clustering.** Given significance results (p-values) for all pairs of metrics, we assign ranks as follows. Starting with the highest-scoring metric, we move down the list of metrics in descending order by score, and assign rank 1 to all metrics until we encounter the first metric that is significantly different from any that have been visited so far. That metric is assigned rank 2, and the process is repeated. This continues until all metrics have been assigned

a rank. Note that this is a greedy algorithm, and hence it can place two metrics that are statistically indistinguishable in different clusters.

## 5.2 Implementation Details

The code for running the meta-evaluation is available in the MT Metrics Eval library.[11]

To calculate $p$-values for SPA, we use a paired permutation test (Noreen, 1989) with 1k resamples.

In previous years' shared tasks, tasks were categorized based on whether they included additional reference translations in the overall system ranking. Following last year's proposal, we always include the additional reference in the overall ranking. This year, this only applies to en→de which is the only language pair with more than one reference translation (see Table 7).

Out of all the submitted MT systems, MSLC consistently scores well below the other systems for all language pairs and was identified as an outlier and removed from the correlation calculation.

## 6 Main Results

As we have described in Section 5, the final statistic used to rank the metrics is defined as the average of the results from the six main tasks (system-level and segment-level tasks in different language pairs). Table 1 shows the official scores and rankings of all baselines and primary submissions. Table 9 shows the scores and rankings of each individual task at system level and segment level, respectively. Similar to last year's results, neural metrics perform significantly better than lexical metrics. Of the 26 evaluated metrics, BLEU, SPBLEU and CHRF are ranked 23rd, 22nd and 20th respectively. Fine-tuned neural metrics, like XCOMET and METRICX-23 are the highest ranked non-ensemble metrics. The ensemble submission METAMETRIC_MT is in the same significance cluster as XCOMET and METRICX-24-HYBRID, but relies heavily on the 2023 version of METRICX-24-HYBRID. Like last year, QE metrics perform very well, with METRICX-24-HYBRID-QE and GEMBA_ESA sharing the second significance cluster.

Figure 1 shows the correlation scores split by language pair. Interestingly, GEMBA_ESA is performing very well for en→es and ja→zh, while ranked below many metrics for en→de. GEMBA_ESA is

---

[11]https://github.com/google-research/mt-metrics-eval

| Metric | avg-corr | | en-de sys SPA task1 | | en-de seg acc$^*_{eq}$ task2 | | en-es sys SPA task3 | | en-es seg acc$^*_{eq}$ task4 | | ja-zh sys SPA task5 | | ja-zh seg acc$^*_{eq}$ task6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MetaMetrics-MT | **1** | **0.725** | 2 | 0.883 | **1** | **0.542** | **1** | **0.804** | 2 | 0.686 | 2 | 0.873 | **1** | **0.561** |
| MetricX-24-Hybrid | **1** | **0.721** | 2 | 0.874 | 2 | 0.532 | 2 | 0.799 | 3 | 0.685 | **1** | **0.897** | 2 | 0.539 |
| XCOMET | **1** | **0.719** | **1** | **0.905** | 2 | 0.530 | 2 | 0.791 | **1** | **0.688** | **1** | **0.890** | 5 | 0.510 |
| MetricX-24-Hybrid-QE* | 2 | 0.714 | 2 | 0.878 | 3 | 0.526 | 2 | 0.789 | 4 | 0.685 | 2 | 0.875 | 3 | 0.530 |
| gemba_esa* | 2 | 0.711 | 4 | 0.793 | 5 | 0.507 | **1** | **0.838** | 5 | 0.683 | **1** | **0.908** | 2 | 0.539 |
| XCOMET-QE* | 3 | 0.695 | **1** | **0.889** | 4 | 0.520 | **1** | **0.801** | 2 | 0.687 | 4 | 0.808 | 10 | 0.463 |
| COMET-22 | 3 | 0.688 | 2 | 0.879 | 8 | 0.482 | 2 | 0.778 | 5 | 0.683 | 4 | 0.813 | 6 | 0.496 |
| BLEURT-20 | 3 | 0.686 | 2 | 0.881 | 7 | 0.486 | 3 | 0.695 | 6 | 0.681 | **1** | **0.887** | 8 | 0.484 |
| MetaMetrics-MT-QE* | 3 | 0.684 | 2 | 0.860 | 6 | 0.497 | 3 | 0.711 | 2 | 0.686 | 3 | 0.837 | 4 | 0.516 |
| bright-qe* | 4 | 0.681 | 3 | 0.816 | 6 | 0.500 | 2 | 0.792 | **1** | **0.689** | 4 | 0.805 | 8 | 0.484 |
| BLCOM_1 | 4 | 0.664 | 3 | 0.840 | 10 | 0.455 | 3 | 0.680 | 6 | 0.681 | 3 | 0.843 | 7 | 0.488 |
| sentinel-cand-mqm* | 5 | 0.650 | 3 | 0.822 | 4 | 0.517 | 2 | 0.785 | 4 | 0.683 | 7 | 0.610 | 8 | 0.481 |
| PrismRefMedium | 5 | 0.646 | 4 | 0.776 | 14 | 0.434 | 3 | 0.652 | 7 | 0.680 | 2 | 0.872 | 10 | 0.462 |
| PrismRefSmall | 5 | 0.642 | 4 | 0.772 | 14 | 0.433 | 4 | 0.634 | 8 | 0.680 | 2 | 0.875 | 11 | 0.457 |
| CometKiwi* | 5 | 0.640 | 5 | 0.732 | 9 | 0.467 | 3 | 0.693 | 4 | 0.684 | 5 | 0.776 | 7 | 0.490 |
| damonmonli | 5 | 0.635 | 5 | 0.696 | 12 | 0.443 | 4 | 0.607 | 6 | 0.682 | **1** | **0.911** | 9 | 0.472 |
| YiSi-1 | 6 | 0.630 | 4 | 0.759 | 13 | 0.436 | 4 | 0.609 | 7 | 0.681 | 3 | 0.835 | 11 | 0.458 |
| BERTScore | 7 | 0.617 | 4 | 0.749 | 14 | 0.435 | 4 | 0.587 | 6 | 0.682 | 4 | 0.799 | 12 | 0.451 |
| MEE4 | 7 | 0.609 | 5 | 0.731 | 13 | 0.437 | 7 | 0.504 | 4 | 0.683 | 2 | 0.855 | 13 | 0.446 |
| chrF | 8 | 0.608 | 4 | 0.750 | 15 | 0.431 | 5 | 0.581 | 8 | 0.680 | 5 | 0.767 | 16 | 0.436 |
| chrfS | 8 | 0.606 | 4 | 0.742 | 14 | 0.434 | 6 | 0.549 | 6 | 0.682 | 4 | 0.788 | 14 | 0.444 |
| spBLEU | 9 | 0.593 | 4 | 0.741 | 17 | 0.431 | 6 | 0.523 | 7 | 0.680 | 6 | 0.744 | 16 | 0.436 |
| BLEU | 9 | 0.589 | 4 | 0.736 | 16 | 0.431 | 6 | 0.512 | 8 | 0.680 | 6 | 0.740 | 17 | 0.435 |
| XLsimMqm* | 10 | 0.515 | 6 | 0.612 | 11 | 0.450 | 8 | 0.359 | 7 | 0.681 | 7 | 0.548 | 15 | 0.438 |
| sentinel-src-mqm* | 10 | 0.513 | 7 | 0.406 | 18 | 0.429 | 5 | 0.580 | 8 | 0.680 | 8 | 0.546 | 17 | 0.435 |
| sentinel-ref-mqm | 10 | 0.513 | 7 | 0.405 | 18 | 0.429 | 4 | 0.581 | 8 | 0.680 | 8 | 0.545 | 17 | 0.435 |

Table 9: Correlation results per task for the main language pairs. See §5 for descriptions of soft pairwise accuracy (SPA) and acc$^*_{eq}$. Rows are sorted by the overall average correlation across all 6 tasks (leftmost column). Starred metrics are reference-free, and underlined metrics are baselines.

a prompt-based metric and not fine-tuned for any metric task. Both en→es and ja→zh are new language pairs, and no fine-tuning data exists which might have played in disadvantage for all fine-tuned metrics.

We continue to be interested in metrics' abilities to generalise across domains. In Figure 2, we present the performance of each metric across different domains. Similar to last year, we observe that neural metrics perform better than lexical overlap metrics across all four domains. Figure 3 shows the average correlations of metrics when grouped separately by system-level and segment-level tasks. There is a high correlation between the rankings of both granularities.

## 7 Beyond accuracy and correlation

Last year, we conducted two additional analyses beyond correlation with human scores to find the threshold of metrics' score differences correspond to statistical significance of MT system rankings demonstrated by human annotators and the metrics themselves. Despite the better correlation with human judgements achieved by new neural metrics, BLEU remains as the most used metric in the MT research community. One of the reasons is that MT researchers have established some "shared understanding" about the relationship between BLEU and the actual translation quality, and similar intuitions about new metrics have yet to crystallize. Our analyses beyond correlation provided an interpretation of the metrics' score differences. Hence, we are continuing such analyses to support building an intuitive sense of metric score meanings and encourage broader adoption of new automatic MT evaluation metrics. As a reminder, our results should *NOT* be used as arguments to forego significance tests or appropriate human evaluation.

Figure 1: Average metrics' meta-evaluation scores in tasks grouped by language pair.



Figure 2: Average metrics' correlation with human in tasks grouped by domain.



Figure 3: Average metrics' correlation with human in tasks grouped by granularity level.

## 7.1 Correspondence to MQM scores significance

We first study the relationship between statistically significant differences in human scores and the magnitude of metric differences as in Lo et al. (2023a). We run a two-sided paired t-test with an equal variance assumption for each system pair on segment-level MQM scores. After that, we fit the corresponding metric score differences and the p-values of the t-test on the MQM scores to an isotonic regression (Robertson et al., 1988), that predicts whether the human MQM score difference will be significant given the metric's score difference. Isotonic regression produces a non-decreasing function where the classifier output can be interpreted as a confidence level.[12] We set $p_{mqm} < 0.05$ as the significance level of MQM scores. Thus, the output of the isotonic regression function can be viewed as $Pr(p_{mqm} < 0.05 | \Delta M)$ where $p_{mqm}$ is the p-value of the t-test on the MQM scores for each system pair and $\Delta M$ is the metric score difference.

Figure 4 shows the (log) p-value of two-sided paired t-test on the MQM scores against the corre-

---

[12]https://scikit-learn.org/stable/modules/isotonic.html

sponding BLEU and COMET-22 score difference for each system pair in en→de. Figures 6-10 in appendix D, show the same analyses for all metrics and language pairs. For each metric, we can choose a particular level of confidence (i.e., a point along the y-axis on the right) to give metric score difference cut-offs (i.e., a point along the x-axis) that this metric difference reflects significant MQM score differences. Drawing a horizontal line from the confidence level, say 80%, to the red line enables us to find the minimum metric difference cut-off required at the corresponding x-value down from the red line, i.e. 5.4 for BLEU in Figure 4. Using this lookup method, Table 10 shows the cut-offs of $\Delta M$ when $Pr(p_{mqm} < 0.05|\Delta M) = 0.8$ for each metric and language pair.

We run the leave-one-system-out cross validation and Table 10 shows that the range of precision in the cross validation are consistently high across metrics, except for BLEU, BRIGHT-QE, COMETKIWI, MEE4, METAMETRICS_MT_MQM_QE_KENDALL.SEG.S, SPBLEU and XLSIMMQM. This means the metric cut-offs we find using the regression model are reliable.

Contrary to the shared understanding that 2 BLEU improvement represents "significant" or "notable by human" improvement in the actual translation quality, our analyses show that 5.4 BLEU improvement is required to be confident (80%) that the MQM scores would be different with statistical significance for en→de and that threshold would be as high as 11 BLEU for en→es. Table 10 serves as a reference between BLEU differences and differences in some of the modern metrics and assists metric users in understanding scores provided by modern metrics. For example, when evaluating ja→zh translation quality, we see that a BLEU difference of 1.4 corresponds to 80% confidence that the metric's ranking of the two MT systems will match the decision made by human annotators with a significant difference. Meanwhile, a COMET-22 score difference of 0.021 would have the same 80% chance of human judged significant difference.

### 7.2 Correspondence to metric scores significance

We run a study similar to that in the previous subsection but on the relations between statistically significant differences in metric scores and the magnitude of metric differences as inspired by Marie (2022). Instead of the two-sided t-test on MQM, the p-values are now obtained by running statis-

tical significance tests with bootstrap resampling on the metric scores for each system pair. We fit the corresponding metric score differences and the p-values of the significance test to an isotonic regression for predicting whether the translation quality improvement as indicated by the metric will be significant given the metric score difference. We set $p_M < 0.05$ and thus, the output of the isotonic regression function is now $Pr(p_M < 0.05|\Delta M)$, where $p_M$ is the p-value of the significance test on the metric scores for each system pair and $\Delta M$ is the metric score difference.

Figure 5 shows the (log) p-value of the significance test with bootstrap resampling on the metric scores for BLEU and COMET-22 score difference of each system pair in en→de. Additional figures (Figures 11-15 in appendix Appendix D) show the same analyses for all metrics and language pairs. Using the same lookup method described in the previous subsection, Table 11 shows the cut-offs of $\Delta M$ when $Pr(p_M < 0.05|\Delta M) = 0.8$ for each metric and language pair.

We run the leave-one-system-out cross validation, and Table 11 shows that the range of precision in the cross validation are consistently high across metrics. This means the metric cut-offs we find using the regression model are reliable.

Table 11 serves as a reference of metric differences that correspond to statistical significance with high confidence. For example, when evaluating en→de translation quality, we see that a BLEU difference of 0.97 corresponds to 80% confidence the difference is statistically significant. Meanwhile, a COMET-22 score difference of 0.0043 would have the same 80% chance of statistical significance. Our results, agreeing with Marie (2022), show that to claim significant differences ($p_M < 0.05$) in BLEU with high confidence (80%), the differences should be much higher than the shared understanding of 0.5 BLEU, ranging from 0.89 to 0.97 for the three language pairs.

Closely related to this analysis, Kocmi et al. (2024b) investigated the agreement between human evaluations and metric differences, employing pairwise accuracy as the meta-evaluation metric. Assuming an 80% agreement rate with human judgments, their findings align closely with ours for pretrained metrics but not for metrics such as BLEU or ChrF. For instance, COMET-22 requires a score difference of 0.0056 to achieve 80% accuracy with humans, compared to our range of 0.0043–0.0055. Similarly, CometKiwi requires a

Figure 4: Log p-value of two-sided paired t-test on MQM scores ($p_{mqm}$) against the metric (left: BLEU, right: COMET-22) score difference for each system pair in en→de. The red line is the isotonic regression fit to all data points, representing $Pr(p_{mqm} < 0.05|\Delta M)$. Note: for readability, values of $p_{mqm}$ are rounded up to 0.0001 when they are less than 0.0001.

| Metric | en→de | | en→es | | ja→zh | |
|---|---|---|---|---|---|---|
| | min $\Delta M$ | c.v. precision | min $\Delta M$ | c.v. precision | min $\Delta M$ | c.v. precision |
| BERTSCORE | 0.0099 | [50-100%] | 0.018 | [50-100%] | 0.013 | [64-100%] |
| BLCOM_1 | 0.022 | [75-100%] | 0.034 | [50-100%] | 0.021 | [62-100%] |
| BLEU | 5.4 | [67-100%] | 11 | [0-100%] | 1.4 | [50-100%] |
| BLEURT-20 | 0.021 | [62-100%] | 0.014 | [60-100%] | 0.029 | [80-100%] |
| BRIGHT-QE | 0.018 | [20-100%] | 0.049 | [50-100%] | 0.061 | [62-100%] |
| CHRF | 3.0 | [67-100%] | 2.1 | [57-100%] | 3.5 | [78-100%] |
| CHRFS | 0.023 | [50-100%] | 0.043 | [50-100%] | 0.021 | [60-100%] |
| COMET-22 | 0.018 | [50-100%] | 0.017 | [60-100%] | 0.021 | [60-100%] |
| COMETKIWI | 0.024 | [17-100%] | 0.027 | [33-100%] | 0.050 | [67-100%] |
| DAMONMONLI | 0.84 | [27-100%] | 0.064 | [50-100%] | 0.51 | [88-100%] |
| GEMBA_ESA | 4.5 | [70-100%] | 1.5 | [67-100%] | 4.8 | [86-100%] |
| MEE4 | 0.019 | [25-100%] | 0.028 | [33-100%] | 0.019 | [55-100%] |
| metametrics_mt_mqm_hybrid_kendall | 0.029 | [53-100%] | 0.066 | [60-100%] | 0.066 | [70-100%] |
| metametrics_mt_mqm_qe_kendall.seg.s | 0.016 | [14-100%] | 0.025 | [50-100%] | 0.031 | [67-100%] |
| METRICX-24-HYBRID | 0.52 | [73-100%] | 0.95 | [62-100%] | 0.60 | [75-100%] |
| METRICX-24-HYBRID-QE | 0.44 | [62-100%] | 0.39 | [67-100%] | 0.63 | [78-100%] |
| PRISMREFMEDIUM | 0.073 | [67-100%] | 0.12 | [50-100%] | 0.14 | [56-100%] |
| PRISMREFSMALL | 0.10 | [67-100%] | 0.15 | [50-100%] | 0.15 | [56-100%] |
| SENTINEL-CAND-MQM | 0.066 | [50-100%] | 0.13 | [50-100%] | 0.088 | [55-100%] |
| SENTINEL-REF-MQM | — | — | — | — | — | — |
| SENTINEL-SRC-MQM | — | — | — | — | — | — |
| SPBLEU | 4.3 | [50-100%] | 9.1 | [0-100%] | 4.0 | [75-100%] |
| XCOMET | 0.022 | [53-100%] | 0.025 | [67-100%] | 0.046 | [78-100%] |
| XCOMET-QE | 0.013 | [50-100%] | 0.029 | [50-100%] | 0.062 | [67-100%] |
| XLSIMMQM | 0.018 | [100-100%] | 0.0012 | [57-100%] | 0.004 | [43-100%] |
| YISI-1 | 0.0063 | [60-100%] | 0.0098 | [56-100%] | 0.012 | [75-100%] |

Table 10: Minimum $\Delta M$ when $Pr(p_{mqm} < 0.05|\Delta M) = 0.8$ for each metric in different language pairs round to 2 significant figures, and the range of precision for the isotonic regression model in leave-one-system-out cross validation.

difference of 0.0053, while our results range from 0.0037 to 0.0056. Conversely, for BLEU, their analysis suggests an expected improvement of 2.34 BLEU points for 80% agreement, whereas our analysis indicates a need for an improvement of 0.89–0.97 BLEU points. However, it is important to note that we are comparing distinct metrics, and that confidence levels are not directly comparable to agreement rates.

We have to emphasize again that our result should *NOT* be interpreted as evidence to forego significance tests or appropriate human evaluation. Instead, we are only providing assistance to build an intuition on the meaning of the scores provided by the new metrics to encourage the transition away from lexical metrics towards more recent and stronger metrics.

Figure 5: Log p-value of significance test with bootstrap resampling ($p_M$) on system-level metric scores against each metric (left: BLEU, right: COMET-22) score difference for each system pair in en→de. The red line is the isotonic regression fit to all data points, representing $Pr(p_M < 0.05|\Delta M)$. Note: for readability, values of $p_M$ are rounded up to 0.0001 when they are less than 0.0001.

| | en→de | | en→es | | ja→zh | |
|---|---|---|---|---|---|---|
| Metric | min $\Delta M$ | c.v. precision | min $\Delta M$ | c.v. precision | min $\Delta M$ | c.v. precision |
| BERTSCORE | 0.0028 | [92-100%] | 0.0028 | [100-100%] | 0.0044 | [100-100%] |
| BLCOM_1 | 0.0039 | [100-100%] | 0.0055 | [100-100%] | 0.0044 | [100-100%] |
| BLEU | 0.97 | [100-100%] | 0.93 | [100-100%] | 0.89 | [91-100%] |
| BLEURT-20 | 0.0056 | [96-100%] | 0.0053 | [94-100%] | 0.0068 | [95-100%] |
| BRIGHT-QE | 0.0041 | [89-100%] | 0.0078 | [94-100%] | 0.024 | [95-100%] |
| CHRF | 0.83 | [96-100%] | 0.77 | [94-100%] | 0.89 | [100-100%] |
| CHRFS | 0.0051 | [91-100%] | 0.0054 | [95-100%] | 0.0055 | [95-100%] |
| COMET-22 | 0.0043 | [96-100%] | 0.0055 | [86-100%] | 0.0046 | [95-100%] |
| COMETKIWI | 0.0037 | [100-100%] | 0.0048 | [82-100%] | 0.0056 | [100-100%] |
| DAMONMONLI | 0.20 | [94-100%] | 0.17 | [82-100%] | 0.41 | [90-100%] |
| GEMBA_ESA | 0.82 | [92-100%] | 0.85 | [91-100%] | 1.4 | [100-100%] |
| MEE4 | 0.0042 | [95-100%] | 0.0051 | [86-100%] | 0.0057 | [95-100%] |
| metametrics_mt_mqm_hybrid_kendall | 0.0067 | [92-100%] | 0.0081 | [89-100%] | 0.013 | [90-100%] |
| metametrics_mt_mqm_qe_kendall.seg.s | 0.0038 | [89-100%] | 0.0050 | [80-100%] | 0.0089 | [95-100%] |
| METRICX-24-HYBRID | 0.11 | [100-100%] | 0.15 | [100-100%] | 0.14 | [95-100%] |
| METRICX-24-HYBRID-QE | 0.087 | [90-100%] | 0.14 | [100-100%] | 0.12 | [100-100%] |
| SENTINEL-CAND-MQM | 0.011 | [96-100%] | 0.013 | [95-100%] | 0.030 | [95-100%] |
| SENTINEL-REF-MQM | — | — | — | — | — | — |
| SENTINEL-SRC-MQM | — | — | — | — | — | — |
| SPBLEU | 0.96 | [96-100%] | 1.1 | [95-100%] | 1.0 | [100-100%] |
| PRISMREFMEDIUM | 0.019 | [95-100%] | 0.02 | [100-100%] | 0.036 | [90-100%] |
| PRISMREFSMALL | 0.023 | [96-100%] | 0.022 | [100-100%] | 0.042 | [95-100%] |
| XCOMET | 0.0051 | [100-100%] | 0.0065 | [86-100%] | 0.010 | [95-100%] |
| XCOMET-QE | 0.0044 | [96-100%] | 0.0058 | [94-100%] | 0.0099 | [100-100%] |
| XLSIMMQM | 0.0036 | [82-100%] | 0.0013 | [90-100%] | 0.0019 | [79-100%] |
| YISI-1 | 0.0010 | [91-100%] | 0.0014 | [90-100%] | 0.0051 | [100-100%] |

Table 11: Minimum $\Delta M$ when $Pr(p_M < 0.05|\Delta M) = 0.8$ for each metric in different language pairs round to 2 significant figures, and the range of precision for the isotonic regression model in leave-one-system-out cross validation.

## 8 ESA Human Evaluation

In addition to our MQM annotations and as a contrastive evaluation to cover more language pairs, we look into the performance of metrics when compared to the human evaluation campaign conducted by the WMT24 General MT Shared Task (Kocmi et al., 2024a), which ran human evaluation for nine language pairs.

In contrast to previous years, WMT24 redefined their human evaluation process and developed a new method called Error Span Analysis (ESA, Kocmi et al. (2024c)), a method that simplifies MQM by asking annotators only to mark error spans and classify them either as minor or major severity. In addition to that, the annotator is asked to mark the whole segment with a score of 0–100 in the SQM fashion. As Kocmi et al. (2024c) claim, the method is cheaper than MQM to annotate, yet

it produces closer human judgment to MQM annotations than the formerly used DA+SQM (Kocmi et al., 2023) due to being less affected by fluency.

We present system-level accuracy results for both MQM and ESA in Table 15. There are many factors that could affect the ranking. Apart from using a different human annotation protocol, MQM compares 3 language pairs whereas ESA compares 9 language pairs, containing also two low-resource pairs: Czech→Ukrainian and English→Icelandic. There is an overlap of only one language pair between the two: English→Spanish.

Most of the metrics have a similar ranking for both MQM and ESA; however, there are two metrics with largely different rankings: GEMBA_ESA and metametrics_mt_mqm_qe_kendall.seg.s, whose rankings are significantly lower under ESA than for MQM. The likely explanation for GEMBA_ESA is that ESA doesn't produce ties, in contrast to MQM, whereas GEMBA_ESA produces them regularly. As for the latter metric, we don't see any clear pattern except for having low performance for Czech→Ukrainian.

## 9 Challenge Sets Sub-task

For the third year, the Metrics Shared Task included a sub-task involving challenge sets. This sub-task is inspired by the *Build it or break it: The Language Edition* shared task (Ettinger et al., 2017) which aimed at testing the generalizability of NLP systems beyond the distributions of their training data. Whereas the standard evaluation of the shared task is conducted on test sets containing generic text from real-world content, the challenge set evaluation is based on test sets designed with the aim of revealing the abilities or the weaknesses of the metrics or evaluating particular translation phenomena. In order to shed light on different perspectives on evaluation, the sub-task takes place in a decentralized manner, since contrary to the main metric task, the test sets are not provided by the organizers but by different research teams, who are also responsible for analysing and presenting the results.

This subtask is made of three consecutive phases; 1) the *Breaking Round*, 2) the *Scoring Round* and 3) the *Analysis Round*:

1. In the *Breaking Round*, every challenge set participant (*Breaker*) submits their challenge set $S$ composed of examples for different phenomena, where every example $(s, t, r) \in S$ contains one source sentence $s$, one translation hypothesis $t$ and one reference $r$.

2. In the *Scoring Round*, The metrics participants from the main task (the *Builders*) are asked to score with their metrics the translations in the given test set. Also, in this phase, the metrics task organizers score all data with the baseline metrics.

3. Finally, after having gathered all metric scores, the organizers return the respective scored translations to the *Breakers* for the *Analysis round*, where they employ their own evaluation for the performance of the metrics with regard to the phenomena they intended to test.

This year there were 4 submissions, covering a wide range of phenomena and 23 different language pairs, which supersede the official language pairs of the Metrics Shared Task. An overview of the submitted challenge sets can be seen in Table 12. A short description of every submission follows:

**AfriMTE Challenge Set** The AFRIMTE challenge set (Wang et al., 2024b) aims to evaluate the capabilities of metrics for machine translation on low-resource languages, primarily assessing cross-lingual transfer learning and generalization across a wide range of under-resourced African languages. The challenge set concentrates on the subsets of the FLORES-200 dataset (NLLB-Team et al., 2022) and covers 13 language pairs. Specifically, there are Darija-French, English-Egyptian Arabic, English-French, English-Hausa, English-Igbo, English-Kikuyu, English-Luo, English-Somali, English-Swahili, English-Twi, English-isiXhosa, English-Yoruba, and Yoruba-English. Originally, AFRIMTE (Wang et al., 2024a) provides both fine-grained word-level error annotations and sentence-level Direct Assessment scoring for translation adequacy and fluency. For this year's challenge set sub-task, we utilize the translation adequacy test set from AFRIMTE as the African Challenge set to evaluate the sentence-level scoring performance of metrics. The analysis of the task submissions (Wang et al., 2024b) has yielded several insights. First, language-specific adaptation, cross-lingual transfer learning, and larger language model sizes significantly enhance metric performance. Second, moderately-sized supervised models can attain robust performance when augmented with language adaptation techniques tailored to

| Challenge Set | Directions | Phenomena | Items | Citation | Link (https://github.com/) |
|---|---|---|---|---|---|
| AfriMTE | 13 | African languages | 2,815 | Wang et al. (2024b) | masakhane-io/africomet |
| BioMQM | 11 | biomedical domain | 4,641 | Zouhar et al. (2024) | thompsonb/bio-mqm-dataset |
| DFKI | 2 | linguistic phenomena | 137,000 | Avramidis et al. (2024) | DFKI-NLP/mt-testsuite |
| MSLC24 | 3 | low quality MT | 964 | Knowles et al. (2024) | nrc-cnrc/MSLC |

Table 12: Overview of the participation at the metrics challenge sets sub-task.

low-resource African languages during pretraining. Last, submissions demonstrate promising outcomes for language pairs such as Darija-French, English-Egyptian Arabic, and English-Swahili. However, considerable challenges remain for extremely low-resource languages like English-Luo and English-Twi, underscoring critical areas for future research and improvement in machine translation metrics for African languages.

**BioMQM** Recent work (Zouhar et al., 2024) has compared trained versus untrained metric performance on the WMT domains compared to the biomedical domain and shown that trained metrics appear to be over-fitting on the domains used in the WMT Metrics Shared Tasks. This is likely due to trained metrics using prior WMT metrics datasets, and then being evaluated on very similar data in the latest WMT Metrics Shared Task. Zouhar et al. (2024) released a biomedical dataset (BioMQM) consisting of source sentences and translations from Yeganova et al. (2021) along with new translations and MQM annotations. We produce scores on the BioMQM for the latest metrics (all those submitted to this Metrics Shared Task, plus the baseline metrics) and release them for future analysis.[13]

**DFKI Challenge Set** This year's submission by DFKI (Avramidis et al., 2024) expands the linguistically motivated challenge set of previous years (Avramidis et al., 2023; Avramidis and Macketanz, 2022), including 137,000 items in overall, extracted from 100 MT systems for the two language directions (en→de, en→ru), covering more than 100 linguistically-motivated phenomena organized in 14 linguistic categories. The metrics with the statistically significant best performance with regard to our linguistically motivated analysis are METRICX-24-HYBRID and METRICX-24 for en→de and METRICX-24 for en→ru, whereas METAMETRICS and XCOMET are in the next rank-

ing positions in both language pairs. Metrics are more accurate in detecting linguistic errors among LLM translations than in translations based on the encoder-decoder NMT architecture. Some of the most difficult phenomena for the metrics to score are the transitive past progressive, the multiple connectors, the ditransitive simple future I for en→de and pseudogapping, contact clause and cleft sentences for en→ru. The LLM-based metric GEMBA, despite the overall low performance, has the best performance on scoring German negation errors.

**MSLC24 Challenge Set** Building on the Metric Score Landscape Challenge (MSLC23; Lo et al., 2023b), which aims to provide a view of metric performance on a broader range of MT quality, MSLC24 includes a collection of low- to medium-quality MT systems' output on the news portion of the WMT24 General MT Shared Task test set, as well as some specific phenomena that may result in unexpected behaviors from some metrics, such as empty strings in source/reference/hypothesis, wrong/mixed language output and different language variants. MSLC24 focuses on three language pairs (English→German, English→Spanish and Japanese→Chinese). The authors also submit the top system in this challenge set to the General Translation task in order to obtain human evaluation. Together with the high quality systems by other participants submitted to the General MT Shared Task, this enables better interpretation of metric scores across a range of different levels of translation quality and analyse metric characteristics beyond just correlation. The results of MSLC24 highlight the importance of examining real-word corner cases and issues of reproducibility in order to more responsibly introduce new metrics to the research community.

## 10 Conclusion

This paper summarizes the results of the WMT24 shared task on automated machine translation evaluation, the Metrics Shared Task. We presented an extensive analysis on how well metrics perform on

---
[13]https://github.com/thompsonb/
bio-mqm-dataset/tree/main/data/WMT24_
Metrics_ChallengeSet

our three main language pairs: English→German, English→Spanish and Japanese→Chinese. The results, based on 6 different tasks, confirm the superiority of neural-based learned metrics over overlap-based metrics like BLEU, SPBLEU or CHRF. These results are confirmed with ESA human judgement. Overall, we did not find any issues for neural fine-tuned metrics when evaluating LLM-based translations. In addition, we continued the challenge set subtask, where participants had to create contrastive test suites for evaluating metrics' ability to capture and penalise specific types of translation errors.

## 11 Ethical Considerations

## 12 Acknowledgments

## References

David Anugraha, Garry Kuwanto, Lucky Susanto, Derry Tanti Wijaya, and Genta Indra Winata. 2024. Metametrics-MT: Tuning machine translation metametrics via human preference calibration. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Eleftherios Avramidis and Vivien Macketanz. 2022. Linguistically motivated evaluation of machine translation metrics based on a challenge set. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 514–529, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Eleftherios Avramidis, Shushen Manakhimova, Vivien Macketanz, and Sebastian Moeller. 2024. Machine translation metrics are better in evaluating linguistic errors on llms than on encoder-decoder systems. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, Florida, USA. Association for Computational Linguistics.

Eleftherios Avramidis, Shushen Manakhimova, Vivien Macketanz, and Sebastian Möller. 2023. Challenging the state-of-the-art machine translation metrics from a linguistic perspective. In *Proceedings of the Eighth Conference on Machine Translation*, pages 713–729, Singapore. Association for Computational Linguistics.

Yanran Chen and Steffen Eger. 2023. MENLI: Robust evaluation metrics from natural language inference. *Transactions of the Association for Computational Linguistics*, 11:804–825.

Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. A statistical analysis of summarization evaluation metrics using resampling methods. *arXiv preprint arXiv:2104.00054*.

Daniel Deutsch, George Foster, and Markus Freitag. 2023. Ties Matter: Meta-Evaluating Modern Metrics with Pairwise Accuracy and Tie Calibration. pages 12914–12929.

Allyson Ettinger, Sudha Rao, Hal Daumé III, and Emily M. Bender. 2017. Towards linguistically generalizable NLP systems: A workshop and shared task. In *Proceedings of the First Workshop on Building Linguistically Generalizable NLP Systems*, pages 1–10, Copenhagen, Denmark. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Mara Finkelstein and Markus Freitag. 2024. MBR and QE finetuning: Training-time distillation of the best and most expensive decoding methods. In *The Twelfth International Conference on Learning Representations*.

Mara Finkelstein, David Vilar, and Markus Freitag. 2024. Introducing the newspalm mbr and qe dataset: Llm-generated high-quality parallel data outperforms traditional web-crawled data. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Marina Fomicheva, Shuo Sun, Erick Fonseca, Chrysoula Zerva, Frédéric Blain, Vishrav Chaudhary, Francisco Guzmán, Nina Lopatina, Lucia Specia, and André F. T. Martins. 2022. MLQE-PE: A multilingual quality estimation and post-editing dataset. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4963–4974, Marseille, France. European Language Resources Association.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High Quality Rather than High Model Probability: Minimum Bayes Risk Decoding with Neural Metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1180–1189, Lille, France. PMLR.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xCOMET: Transparent Machine Translation Evaluation through Fine-grained Error Detection. *arXiv preprint arXiv:2310.10482*.

Juraj Juraska, Daniel Deutsch, Mara Finkelstein, and Markus Freitag. 2024. Metricx-24: The google submission to the WMT 2024 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Rebecca Knowles, Samuel Larkin, and Chi-kiu Lo. 2024. MSLC24: Further challenges for metrics on a wide landscape of translation quality. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024a. Findings of the WMT24 general machine translation shared task: the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Masaaki Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, Mariya Shmatova, and Jun Suzuki. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.

Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024b. Navigating the metrics maze: Reconciling score magnitudes and accuracies. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1999–2014, Bangkok, Thailand. Association for Computational Linguistics.

Tom Kocmi, Vilém Zouhar, Eleftherios Avramidis, Roman Grundkiewicz, Marzena Karpinska, Maja Popović, Mrinmaya Sachan, and Mariya Shmatova. 2024c. Error span annotation: A balanced approach for human evaluation of machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Geza Kovacs, Daniel Deutsch, and Markus Freitag. 2024. Mitigating metric bias in minimum bayes risk decoding. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

AI @ Meta Llama Team. 2024. The llama 3 herd of models.

Chi-kiu Lo. 2019. YiSi - a unified semantic MT quality evaluation and estimation metric for languages with different levels of available resources. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 507–513, Florence, Italy. Association for Computational Linguistics.

Chi-kiu Lo, Rebecca Knowles, and Cyril Goutte. 2023a. Beyond correlation: Making sense of the score differences of new MT evaluation metrics. In *Proceedings of Machine Translation Summit XIX, Vol. 1: Research Track*, pages 186–199, Macau SAR, China. Asia-Pacific Association for Machine Translation.

Chi-kiu Lo, Samuel Larkin, and Rebecca Knowles. 2023b. Metric score landscape challenge (MSLC23): Understanding metrics' performance on a wider landscape of translation quality. In *Proceedings of the Eighth Conference on Machine Translation*, pages 776–799, Singapore. Association for Computational Linguistics.

Arle Lommel, Hans Uszkoreit, and Aljoscha Burchardt. 2014. Multidimensional Quality Metrics (MQM) : A Framework for Declaring and Describing Translation Quality Metrics. *Tradumàtica*, pages 0455–463.

Benjamin Marie. 2022. Yes, we need statistical significance testing. towardsai.net https://pub.towardsai.net/yes-we-need-statistical-significance-testing-927a8d21f9f0.

Ananya Mukherjee and Manish Shrivastava. 2023a. MEE4 and XLsim : IIIT HYD's submissions' for WMT23 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 800–805, Singapore. Association for Computational Linguistics.

Ananya Mukherjee and Manish Shrivastava. 2023b. MEE4 and XLsim: IIIT HYD's Submissions for WMT23 Metrics Shared Task. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.

Ananya Mukherjee and Manish Shrivastava. 2024. chrf-s: Semantics is all you need. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

NLLB-Team, Marta Ruiz Costa-jussà, James Cross, Onur cCelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Alison Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon L. Spruit, C. Tran, Pierre Yves Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzm'an, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *ArXiv*, abs/2207.04672.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume

Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No Language Left Behind: Scaling Human-Centered Machine Translation. *arXiv preprint arXiv:2207.04672*.

Eric W Noreen. 1989. Computer intensive methods for hypothesis testing: An introduction. *Wiley, New York*, 19:21.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Stefano Perrella, Lorenzo Proietti, Alessandro Scirè, Edoardo Barba, and Roberto Navigli. 2024. Guardians of the machine translation meta-evaluation: Sentinel metrics fall in! In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16216–16244, Bangkok, Thailand. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. Learning compact metrics for MT. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 751–762, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C.

de Souza, and André F. T. Martins. 2023. Scaling up COMETKIWI: Unbabel-IST 2023 Submission for the Quality Estimation Shared Task. In *Proceedings of the eighth Conference on Machine Translation*, Singapore. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Parker Riley, Daniel Deutsch, George Foster, Viresh Ratnakar, Ali Dabirmoghaddam, and Markus Freitag. 2024. Finding replicable human evaluations via stable ranking probability. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4908–4919, Mexico City, Mexico. Association for Computational Linguistics.

T. Robertson, F.T. Wright, and R. Dykstra. 1988. *Order Restricted Statistical Inference*. Probability and Statistics Series. Wiley.

Davor Runje and Sharath M Shankaranarayana. 2023. Constrained monotonic neural networks. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 29338–29353. PMLR.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Brian Thompson, Nitika Mathur, Daniel Deutsch, and Huda Khayrallah. 2024. Improving statistical significance in human evaluation of automatic metrics via soft pairwise accuracy.

Brian Thompson and Matt Post. 2020a. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 90–121, Online. Association for Computational Linguistics.

Brian Thompson and Matt Post. 2020b. Paraphrase generation as zero-shot multilingual translation: Disentangling semantic similarity from lexical and syntactic diversity. In *Proceedings of the Fifth Conference on Machine Translation*, pages 561–570, Online. Association for Computational Linguistics.

Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 118–128, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Jiayi Wang, David Adelani, Sweta Agrawal, Marek Masiak, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Xuanli He, Sofia Bourhim, Andiswa Bukula, Muhidin Mohamed, Temitayo Olatoye, Tosin Adewumi, Hamam Mokayed, Christine Mwase, Wangui Kimotho, Foutse Yuehgoh, Anuoluwapo Aremu, Jessica Ojo, Shamsuddeen Muhammad, Salomey Osei, Abdul-Hakeem Omotayo, Chiamaka Chukwuneke, Perez Ogayo, Oumaima Hourrane, Salma El Anigri, Lolwethu Ndolela, Thabiso Mangwana, Shafie Mohamed, Hassan Ayinde, Oluwabusayo Awoyomi, Lama Alkhaled, Sana Alazzawi, Naome Etori, Millicent Ochieng, Clemencia Siro, Njoroge Kiragu, Eric Muchiri, Wangari Kimotho, Toadoum Sari Sakayo, Lyse Naomi Wamba, Daud Abolade, Simbiat Ajao, Iyanuoluwa Shode, Ricky Macharm, Ruqayya Iro, Saheed Abdullahi, Stephen Moore, Bernard Opoku, Zainab Akinjobi, Abeeb Afolabi, Nnaemeka Obiefuna, Onyekachi Ogbu, Sam Ochieng', Verrah Otiende, Chinedu Mbonu, Yao Lu, and Pontus Stenetorp. 2024a. AfriMTE and AfriCOMET: Enhancing COMET to embrace under-resourced African languages. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5997–6023, Mexico City, Mexico. Association for Computational Linguistics.

Jiayi Wang, David Ifeoluwa Adelani, and Pontus Stenetorp. 2024b. Evaluating WMT 2024 metrics shared task submissions on afrimte (the african challenge set). In *Proceedings of the Ninth Conference on Machine Translation*, Miami, Florida, USA. Association for Computational Linguistics.

Johnny Wei, Tom Kocmi, and Christian Federmann. 2022. Searching for a higher power in the human evaluation of MT. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 129–139, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Genta Indra Winata, David Anugraha, Lucky Susanto, Garry Kuwanto, and Derry Tanti Wijaya. 2024. Metametrics: Calibrating metrics for generation tasks using human preferences. *arXiv preprint arXiv:2410.02381*.

Lana Yeganova, Dina Wiemann, Mariana Neves, Federica Vezzani, Amy Siu, Inigo Jauregi Unanue, Maite Oronoz, Nancy Mah, Aurélie Névéol, David Martinez, Rachel Bawden, Giorgio Maria Di Nunzio, Roland Roller, Philippe Thomas, Cristian Grozea, Olatz Perez-de Viñaspre, Maika Vicente Navarro, and Antonio Jimeno Yepes. 2021. Findings of the WMT 2021 biomedical translation shared task: Summaries

of animal experiments as new test set. In *Proceedings of the Sixth Conference on Machine Translation*, pages 664–683, Online. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Vilém Zouhar, Shuoyang Ding, Anna Currey, Tatyana Badeka, Jenyuan Wang, and Brian Thompson. 2024. Fine-tuned machine translation metrics struggle in unseen domains. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 488–500, Bangkok, Thailand. Association for Computational Linguistics.

# A   Correlations with MQM for all metrics

Table 13 contains the results for all metrics (including contrastive submissions) on the 6 standard tasks described in Table 8.

| Metric | avg-corr | | en-de sys SPA task1 | | en-de seg acc$^*_{eq}$ task2 | | en-es sys SPA task3 | | en-es seg acc$^*_{eq}$ task4 | | ja-zh sys SPA task5 | | ja-zh seg acc$^*_{eq}$ task6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *MetricX-24* | **1** | **0.725** | 2 | 0.873 | 2 | 0.534 | 2 | 0.789 | 3 | 0.685 | **1** | **0.921** | 2 | 0.547 |
| MetaMetrics-MT | **1** | **0.725** | 2 | 0.882 | **1** | **0.542** | 2 | 0.805 | 2 | 0.686 | 3 | 0.872 | **1** | **0.561** |
| *metametrics_mt_mqm_kendall* | **1** | **0.724** | 2 | 0.882 | **1** | **0.542** | 2 | 0.804 | 2 | 0.686 | 3 | 0.871 | **1** | **0.561** |
| *metametrics_mt_mqm_same_source_targ* | 2 | 0.723 | **1** | **0.883** | **1** | **0.542** | 2 | 0.803 | 2 | 0.686 | 3 | 0.874 | 2 | 0.550 |
| MetricX-24-Hybrid | 2 | 0.720 | 2 | 0.873 | 2 | 0.532 | 2 | 0.796 | 3 | 0.685 | 2 | 0.895 | 3 | 0.539 |
| XCOMET | 2 | 0.719 | **1** | **0.906** | 3 | 0.530 | 2 | 0.788 | **1** | **0.688** | 2 | 0.890 | 7 | 0.510 |
| MetricX-24-Hybrid-QE* | 3 | 0.714 | 2 | 0.880 | 4 | 0.526 | 2 | 0.790 | 4 | 0.685 | 3 | 0.875 | 4 | 0.530 |
| gemba_esa* | 3 | 0.712 | 4 | 0.793 | 6 | 0.507 | **1** | **0.838** | 5 | 0.683 | **1** | **0.909** | 3 | 0.539 |
| *MetricX-24-QE** | 3 | 0.710 | 2 | 0.880 | 3 | 0.528 | 3 | 0.772 | 3 | 0.685 | 3 | 0.875 | 5 | 0.522 |
| *CometKiwi-XXL** | 3 | 0.703 | 3 | 0.839 | 9 | 0.481 | **1** | **0.843** | 8 | 0.680 | 2 | 0.881 | 8 | 0.494 |
| XCOMET-QE* | 4 | 0.695 | **1** | **0.890** | 5 | 0.520 | 2 | 0.801 | 2 | 0.687 | 5 | 0.809 | 12 | 0.463 |
| COMET-22 | 4 | 0.689 | 2 | 0.877 | 9 | 0.482 | 2 | 0.782 | 5 | 0.683 | 5 | 0.815 | 8 | 0.496 |
| *metametrics_mt_mqm_qe_same_source_t** | 4 | 0.688 | 2 | 0.860 | 7 | 0.497 | 4 | 0.709 | 2 | 0.686 | 4 | 0.853 | 5 | 0.524 |
| BLEURT-20 | 4 | 0.686 | 2 | 0.879 | 8 | 0.486 | 4 | 0.696 | 6 | 0.681 | 2 | 0.888 | 10 | 0.484 |
| MetaMetrics-MT-QE* | 5 | 0.685 | 2 | 0.859 | 7 | 0.497 | 4 | 0.710 | 2 | 0.686 | 5 | 0.839 | 6 | 0.516 |
| bright-qe* | 5 | 0.682 | 3 | 0.817 | 7 | 0.500 | 2 | 0.794 | **1** | **0.689** | 5 | 0.806 | 10 | 0.484 |
| BLCOM_1 | 6 | 0.664 | 3 | 0.842 | 11 | 0.455 | 4 | 0.679 | 6 | 0.681 | 4 | 0.840 | 9 | 0.488 |
| sentinel-cand-mqm* | 7 | 0.649 | 3 | 0.820 | 5 | 0.517 | 2 | 0.786 | 4 | 0.683 | 9 | 0.609 | 10 | 0.481 |
| PrismRefMedium | 7 | 0.646 | 4 | 0.776 | 15 | 0.434 | 4 | 0.651 | 8 | 0.680 | 3 | 0.872 | 12 | 0.462 |
| PrismRefSmall | 7 | 0.643 | 4 | 0.774 | 15 | 0.433 | 5 | 0.635 | 8 | 0.680 | 3 | 0.874 | 13 | 0.457 |
| CometKiwi* | 7 | 0.640 | 5 | 0.731 | 10 | 0.467 | 4 | 0.695 | 4 | 0.684 | 6 | 0.775 | 9 | 0.490 |
| damonmonli | 7 | 0.635 | 5 | 0.695 | 13 | 0.443 | 5 | 0.607 | 6 | 0.682 | **1** | **0.912** | 11 | 0.472 |
| YiSi-1 | 8 | 0.630 | 4 | 0.758 | 14 | 0.436 | 5 | 0.610 | 7 | 0.681 | 5 | 0.836 | 13 | 0.458 |
| *monmonli* | 8 | 0.624 | 5 | 0.681 | 14 | 0.437 | 5 | 0.583 | 7 | 0.681 | 2 | 0.891 | 11 | 0.470 |
| BERTScore | 9 | 0.617 | 4 | 0.749 | 15 | 0.435 | 5 | 0.585 | 6 | 0.682 | 6 | 0.798 | 14 | 0.451 |
| MEE4 | 9 | 0.609 | 5 | 0.731 | 14 | 0.437 | 7 | 0.498 | 4 | 0.683 | 3 | 0.856 | 15 | 0.446 |
| chrF | 10 | 0.607 | 4 | 0.751 | 17 | 0.431 | 5 | 0.579 | 9 | 0.680 | 7 | 0.765 | 18 | 0.436 |
| chrfS | 10 | 0.606 | 4 | 0.742 | 15 | 0.434 | 6 | 0.549 | 6 | 0.682 | 6 | 0.788 | 16 | 0.444 |
| spBLEU | 11 | 0.593 | 4 | 0.741 | 19 | 0.431 | 6 | 0.524 | 8 | 0.680 | 8 | 0.745 | 18 | 0.436 |
| BLEU | 11 | 0.589 | 4 | 0.736 | 18 | 0.431 | 7 | 0.513 | 9 | 0.680 | 8 | 0.739 | 19 | 0.435 |
| *BLCOM* | 12 | 0.537 | 6 | 0.619 | 16 | 0.433 | 3 | 0.730 | 8 | 0.680 | 10 | 0.325 | 19 | 0.435 |
| sentinel-ref-mqm | 12 | 0.523 | 6 | 0.495 | 20 | 0.429 | 6 | 0.514 | 9 | 0.680 | 9 | 0.583 | 19 | 0.435 |
| sentinel-src-mqm* | 12 | 0.522 | 6 | 0.496 | 20 | 0.429 | 7 | 0.512 | 9 | 0.680 | 9 | 0.581 | 19 | 0.435 |
| *XLsimDA** | 12 | 0.514 | 6 | 0.614 | 12 | 0.450 | 8 | 0.357 | 7 | 0.681 | 9 | 0.548 | 17 | 0.438 |
| XLsimMqm* | 12 | 0.514 | 6 | 0.614 | 12 | 0.450 | 8 | 0.357 | 7 | 0.681 | 9 | 0.547 | 17 | 0.438 |

Table 13: Soft pairwise accuracy (SPA) and acc$^*_{eq}$ results for all metrics for main language pairs. See §5 for descriptions of SPA and acc$^*_{eq}$. Rows are sorted by the overall average correlation across all 6 tasks (leftmost column). Starred metrics are reference-free, underlined metrics are baselines, and italicized metrics are contrastive submissions.

| Metric | | avg corr | p-values |
|---|---|---|---|
| MetaMetrics-MT | 1 | 0.725 | . 19 07 01 01 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| MetricX-24-Hybrid | 1 | 0.721 | . . 31 01 01 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| XCOMET | 1 | 0.719 | . . . 15 10 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| MetricX-24-Hybrid-QE* | 2 | 0.714 | . . . . 36 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| gemba_esa* | 2 | 0.711 | . . . . . 01 00 01 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| XCOMET-QE* | 3 | 0.695 | . . . . . . 22 14 14 02 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| COMET-22 | 3 | 0.688 | . . . . . . . 20 34 20 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| BLEURT-20 | 3 | 0.686 | . . . . . . . . 43 28 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| MetaMetrics-MT-QE* | 3 | 0.684 | . . . . . . . . . 34 02 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| bright-qe* | 4 | 0.681 | . . . . . . . . . . 06 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |
| BLCOM_1 | 4 | 0.664 | . . . . . . . . . . . 04 02 00 00 01 00 00 00 00 00 00 00 00 00 00 00 00 |
| sentinel-cand-mqm* | 5 | 0.650 | . . . . . . . . . . . . 41 25 21 13 06 01 00 00 00 00 00 00 00 00 00 00 |
| PrismRefMedium | 5 | 0.646 | . . . . . . . . . . . . . 11 35 19 01 00 00 00 00 00 00 00 00 00 00 00 |
| PrismRefSmall | 5 | 0.642 | . . . . . . . . . . . . . . 43 30 03 00 00 00 00 00 00 00 00 00 00 00 |
| CometKiwi* | 5 | 0.640 | . . . . . . . . . . . . . . . 33 17 03 00 01 00 00 00 00 00 00 00 00 |
| damonmonli | 5 | 0.635 | . . . . . . . . . . . . . . . . 34 06 01 02 01 00 00 00 00 00 00 00 |
| YiSi-1 | 6 | 0.630 | . . . . . . . . . . . . . . . . . 01 00 00 00 00 00 00 00 00 00 00 |
| BERTScore | 7 | 0.617 | . . . . . . . . . . . . . . . . . . 14 04 03 00 00 00 00 00 00 |
| MEE4 | 7 | 0.609 | . . . . . . . . . . . . . . . . . . . 41 26 00 01 00 00 00 00 |
| chrF | 8 | 0.608 | . . . . . . . . . . . . . . . . . . . . 36 00 00 00 00 00 00 |
| chrfS | 8 | 0.606 | . . . . . . . . . . . . . . . . . . . . . 00 01 00 00 00 00 |
| spBLEU | 9 | 0.593 | . . . . . . . . . . . . . . . . . . . . . . 25 00 00 00 00 |
| BLEU | 9 | 0.589 | . . . . . . . . . . . . . . . . . . . . . . . 00 00 00 00 |
| XLsimMqm* | 10 | 0.515 | . . . . . . . . . . . . . . . . . . . . . . . . 45 49 |
| sentinel-src-mqm* | 10 | 0.513 | . . . . . . . . . . . . . . . . . . . . . . . . . 53 |
| sentinel-ref-mqm | 10 | 0.513 | . . . . . . . . . . . . . . . . . . . . . . . . . . |

Table 14: Results of pairwise metric significance tests for primary submissions using permutation resampling. Each value gives the $100 \times$ estimated probability of the null hypothesis that the average correlation of the metric in the current row is $\leq$ the average correlation of the metric in the current column. Starred metrics are reference-free, and underlined metrics are baselines.

## B  Significance comparisons for main results

Table 14 contains the results of pairwise comparisons for the results in Table 1.

## C  Correlations with WMT ESA for all metrics

Table 15 shows the correlations of the metrics to the ESA scores (see Section 8 for which those scores are available). The overall ranking is sorted by the average correlation, which is the average over all tasks across all language pairs. Metrics that did not participate in all tasks do not have an average correlation, and are displayed at the end of the table.

The system-level ESA scores that were used to calculate SPA here differ slightly from those in the General MT Shared Task. Namely, the General Task calculates scores by macro-averaging over domains (each domain receives equal weight), whereas we perform a standard micro-average (each segment receives equal weight).

Table 15: Correlations of metrics to the ESA annotations that were collected as part of the General MT Shared Task. The metrics are sorted by the average correlation across all of the correlations and language pairs. Metrics in italics are contrastive submissions and underlined metrics are baselines. QE metrics are marked by an asterisk.

| Metric | avg-corr | cs-uk sys SPA task1 | cs-uk seg acc*_eq task2 | en-cs sys SPA task3 | en-cs seg acc*_eq task4 | en-es sys SPA task5 | en-es seg acc*_eq task6 | en-hi sys SPA task7 | en-hi seg acc*_eq task8 | en-is sys SPA task9 | en-is seg acc*_eq task10 | en-ja sys SPA task11 | en-ja seg acc*_eq task12 | en-ru sys SPA task13 | en-ru seg acc*_eq task14 | en-uk sys SPA task15 | en-uk seg acc*_eq task16 | en-zh sys SPA task17 | en-zh seg acc*_eq task18 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *MetricX-24* | 1 0.708 | 3 0.890 | 1 0.482 | 1 0.896 | 1 0.585 | 2 0.834 | 1 0.503 | 1 0.938 | 2 0.567 | 3 0.855 | 1 0.670 | 3 0.791 | 1 0.558 | 1 0.932 | 1 0.537 | 1 0.872 | 2 0.447 | 2 0.826 | 1 0.569 |
| *MetricX-24-Hybrid* | 1 0.706 | 3 0.884 | 1 0.483 | 1 0.886 | 2 0.582 | 1 0.846 | 1 0.496 | 1 0.953 | 1 0.571 | 3 0.847 | 2 0.661 | 3 0.793 | 1 0.557 | 1 0.939 | 1 0.536 | 1 0.880 | 2 0.443 | 2 0.808 | 1 0.568 |
| metametrics_mt_mqm_kendall | 2 0.703 | 5 0.819 | 2 0.483 | 2 0.886 | 3 0.575 | 1 0.860 | 1 0.502 | 1 0.930 | 2 0.564 | 4 0.846 | 2 0.664 | 3 0.787 | 3 0.549 | 2 0.928 | 2 0.536 | 2 0.850 | 2 0.449 | 1 0.853 | 2 0.564 |
| metametrics_mt_mqm_same_source_targ | 2 0.702 | 4 0.820 | 1 0.483 | 2 0.887 | 3 0.575 | 2 0.858 | 1 0.502 | 1 0.928 | 2 0.564 | 4 0.846 | 2 0.664 | 3 0.789 | 3 0.549 | 2 0.926 | 1 0.536 | 2 0.845 | 2 0.449 | 1 0.854 | 2 0.564 |
| MetaMetrics-MT | 3 0.702 | 4 0.821 | 1 0.483 | 2 0.885 | 3 0.575 | 1 0.860 | 1 0.502 | 1 0.927 | 2 0.564 | 4 0.846 | 6 0.664 | 3 0.786 | 3 0.549 | 2 0.928 | 1 0.536 | 2 0.849 | 2 0.449 | 2 0.852 | 2 0.564 |
| BLEURT-20 | 3 0.701 | 1 0.960 | 3 0.471 | 1 0.915 | 5 0.563 | 3 0.793 | 3 0.491 | 2 0.906 | 4 0.556 | 3 0.888 | 6 0.616 | 2 0.824 | 2 0.543 | 4 0.944 | 4 0.520 | 3 0.797 | 5 0.434 | 4 0.846 | 5 0.550 |
| XCOMET | 3 0.701 | 6 0.782 | 3 0.472 | 1 0.901 | 3 0.572 | 2 0.856 | 2 0.483 | 2 0.918 | 2 0.567 | 3 0.866 | 3 0.663 | 4 0.778 | 2 0.550 | 1 0.931 | 3 0.531 | 1 0.875 | 1 0.455 | 1 0.847 | 1 0.568 |
| COMET22 | 3 0.700 | 2 0.918 | 2 0.477 | 3 0.862 | 4 0.566 | 3 0.870 | 5 0.498 | 2 0.937 | 2 0.552 | 2 0.896 | 3 0.650 | 6 0.710 | 3 0.548 | 2 0.916 | 3 0.528 | 3 0.817 | 3 0.441 | 1 0.856 | 1 0.566 |
| MetricX-24-Hybrid-QE* | 4 0.690 | 6 0.790 | 4 0.463 | 3 0.875 | 4 0.568 | 4 0.844 | 5 0.479 | 1 0.934 | 1 0.568 | 4 0.835 | 5 0.637 | 2 0.787 | 3 0.550 | 2 0.914 | 3 0.526 | 1 0.872 | 4 0.439 | 3 0.778 | 1 0.558 |
| *MetricX-24-QE** | 4 0.688 | 5 0.818 | 5 0.459 | 2 0.873 | 3 0.572 | 4 0.838 | 4 0.481 | 2 0.904 | 1 0.569 | 4 0.841 | 4 0.643 | 3 0.756 | 2 0.553 | 3 0.908 | 3 0.527 | 2 0.849 | 3 0.450 | 3 0.789 | 2 0.561 |
| XCOMET-QE* | 4 0.686 | 8 0.697 | 6 0.453 | 3 0.875 | 6 0.560 | 1 0.859 | 6 0.471 | 1 0.920 | 3 0.558 | 4 0.838 | 4 0.643 | 5 0.761 | 5 0.537 | 3 0.926 | 5 0.513 | 1 0.884 | 2 0.448 | 2 0.842 | 3 0.557 |
| *CometKiwi-XXL** | 5 0.681 | 2 0.915 | 4 0.453 | 3 0.851 | 8 0.542 | 3 0.777 | 3 0.473 | 3 0.907 | 3 0.561 | 1 0.832 | 3 0.621 | 3 0.786 | 4 0.542 | 3 0.887 | 6 0.514 | 4 0.755 | 6 0.450 | 3 0.809 | 3 0.560 |
| YiSi-1 | 5 0.677 | 3 0.904 | 3 0.460 | 3 0.851 | 7 0.548 | 3 0.824 | 3 0.491 | 3 0.872 | 5 0.530 | 2 0.923 | 6 0.621 | 3 0.856 | 5 0.536 | 3 0.836 | 7 0.498 | 4 0.755 | 3 0.421 | 3 0.779 | 7 0.533 |
| PrismRefSmall | 5 0.676 | 2 0.924 | 4 0.471 | 4 0.804 | 7 0.549 | 2 0.816 | 6 0.473 | 3 0.835 | 4 0.556 | 2 0.956 | 7 0.614 | 4 0.846 | 4 0.545 | 3 0.841 | 6 0.509 | 4 0.731 | 8 0.408 | 4 0.759 | 4 0.555 |
| PrismRefMedium | 6 0.673 | 3 0.897 | 3 0.472 | 4 0.797 | 7 0.549 | 3 0.804 | 6 0.471 | 3 0.842 | 6 0.561 | 2 0.920 | 6 0.617 | 4 0.820 | 4 0.542 | 3 0.827 | 7 0.507 | 3 0.735 | 7 0.414 | 3 0.763 | 4 0.554 |
| chrfS | 7 0.666 | 3 0.903 | 5 0.458 | 6 0.762 | 9 0.534 | 3 0.785 | 4 0.481 | 3 0.851 | 5 0.528 | 3 0.884 | 7 0.612 | 1 0.895 | 4 0.541 | 3 0.832 | 8 0.493 | 5 0.721 | 7 0.416 | 4 0.767 | 8 0.529 |
| chrF | 8 0.662 | 3 0.903 | 6 0.454 | 5 0.781 | 10 0.530 | 4 0.750 | 5 0.479 | 3 0.852 | 6 0.523 | 3 0.902 | 7 0.615 | 2 0.865 | 6 0.528 | 3 0.858 | 9 0.486 | 5 0.741 | 8 0.412 | 4 0.713 | 9 0.522 |
| BERTScore | 9 0.656 | 6 0.826 | 6 0.455 | 6 0.742 | 11 0.528 | 4 0.750 | 7 0.480 | 6 0.820 | 6 0.524 | 3 0.868 | 8 0.591 | 1 0.869 | 6 0.531 | 3 0.810 | 11 0.485 | 5 0.716 | 7 0.413 | 5 0.773 | 8 0.531 |
| spBLEU | 9 0.652 | 2 0.915 | 9 0.440 | 5 0.767 | 11 0.527 | 4 0.732 | 8 0.473 | 3 0.838 | 7 0.515 | 2 0.906 | 8 0.601 | 7 0.842 | 7 0.516 | 3 0.833 | 11 0.473 | 4 0.736 | 9 0.408 | 4 0.701 | 10 0.518 |
| sentinel-cand-mqm* | 9 0.649 | 7 0.741 | 8 0.446 | 4 0.814 | 8 0.540 | 3 0.753 | 8 0.460 | 4 0.822 | 8 0.510 | 5 0.831 | 8 0.607 | 7 0.665 | 7 0.509 | 3 0.922 | 6 0.506 | 5 0.859 | 5 0.434 | 4 0.729 | 7 0.536 |
| *CometKiwi** | 10 0.641 | 8 0.683 | 7 0.447 | 7 0.784 | 7 0.545 | 4 0.750 | 6 0.470 | 4 0.908 | 8 0.558 | 5 0.808 | 9 0.590 | 6 0.719 | 5 0.530 | 5 0.702 | 8 0.490 | 4 0.737 | 5 0.433 | 3 0.832 | 7 0.556 |
| BLEU | 10 0.637 | 3 0.893 | 11 0.415 | 6 0.743 | 13 0.513 | 5 0.704 | 7 0.469 | 4 0.823 | 9 0.500 | 2 0.891 | 10 0.567 | 6 0.846 | 7 0.514 | 4 0.807 | 12 0.460 | 5 0.713 | 10 0.402 | 4 0.698 | 10 0.517 |
| damonmonli | 10 0.633 | 3 0.892 | 6 0.455 | 6 0.754 | 14 0.508 | 5 0.695 | 8 0.463 | 4 0.878 | 9 0.518 | 5 0.779 | 11 0.550 | 3 0.835 | 6 0.519 | 6 0.644 | 10 0.478 | 9 0.683 | 9 0.404 | 9 0.826 | 10 0.517 |
| metametrics_mt_mqm_qe_same_source_r* | 11 0.630 | 9 0.596 | 6 0.429 | 4 0.817 | 12 0.521 | 5 0.761 | 9 0.458 | 3 0.828 | 5 0.529 | 3 0.847 | 8 0.603 | 6 0.688 | 8 0.509 | 4 0.802 | 9 0.483 | 5 0.690 | 6 0.422 | 5 0.818 | 5 0.545 |
| MetaMetrics-MT-QE* | 11 0.630 | 9 0.594 | 10 0.429 | 4 0.817 | 12 0.521 | 3 0.762 | 9 0.458 | 3 0.826 | 5 0.529 | 3 0.847 | 8 0.603 | 6 0.687 | 8 0.509 | 4 0.801 | 9 0.483 | 5 0.690 | 6 0.422 | 2 0.821 | 5 0.545 |
| monmonli | 12 0.609 | 6 0.780 | 12 0.332 | 3 0.861 | 17 0.409 | 6 0.621 | 13 0.367 | 3 0.841 | 11 0.388 | 4 0.827 | 12 0.535 | 7 0.752 | 10 0.370 | 5 0.751 | 14 0.372 | 2 0.847 | 6 0.422 | 2 0.833 | 11 0.355 |
| gemba_esa* | 12 0.601 | 10 0.477 | 11 0.413 | 7 0.589 | 16 0.480 | 6 0.547 | 11 0.429 | 5 0.707 | 10 0.481 | 7 0.523 | 13 0.493 | 5 0.644 | 8 0.486 | 7 0.475 | 13 0.444 | 6 0.464 | 11 0.382 | 6 0.429 | 12 0.471 |
| XLsimDA* | 13 0.496 | 10 0.477 | 11 0.413 | 7 0.589 | 16 0.480 | 6 0.543 | 11 0.429 | 5 0.709 | 10 0.481 | 7 0.521 | 13 0.493 | 9 0.645 | 9 0.486 | 7 0.478 | 13 0.444 | 6 0.462 | 12 0.382 | 6 0.428 | 12 0.471 |
| XLsimMqm* | 13 0.496 | 10 0.475 | 11 0.413 | 7 0.587 | 16 0.480 | 6 0.543 | 11 0.429 | 5 0.709 | 10 0.481 | 7 0.521 | 13 0.493 | 9 0.645 | 9 0.486 | 7 0.478 | 13 0.444 | 6 0.462 | 12 0.382 | 6 0.428 | 12 0.471 |
| sentinel-ref-mqm | 14 0.330 | 9 0.547 | 13 0.176 | 8 0.466 | 18 0.070 | 5 0.656 | 14 0.170 | 5 0.745 | 12 0.046 | 7 0.358 | 15 0.059 | 8 0.419 | 11 0.034 | 5 0.739 | 15 0.137 | 6 0.402 | 12 0.278 | 5 0.617 | 14 0.025 |
| sentinel-src-mqm* | 14 0.330 | 9 0.542 | 13 0.176 | 8 0.468 | 18 0.070 | 3 0.761 | 12 0.170 | 5 0.739 | 12 0.046 | 7 0.361 | 15 0.059 | 8 0.420 | 11 0.034 | 5 0.737 | 15 0.137 | 6 0.404 | 12 0.278 | 5 0.617 | 14 0.025 |
| *BLCOM* | – | – | – | – | – | 6 0.537 | 10 0.445 | – | – | – | – | – | – | – | – | – | – | – | – |
| BLCOM_1 | – | – | – | – | – | 2 0.867 | 2 0.497 | – | – | – | – | – | – | – | – | – | – | – | – |
| MEE4 | – | – | – | – | – | 2 0.800 | 2 0.480 | – | – | – | – | – | – | – | – | – | – | – | – |
| bright-qe* | – | – | – | – | – | 3 0.761 | 12 0.396 | – | – | – | – | – | – | – | – | – | – | – | – |

Table 15: Correlations of metrics to the ESA annotations that were collected as part of the General MT Shared Task. The metrics are sorted by the average correlation across all of the correlations and language pairs. Metrics in italics are contrastive submissions and underlined metrics are baselines. QE metrics are marked by an asterisk.

# D Additional figures

Figures 6-10 show the (log) p-value of two-sided paired t-test on the MQM scores against the score difference of each metric for each system pair in each language pair. Figures 11-15 show the (log) p-value of significance test with bootstrap resampling on the metric scores against the score difference of that metric for each system pair in each language pair.



Figure 6: Log p-value of two-sided paired t-test on MQM scores ($p_{mqm}$) against the score difference of each metric (top to bottom: BERTSCORE, BLCOM_1, BLEU, BLEURT-20, BRIGHT-QE) for each system pair in each language pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_{mqm} < 0.05|\Delta M)$. Note: for readability, values of $p_{mqm}$ are rounded up to 0.0001 when they are less than 0.0001.

Figure 7: Log p-value of two-sided paired t-test on MQM scores ($p_{mqm}$) against the score difference of each metric (top to bottom: CHRF, CHRFS, COMET-22, COMETKIWI, DAMONMONLI, GEMBA_ESA) for each system pair in each language pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_{mqm} < 0.05|\Delta M)$. Note: for readability, values of $p_{mqm}$ are rounded up to 0.0001 when they are less than 0.0001.

Figure 8: Log p-value of two-sided paired t-test on MQM scores ($p_{mqm}$) against the score difference of each metric (top to bottom: MEE4, METAMETRICS_MT_MQM_HYBRID_KENDALL, METAMETRICS_MT_MQM_QE_KENDALL.SEG.S, METRICX-24-HYBRID, METRICX-24-HYBRID-QE) for each system pair in eachlanguage pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_{mqm} < 0.05|\Delta M)$. Note: for readability, values of $p_{mqm}$ are rounded up to 0.0001 when they are less than 0.0001.

Figure 9: Log p-value of two-sided paired t-test on MQM scores ($p_{mqm}$) against the score difference of each metric (top to bottom: PRISMREFMEDIUM, PRISMREFSMALL, SENTINEL-CAND-MQM, SENTINEL-REF-MQM, SENTINEL-SRC-MQM, SPBLEU) for each system pair in eachlanguage pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_{mqm} < 0.05|\Delta M)$. Note: for readability, values of $p_{mqm}$ are rounded up to 0.0001 when they are less than 0.0001.

75

Figure 10: Log p-value of two-sided paired t-test on MQM scores ($p_{mqm}$) against the score difference of each metric (top to bottom: XCOMET, XCOMET-QE. XLsimMqm, YiSi-1) for each system pair in eachlanguage pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_{mqm} < 0.05|\Delta M)$. Note: for readability, values of $p_{mqm}$ are rounded up to 0.0001 when they are less than 0.0001.

76

Figure 11: Log p-value of significance test with bootstrap resampling ($p_M$) on system-level metric scores against each metric (top to bottom: BERTSCORE, BLCOM_1, BLEU, BLEURT-20, BRIGHT-QE, CHRF) score difference for each system pair in each language pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_M < 0.05|\Delta M)$. Note: for readability, values of $p_M$ are rounded up to 0.0001 when they are less than 0.0001.

Figure 12: Log p-value of significance test with bootstrap resampling ($p_M$) on system-level metric scores against each metric (top to bottom: CHRFS, COMET-22, COMETKIWI, DAMONMONLI, GEMBA_ESA, MEE4) score difference for each system pair in each language pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_M < 0.05|\Delta M)$. Note: for readability, values of $p_M$ are rounded up to 0.0001 when they are less than 0.0001.

Figure 13: Log p-value of significance test with bootstrap resampling ($p_M$) on system-level metric scores against each metric (top to bottom: METAMETRICS_MT_MQM_HYBRID_KENDALL, METAMETRICS_MT_MQM_QE_KENDALL.SEG.S, METRICX-24-HYBRID, METRICX-24-HYBRID-QE, PRISMREFMEDIUM, PRISMREFSMALL) score difference for each system pair in each language pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_M < 0.05|\Delta M)$. Note: for readability, values of $p_M$ are rounded up to 0.0001 when they are less than 0.0001.

Figure 14: Log p-value of significance test with bootstrap resampling ($p_M$) on system-level metric scores against each metric (top to bottom: SENTINEL-CAND-MQM, SENTINEL-REF-MQM, SENTINEL-SRC-MQM, SPBLEU, XCOMET, XCOMET-QE) score difference for each system pair in each language pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_M < 0.05|\Delta M)$. Note: for readability, values of $p_M$ are rounded up to 0.0001 when they are less than 0.0001.

Figure 15: Log p-value of significance test with bootstrap resampling ($p_M$) on system-level metric scores against each metric (top to bottom: XLSIMMQM, YISI-1) score difference for each system pair in each language pair (left to right: en→de, en→es, ja→zh). The red line is the isotonic regression fit to all data points, representing $Pr(p_M < 0.05|\Delta M)$. Note: for readability, values of $p_M$ are rounded up to 0.0001 when they are less than 0.0001.

# Findings of the Quality Estimation Shared Task at WMT 2024
# Are LLMs Closing the Gap in QE?

**Chrysoula Zerva**[(1,2)]*, **Frédéric Blain**[(3)]*, **José G. C. de Souza**[(4)], **Diptesh Kanojia**[(5)],
**Sourabh Deoghare**[(6)], **Nuno M. Guerreiro**[(1,2,4,10)], **Giuseppe Attanasio**[(1)], **Ricardo Rei**[(2,4,7)],
**Constantin Orăsan**[(5)], **Matteo Negri**[(8)], **Marco Turchi**[(11)], **Rajen Chatterjee**[(9)],
**Pushpak Bhattacharyya**[(6)], **Markus Freitag**[(12)], **André F. T. Martins**[(1,2,4)]

[(1)]Instituto de Telecomunicações, [(2)]Instituto Superior Técnico, Universidade de Lisboa, [(3)]Tilburg University,
[(4)]Unbabel, [(5)]University of Surrey, [(6)]IIT Bombay, [(7)]INESC-ID, [(8)]FBK, [(9)]Apple Inc.
[(10)]MICS, CentraleSupélec, Université Paris-Saclay, [(11)]Zoom Video Communications, [(12)]Google Inc.
wmt-qe-organizers@googlegroups.com

## Abstract

We report the results of the WMT 2024 shared task on Quality Estimation, in which the challenge is to predict the quality of the output of neural machine translation systems without access to reference translations. In this edition, we continue to focus both on predicting sentence-level scores and on detecting error spans. Further, we expanded our scope to assess the potential for quality estimation to help in the correction of translated outputs, hence including an automated post-editing (APE) task.

We publish new test sets with human annotations that target two directions: providing new Multidimensional Quality Metrics (MQM) annotations for three multi-domain language pairs (English to German, Spanish and Hindi) and extending the annotations on Indic languages, providing direct assessments and post edits for translation from English into Hindi, Gujarati, Tamil and Telugu. We also perform a detailed analysis of the behaviour of different models with respect to different phenomena, including gender bias, idiomatic language, and numerical and entity perturbations. We received submissions based on both traditional encoder-based approaches and large language models (LLMs) and attempted to draw some comparisons in terms of performance and robustness to different phenomena.

## 1 Introduction

This edition of the shared task on Quality Estimation (QE) for machine translation builds upon previous iterations and findings, to further benchmark methods for estimating the quality of neural Machine Translation (MT) output at run-time, *i.e.* without relying on reference translations. The shared task introduces (sub)tasks that assess translation quality from multiple perspectives, examining errors both at a higher level (segment scores)

and with a more fine-grained view (error spans). Additionally, we expand our scope to generating corrected outputs through Automatic Post-Editing (APE).

Recently we have observed a gradual shift in the QE paradigms and methodologies, enabled by the advancement of neural metrics as well as large language models. Specifically, we have seen consistently strong performance across different language pairs and setups at sentence-level QE (Specia et al., 2021; Zerva et al., 2022; Blain et al., 2023), alongside increased efforts towards more finer-grained, explainable, and actionable evaluation of translations that focusses on error identification and explanation (Blain et al., 2023; Fernandes et al., 2023b; Guerreiro et al., 2023). The proliferation of LLM applications has led to significant performance improvements in MT, elevating the importance of advancing methodologies for quality estimation, and at the same time, it has allowed for novel perspectives and tasks related to quality estimation (Fabbri et al., 2022).

In light of the above, in this edition, we emphasise –beyond multilingual quality estimation– the analysis of the behaviour and abilities of submitted models with respect to different linguistic phenomena as well as their robustness to different error types and biases. Furthermore, we attempt to explore the degree to which quality estimation signals can be leveraged to improve translation quality via downstream automatic post-editing (Chatterjee et al., 2018b; Deoghare et al., 2023). We thus **bring APE under the QE umbrella** to make it easier for participants to develop QE systems and explore different techniques to apply it in APE shared task. These considerations collectively contribute to progress toward trustworthy and dependable QE systems that could facilitate real-time, reliable assessments of translation quality, as well as inform APE systems towards generating a corrected trans-

---

*Main organisers

lation.

In this edition of the shared task, we further expand the provided resources for sentence-level and fine-grained QE, providing new test sets and expanding to new language pairs. Following the previous editions, we provide annotations for *direct assessments* (DA; English-Tamil, English-Hindi, English-Telugu and English-Gujarati), *post-edits* (PE; English-Tamil and English-Hindi) and *Multidimensional Quality Metrics* (MQM; English-Hindi, English-Spanish and English-German) (Lommel et al., 2014). We describe in detail the annotation process and provide statistics for the new resources in Section 3.

Overall, in addition to advancing the state-of-the-art at all prediction levels, our main goals are:

- To extend the languages covered in our datasets and provide new test sets emphasising low- and medium-resource languages and zero-shot approaches;

- To continue investigating the potential of fine-grained quality estimation;

- To study the robustness of QE approaches to different linguistic phenomena, error types and biases;

- To continue monitoring the computational efficiency of proposed approaches for sustainability purposes; and

- To study whether we can leverage QE signals to improve translation quality via downstream APE task.

We thus designed three tasks this year:

**Task 1** The core QE task, which consists of separate sentence-level sub-tasks for different language pairs (§**??**). The goal is to predict a quality score for each segment in a given test set, which can be a variant of DA (§3.2) or MQM (§3.3).

**Task 2** The fine-grained error prediction task, where participants were asked to detect error spans alongside error severities (*Major* versus *Minor*) (§2.2).

**Task 3** A newly introduced task, which requires participants to combine quality estimation and automatic post-editing in order to correct the output of machine translation. (§2.3).

The tasks make use of large datasets annotated by professional translators with either $0 - 100$ DA scoring, post-editing or MQM annotations. We provide new training, development and test data for Task 3 as well as fresh new test sets for Tasks 1 and 2. The datasets and models released are publicly available[1].

Besides the data made available through the QE shared task, participants were also allowed to explore any additional data and resources deemed relevant, across tasks. In addition, LLMs could also be used both to extend resources and to complement predictions.

The shared task uses *CodaBench* as a submission platform, where each sub-task corresponds to a separate competition instance. Participants (Section 5) could submit up to a total of 10 submissions per sub-task. Results for all tasks, evaluated according to standard metrics, are given in Section 6. Baseline systems were trained by the task organisers and entered into the platform to provide a basis for comparison (Section 4). We provide an additional evaluation focussed on robustness against different phenomena and biases in Section 7. A discussion on the main findings from this year's task is presented in Section 8.

## 2 Quality Estimation tasks

In what follows, we briefly describe each sub-task, including the datasets provided for them.

### 2.1 Task 1: Predicting translation quality

The ability to accurately estimate the quality of translations on-the-fly, i.e., without access to human references, is at the core of the QE shared task. This year, we focus on sentence-level quality, attempting to disentangle finer-grained analysis or post-edits that are tackled in Tasks 2 and 3.

Similar to the last edition, the data was produced as follows:

1. DA sentence level scores: The quality of each source-translation pair is annotated by at least 3 independent expert annotators, using DA on a scale 0-100.

2. MQM annotation: Each source-translation pair is evaluated by at least 1 expert annotator, and errors identified in texts are highlighted

---

[1] https://github.com/WMT-QE-Task/wmt-qe-2023-data

and classified in terms of severity (minor, major, critical) and type (grammar correctness, omission, style, mistranslation, among others).

The DA and MQM sentence level annotations were further processed to obtain normalised quality scores that have the same direction between high and low quality. We provide more details on the required pre-processing in §2.1.1.

### 2.1.1 Sentence-level quality prediction

Similarly to the previous year, we used a single competition instance both for DA and MQM-derived annotations aiming to motivate the submission of models that are robust to both annotation formats. Hence, we also aligned the scores by processing and normalising them as follows:

- For the **DA** scores we standardize the scores with respect to each annotator and then compute the mean average of standardized scores for each sentence.

- For the **MQM** scores we need to first compute the overall score from the individual errors. Hence for each annotator, we first compute the sentence-level score as:

$$MQM^{sent}(hyp) = \frac{100 - \sum\limits_{e \in hyp} severity(e)}{|hyp|},$$
(1)

where $hyp$ is a hypothesis sentence represented as a sequence of tokens, $e$ is an error annotated in that sentence and the $severity$ is computed but adding:

+ 1 point for minor errors
+ 5 points for major errors
+ 10 points for critical errors

To align with DA annotations, we subtract the summed penalties from 100 (perfect score) and we then divide by the sentence length (computed as number of words). We then normalise per annotator as in the DA case and compute the mean average in the case of multiple annotators.

Regarding evaluation, systems in this task (both for DA and MQM) are **evaluated against the true z-normalised sentence scores using Spearman's rank correlation coefficient $\rho$ as the primary metric**. This is what was used for ranking system

submissions. Pearson's correlation coefficient, $r$, and Kendall $\tau$ were also computed as secondary metrics but not used for the final ranking of systems.

### 2.1.2 Finer-grained Evaluation and Challenge Sets

To assess the robustness and capabilities of automatic machine translation evaluation systems, we created a challenge set focusing on five different phenomena for the En-De and En-Es language pairs. Each category tests a particular aspect of translation quality that may have impact in real-world applications. The challenge set aims to determine whether evaluation systems can distinguish between correct translations—which we designate as hyp—and those containing subtle but relevant variations—which we designate as con.

**Currency and date formatting** This set tests the detection of format changes in currency symbols and date expressions. The hyp preserves the original source format (e.g., keeping "$100" or "MM/D-D/YYYY"), while the con presents localized versions (e.g., "100 USD" or "DD/MM/YYYY"). Note that here it is the case that con is also a good-quality translation.

**Word order** This category examines the handling of word order variations. The hyp consists of monotonic translations that closely follow the source sentence order, while the con presents non-monotonic translations that rearrange words while preserving meaning. Evaluation models might have a preference towards one or the other, even though both preserve the meaning of the source.

**Detached translations and omissions** This set focuses on critical divergences from the source text. The hyp provides accurate and complete translations of the source. In contrast, the con includes examples where translations start correctly but then veer into unrelated topics or omit substantial portions of the source text. Evaluation systems are expected to detect these critical errors.

**Idiomatic translations** This category tests the handling of idiomatic expressions. The hyp presents idiomatic renderings that accurately convey the meaning in the target language, while the con offers literal word-for-word translations that may render the target text non-sensical. Evaluation systems should appropriately score translations that

prioritize conveying the correct meaning over strict word-for-word translation.

We generated data for all the phenomena listed above using GPT-3.5 (`gpt-3.5-turbo-0125`) and GPT-4 (`gpt-4-1106-preview`). Then, we conducted a human annotation study to discard erroneous triples.

**Gender Subset**  The gender subset of the challenge set aims to study QE metrics and gender inflection in grammatical gender languages.

Following Zaranis et al. (2024), we collected unmodified instances from the counterfactual subsets (Es and De) of MT-GenEval (Currey et al., 2022), an evaluation set for sentence-level gender bias in machine translation. In these examples, sources from English Wikipedia mention exactly one human entity and contain intra-sentence lexical clues that help disambiguate the entity's gender identity.[2] Each source is provided with a masculine (M) and a feminine (F) variant (e.g., "She/He is a graduate of Harvard, but rarely applies such skills."). Human references are included as well.

We compiled the gender subset by constructing contrastive pairs as follows. First, we sampled 150 instances from the original MT-GenEval's subset. Fifty unique sources have a female referent and fifty a male referent. From each instance, we created a triplet with the source, the reference with correct gender inflection used as hypothesis, and the reference with wrong gender inflection used as contrast. Then, to isolate the impact of the source content, we created two triplets for each of the remaining fifty instances. The source in the triples is identical except for the gender identity of the entity. This step yields 100 more examples. The gender subset hence counts 200 contrastive triplets in total.

## 2.2 Task 2: Fine-grained error detection

For this task, we focus on finer-grained quality predictions, taking advantage of the detailed information provided in the MQM annotation schema. Specifically, each error span is annotated with error severity (*minor*, *major*, *critical*) as well as error type (see also Figure 1). Following the findings of the previous edition, we focus on the severity annotations and do not use the other error categories annotated in the MQM schema. As a result, we aimed to (1) identify error spans and (2) classify

said error spans as either *minor* or *major*. We note that we merge the critical and major categories, since in this edition we noticed particularly sparse occurrences of critical errors (even less than the previous year). Additionally, in this edition, the annotations included a *neutral* category, which was ignored as it was (1) not occurring for all language pairs and (2) they correspond to subjective opinions/preferences about translation. [3] We point readers to Figure 3 for some statistics on error severity distribution per language pair and domain.

The information used for this task consists of: *i*) start and end index positions for each error span; and *ii*) the simplified error severity. The error spans are identified as continuous sequences of characters within a target hypothesis, allowing for annotations of single white spaces and punctuation marks in order to account for omission and punctuation errors, respectively. Aiming to mimic the human annotations and simplify the task, overlapping error spans are allowed and count towards *recall* of different errors, but overlapping annotations are flattened for both gold and system annotations (see below). Figure 1 shows an example of annotations.

For the evaluation, the primary metric is the **F1-score**, computed on the character level and weighted to allow for half points for correctly identified span but misclassified severity. Precision and recall were also provided as complementary metrics. With respect to overlapping annotations, we allow for multiple character level annotations[4] and consider the best matching annotation per character position. As such, for each segment, we compute recall for the characters in gold annotation text spans by computing the ratio between the overlap with system error spans and the gold error span length and weighting severity mismatches by $0.5$. Respectively, we compute precision with respect to the system error span length and apply the same weighting convention (down-weighting by 0.5 for mismatched error severities). Figure 1 and Table 1 show an example of the aforementioned process [5].

---

[2]We acknowledge a notion of gender identity beyond the binary. However, we include only masculine and feminine examples as they are provided in the original dataset.

[3]Note that the neural errors are also not considered when computing an MQM score.

[4]The gold data was processed to remove identical segments that correspond to the same span but have different error categories, but it preserved any partially overlapping segments that correspond to different error categories and/or severities.

[5]The link to evaluation scripts can be found at: `https://github.com/WMT-QE-Task/qe-eval-scripts/blob/main/wmt24/`

| Systems | Precision | Recall | F1-score |
|---|---|---|---|
| System A | $\frac{1*7+1*28+0.5*6}{7+28+13} = 0.79$ | $\frac{1*7+1*28+0.5*6}{12+28+6} = 0.83$ | 0.81 |
| System B | $\frac{0.5*12+1*28+0.5*6}{12+28+6} = 0.80$ | $\frac{1*12+1*28+0.5*6}{12+28+6} = 0.80$ | 0.80 |

Table 1: Example of Precision and Recall computations for each annotation in the example of Figure 1.



Figure 1: Example of gold annotations (MQM) for Task 2 (top) and respective prediction examples (bottom). Example taken from He-En test set.

## 2.3 Task 3: QE-informed APE

MT Automatic Post-Editing (APE) is the task of automatically correcting errors in a machine-translated text. As pointed out by Chatterjee et al. (2015), from the application point of view, the task is motivated by its possible uses to:

- Enhance MT output by harnessing information that is not available to the decoder or by conducting deeper text analysis, which may be prohibitively expensive during the decoding phase.

- Address systematic errors stemming from an MT system whose decoding process is inaccessible for focused modifications.

- Provide professional translators with improved MT output quality, thereby reducing the need for subsequent human post-editing.

- Tailor the output of a general-purpose MT system to align with the lexicon and style requirements of a specific application domain.

Building on the work of Chatterjee et al. (2018b); Deoghare et al. (2023), which demonstrated the potential of QE to enhance APE systems, this edition of the WMT QE shared task introduced the new QE-informed APE subtask. In this subtask,

we focus on a unified *evaluation and correction* paradigm, taking advantage of the additional information provided by the human post-edits. Participants were encouraged to incorporate signals from QE systems to improve APE performance. The evaluation setup remained consistent with the previous rounds WMT APE shared tasks, requiring participants to automatically correct translations generated by a generic, domain-unadapted "black-box" NMT system. The training data consisted of human post-edits of translations produced by this system. While TER (Snover et al., 2006) and BLEU (Papineni et al., 2002) continued as the primary and secondary evaluation metrics, this year also introduced chrF (Popović, 2015) and COMET[6] for a more comprehensive automatic evaluation of the submitted APE systems.

For this year, English-Hindi and English-Tamil were the selected language pairs, with Hindi and Tamil as the target languages for post-editing. The training, development, and test data encompassed a wide range of domains, including education, legal, healthcare, culture, tourism, reviews, subtitles, and general/news.

## 3 Datasets

Below, we describe the datasets provided to participants for development and testing. Specifically, this year, we provided training data only for Task 3, which was newly introduced (see §3.4).

### 3.1 Training Resources

Overall, participants were encouraged to employ training data from a wide range of sources, including datasets from previous competitions, as well as synthetic or proprietary data.

Proposed training data for DA annotations, following the previous editions, includes the language pairs from the MLQE-PE dataset (Fomicheva et al., 2022), as well as the data from the previous QE editions (Zerva et al., 2022; Blain et al., 2023). Similarly, for the MQM data, we encouraged participants to refer to data from previous editions that

---

[6]https://github.com/Unbabel/COMET .

cover translation into German (En-De), Russian (En-Ru), Hebrew (En-He) and out of Chinese (Zh-En) (Freitag et al., 2021a,b), as well as the Indic-MT eval dataset (Sai B et al., 2023). However, we emphasise that in this edition, we introduce no new training data, treating the translations into Spanish (En-Es) and Hindi (En-Hi) as zero-shot tasks, and only En-De as supervised.

## 3.2 Direct Assessment (DA) Data

For all language pairs, the data provided is selected from publicly available resources.

We expand the Indic language pairs introduced in previous years, providing new unseen test sets of approx 1K segments each for Hindi (Hi; 1000 segments) and Gujarati (Gu; 1012 segments) as target languages from the Indo-Aryan language family as well as Tamil (Ta; 1000 segments) and Telugu (Te; 1000 segments) from the Dravidian language family. Following the previous edition, dataset curation and annotation were performed with the help of professional translators who were native speakers of the target language. The annotators were provided with guidelines which discussed DA score ranges with various error types. Additionally, parallel segments were curated from the following parallel corpora: *i) Anuvaad* parallel corpus[7] (General, Healthcare and Legal domain; *ii)* IITB English-Hindi parallel corpus[8] (Kunchukuttan et al., 2018) (Culture/Tourism domain), and parallel segments scraped from NPTEL[9]; and *iii)* SpokenTutorials[10] (Education domain). The curated segments were selected from the above-mentioned domains to ensure cross-domain impact and performance.

From the *Anuvaad* parallel corpus, we filter parallel segments using LaBSE, and select source sentences with varying token lengths, while the translation was obtained using $1.3B$ parameter NLLB model (Costa-jussà et al., 2022), as discussed in (Blain et al., 2023). During the annotation, weekly validation of randomly selected instances was performed by an unbiased native speaker who provided feedback to further improve annotations during the data curation. After all three annotators performed the DA annotations, we separated the data into training, development, and test



Figure 2: Distribution of DA scores for the Indic language pairs.

sets while filtering for a balanced distribution of DA scores across all sets. We provide the distribution of DA scores for each language pair in Figure 2, where we can see that for all language pairs, we have similar distributions skewed towards high-quality scores. We can also observe that for Tamil, we have fewer segments of very low quality (DA $\leq 20$), but instead, we have larger counts of segments of moderate quality ($20 \leq$ DA $\leq 60$).

## 3.3 MQM Data

As **test data**, we annotated new evaluation sets for three language directions: English-German (En-De), English-Spanish (En-Es) and English-Hindi (En-Hi). The evaluation sets were annotated by professional translators following a MQM typology (Burchardt, 2013) and specific guidelines[11].

The documents used for the evaluation sets are shared with the WMT General MT task and follow the same distribution of domains in that data (*e.g.*, news, social, literary and speech). The full documents were translated using the $54B$ parameters NLLB model (Team et al., 2022)[12] without sentence splitting. We subsequently split segments for annotation and annotated a total of 1511 segments for each translation direction.

The test data distribution according to error severities is shown in Figure 3. The NLLB model used to translate the evaluation sets is clearly stronger for En-De, with less than 100 major and minor errors for each content type. The distribu-

---

[7]https://github.com/project-anuvaad/anuvaad-parallel-corpus
[8]Unreleased parallel segments, to be released here in v3.2: https://www.cfilt.iitb.ac.in/iitb_parallel/
[9]https://nptel.ac.in/
[10]https://spoken-tutorial.org/

[11]http://bit.ly/mqm-guidelines
[12]Model identifier FACEBOOK/NLLB-MOE-54B

87

Figure 3: Distribution of error severities across language pairs and domains/content types.



Figure 4: Distribution of average length (character count) for different severities across language pairs and domains/content types.

tion of major and minor errors changes drastically for En-Es and En-Hi, in particular the number of minor errors for the literary, social and speech domains, with more than 200 minor errors each. In addition, we can see that we have fewer errors for the news domain across all three language pairs, both in terms of minor and major errors. Contrary to frequency, however, Figure 4 shows that error spans identified for En-De are significantly longer on average for both identified error categories.

### 3.4 QE-APE Data

This year we introduce two new language pairs for the APE task: English-Hindi (*En-Hi*) and English-Tamil (*En-Ta*). For each language pair, the train, dev, and test sets respectively consist of $7,000$, $1,000$, and $1,000$ (*source, target, human post-edit*) triplets, where:

- The source (SRC) is an English sentence;

- The target (TGT) is a Hindi/Tamil translation of the source produced by a generic, black-box NMT system unknown to participants.

- The human post-edit (PE) is a manually revised version of the target, which was pro-

duced by native Hindi/Tamil speakers.

The English-Hindi train, dev, and test sets span culture, education, health, tourism, and general domains. Similarly, English-Tamil APE datasets contain sentences from legal, literacy, reviews, subtitles, news, health, and general domains.

We also provide a corpus of artificially generated data as additional training material. It consists of 2.5 million triplets for each language pair derived from the Anuvaad parallel corpus. Specifically, the source, target, and post-edit instances of this synthetic corpus are respectively obtained by combining: i) the original English source sentence from the Anuvaad corpus, ii) its automatic translation into Marathi, iii) the original Marathi target sentence from the Anuvaad corpus. Furthermore, we provide the DA scores for all samples in both train and dev sets. Additionally, the participants were encouraged to use the DA data released in the earlier iteration of the QE shared task for these language pairs.

To get an idea of the task difficulty, we focused on three aspects of the released data, which provided us with information about the possibility of learning useful correction patterns during APE

| | Lang. | Domain | MT type | RR_src | RR_tgt | RR_pe | Basel. BLEU | Basel. TER | δ TER |
|---|---|---|---|---|---|---|---|---|---|
| 2015 | en-es | News | PBSMT | 2.9 | 3.31 | 3.08 | n/a | 23.84 | +0.31 |
| 2016 | en-de | IT | PBSMT | 6.62 | 8.84 | 8.24 | 62.11 | 24.76 | -3.24 |
| 2017 | en-de | IT | PBSMT | 7.22 | 9.53 | 8.95 | 62.49 | 24.48 | -4.88 |
| 2017 | de-en | Medical | PBSMT | 5.22 | 6.84 | 6.29 | 79.54 | 15.55 | -0.26 |
| 2018 | en-de | IT | PBSMT | 7.14 | 9.47 | 8.93 | 62.99 | 24.24 | -6.24 |
| 2018 | en-de | IT | NMT | 7.11 | 9.44 | 8.94 | 74.73 | 16.84 | -0.38 |
| 2019 | en-de | IT | NMT | 7.11 | 9.44 | 8.94 | 74.73 | 16.84 | -0.78 |
| 2019 | en-ru | IT | NMT | 18.25 | 14.78 | 13.24 | 76.20 | 16.16 | +0.43 |
| 2020 | en-de | Wiki | NMT | 0.65 | 0.82 | 0.66 | 50.21 | 31.56 | -11.35 |
| 2020 | en-zh | Wiki | NMT | 0.81 | 1.27 | 1.2 | 23.12 | 59.49 | -12.13 |
| 2021 | en-de | Wiki | NMT | 0.73 | 0.78 | 0.76 | 71.07 | 18.05 | -0.77 |
| 2022 | en-mr | health/tourism/news | NMT | 1.46 | 0.89 | 0.72 | 67.55 | 20.28 | -3.49 |
| 2023 | en-mr | health/tourism/news | NMT | 1.85 | 1.24 | 1.12 | 70.66 | 26.60 | +1.13 |
| 2024 | en-hi | health/tourism/news | NMT | 2.7 | 3.55 | 3.32 | 39.28 | 46.36 | -19.29 |
| 2024 | en-ta | health/tourism/news | NMT | 1.97 | 1.49 | 1.1 | 70.16 | 24.71 | -0.47 |

Table 2: Basic information about the APE shared task data released since 2015- languages, domain, type of MT technology, repetition rate and initial translation quality (TER/BLEU of TGT). The last column (δ TER) indicates, for each evaluation round, the difference in TER between the baseline (*i.e.,* the "*do-nothing*" system) and the top-ranked official submission.

model training and successfully applying them at test time. These are: *i)* repetition rate, *ii)* MT quality, and *iii)* TER distribution in the test set. For the sake of comparison across the nine rounds of the APE task (2015–2023), Table 2 reports, for each dataset, information about the first two aspects. The third aspect, however, will be discussed by referring to Figure 5 and Figure 6.

### 3.4.1 Repetition Rate

The repetition rate (RR), measures the repetitiveness inside a text by looking at the rate of non-singleton $n$-gram types ($n = 1...4$) and combining them using the geometric mean. Larger values indicate a higher text repetitiveness that may suggest a higher chance of learning from the training set correction patterns that are also applicable to the test set. However, over the years, the influence of repetition rate in the data on system performance was found to be marginal.[13]

As shown in Table 2, in this edition, the RR for English-Hindi ranges between 2.7-3.3, and for English-Tamil RR ranges between 1.1-2.0. This difference may contribute to motivating the significantly different APE results observed for the two languages, as evidenced by a substantial TER reduction for English-Hindi ($-19.29$ "δ TER") compared to the "do-nothing" the baseline (see §4.3). Reviewing previous rounds of the APE task, however, suggests that RR remains only a partially in-

formative indicator of task difficulty due to its variable correlation with final results, which may also depend on other factors or on the interaction of multiple factors that are yet to be fully understood.

### 3.4.2 MT Quality

Another complexity indicator is MT quality, which is the initial quality of the machine-translated (TGT) texts to be corrected. We measure it by computing the TER (↓) and BLEU (↑) scores (Basel. TER/BLEU rows in Table 2) using the human post-edits as reference. In principle, higher quality of the original translations leaves the APE systems with less room for improvement since they have, at the same time, less to learn during training and less to correct at the test stage. On one side, training on good (or near-perfect) automatic translations can drastically reduce the number of learned correction patterns. On the other side, testing on similarly good translations can *i)* drastically reduce the number of corrections required and the applicability of the learned patterns, and *ii)* increase the chance of introducing errors, especially when post-editing near-perfect translations. The findings of all previous rounds of the task support this observation, which is corroborated by the high correlation (>0.78) between the initial MT quality ("Basel. TER" in Table 2) and the TER difference between the baseline and the top-ranked submission ("δ TER" in Table 2).

As discussed in Section 6.3, this year seems to confirm the trends observed in the past. For English-Hindi, the baseline TER is quite high (46.36 points), leaving more room for improvement.

---

[13]The analyses carried out over the years produced mixed outcomes, with impressive final results obtained in spite of low repetition rates (Chatterjee et al., 2020) and vice-versa (Chatterjee et al., 2018a, 2019; Akhbardeh et al., 2021).

Whereas English-Tamil falls in medium-high difficulty (20.0<TER<25.0), making the task more challenging. The final gains ("$\delta$ TER" in Table 2) confirm the correlation between the quality of the initial translations and the actual potential of APE.



Figure 5: TER distribution in the APE 2024 English-Hindi test set.



Figure 6: TER distribution in the APE 2024 English-Tamil test set.

### 3.4.3 TER Distribution

A third complexity indicator is the TER distribution (computed against human references) for the translations present in the test sets. Although TER distribution and MT quality can be seen as two sides of the same coin, it's worth remarking that, even at the same level of overall quality, more/less peaked distributions can result in very different testing conditions. Indeed, as shown by previous analyses, harder rounds of the task were typically characterised by TER distributions particularly skewed towards low values (*i.e.,* a larger percentage of test items having a TER between 0 and 10). On one side, the higher the proportion of (near-)perfect test instances requiring few edits or no corrections at all, the higher the probability that APE systems will

perform unnecessary corrections penalised by automatic evaluation metrics. On the other side, less skewed distributions can be expected to be easier to handle as they give automatic systems larger room for improvement (*i.e.,* more test items requiring - at least minimal - revision). In the lack of more focused analyses on this aspect, we can hypothesise that in ideal conditions from the APE standpoint, the peak of the distribution would be observed for "post-editable" translations containing enough errors that leave some margin for focused corrections but not too many errors to be so unintelligible to require a whole re-translation from scratch.[14]

As shown in Figure 5, for English-Hindi the TER distribution follows more or less uniform distribution. The distribution is not too skewed towards near-perfect translation (which would have made it harder to further improve), nor towards the higher end of TER (which would have made it harder to learn error-correction patterns due to too noisy data). These characteristics make it easier to improve translation, which is reflected in the final evaluation results. On the other hand, as shown in Figure 6, for English-Tamil the TER distribution is highly skewed towards near-perfect translations. Around half of the test set falls in 0-5 TER points, making it prone to over-correction, which can be penalised by automatic evaluation metrics. This characteristic makes the English-Tamil test set much more challenging when it comes to gaining further translation quality improvements.

## 4 Baselines

In this edition, we opted to use publicly available, existing models without further tuning. Hence, we use a more unified architecture for Tasks 1 and 2, where all models use a large XLM-RoBERTa pre-trained encoder without additional language tuning (see also Appendix A for hyperparameter details). The specific hyperparameters used are presented in Table 7. For Task 3, we opted for a simple "do nothing" approach as discussed in Section 4.3.

### 4.1 Task 1: Quality Estimation

For the **sentence-level** sub-task, we opted for using CometKiwi 2022 (Rei et al., 2022) which was trained on data from the Metrics and QE shared tasks (combining data from previous years up to

---

[14]For instance, based on the empirical findings reported in (Turchi et al., 2013), TER=0.4 is the threshold that, for human post-editors, separates the "post-editable" translations from those that require complete rewriting from scratch.

2022). Models are publicly available for download[15].

### 4.2 Task 2: Fine-grained Error Detection

For **Task 2** we also used a CometKiwi model, specifically one trained on the multi-task setting, to produce both sentence-level scores and word-level quality estimates. The model, trained on 2022 QE data is publicly available.[16] The word-level estimates are in the form of OK/BAD tags, and for this reason it is necessary to convert the original output to the one required by the Task 2 format. As such we process the word-level predictions as follows:

- Detokenize the sentence

- Annotate continuous BAD tokens as a single text span

- Assume all errors are major

### 4.3 Task 3: QE-informed APE

The official baseline results for **Task 3** are the TER/BLEU/chrF/COMET scores calculated by comparing the raw MT output with human post-edits. This corresponds to the score achieved by a "do-nothing" APE system that leaves all the test segments unmodified.

## 5 Participants

In this section, we present a brief system description gathered from each participant. For each team, we indicate the task(s) and sub-task(s) (*i.e.* language-pair(s)) they participated in, and point to relevant publications, if any.

**Unbabel (T1; all):** The submission for Task 1 follows their work from the previous competition (Rei et al., 2023), which corresponds to an ensemble of multiple checkpoints for the sentence-level subtask, using a weighted averaging of the predicted scores, optimised by language pair. The emphasis is on scaling the size of the pre-trained encoder from InfoXLM to XLM-R XL and XXL.

**Pister Labs (T1; all):** The team opted for an approach where they generated a set of reading comprehension questions and scored each hypothetical translation by evaluating how well it could answer the comprehension question when compared with the reference translation. The overall score for a hypothetical is then a simple average across the questions asked of it. Answers are generated by providing the question and the hypothetical translation to Llama3.1-8B (Dubey et al., 2024). The initial set of reading comprehension questions is generated through few-shot prompting of Llama3.1-70B, and evaluating results on a subsample of 100 training En-De translation pairs with Llama3.1-70B. The four questions with the highest Spearman correlation were then used for final testing. To improve question generation quality, they use techniques from OpenAI and Anthropic's prompting guides, as well as the self-consistency technique.

**HW-TSC (T1; En-Hi, En-Ta, En-Te, En-Gu):** The team employed the CROSS-QE approach (Li et al., 2023) as the basis for further tuning and opted for tuning separate models for each language pair. They used encoder-based models, experimenting with different encoders, which were trained on different combinations of source and translation vectors as input. They focused on improving model performance both in terms of training by employing different data augmentation methods and in terms of inference, exploring better strategies for ensembling checkpoints. In terms of data augmentation, they use a combination of LLMs with specific prompts to generate pseudo-data as well as text editing methods.[17]

**HW-TSC (T2; all):** The team employs a combination of LLMs, hypothesising that the reasoning abilities of large models may be helpful in the fine-grained task. They use the TowerInstruct-7B-v0.2 (Alves et al., 2024) model and the GPT-4o-mini (Islam and Moushi, 2024) model, using prompt engineering and in-context learning to obtain the predictions. Additionally, they employ data augmentation techniques mentioned for Task 1 and find that they can rely on pseudo-data for tuning the models.[18]

---

[15] https://huggingface.co/Unbabel/wmt22-cometkiwi-da

[16] https://huggingface.co/Unbabel/WMT24-QE-task2-baseline

[17] We consider submissions from users s50042889 and zhaoxf4 mentioned in the results page as one submission

[18] We consider submissions from users zhuming, zhaoxf4 and mengyao mentioned in the results page as one submission

**TMU-HIT (T1; En-Hi, En-Ta, En-Te, En-Gu):** The team submitted predictions that rely on LLMs, inspired by (Liu et al., 2023; Enomoto et al., 2024). They designed custom prompts for quality estimation and employed GPT-4o mini (Achiam et al., 2023) to sample assessment scores multiple times using the same prompt. They then experimented with combining the generated scores to compute the final score using either their average or their weighted sum, employing the generation probabilities as weights for the latter. They conducted evaluation experiments in both zero-shot and three-shot settings. Further, they also attempted fine-tuning GPT-4o mini using the training data released for the WMT23 Machine Translation task (Kocmi et al., 2023).

**HW-TSC (T3; all):** (Yu et al., 2024) The team explored two distinct approaches for developing APE systems. For the En-Hi pair, they leveraged the Llama3-8B-Instruct model through continual pre-training on the collected data and then supervised fine-tuning it on the real APE data. For the En-Ta pair, they trained a transformer model from scratch, first focusing on the MT (Machine Translation) task using web-collected data, followed by training on APE data. External MT candidates were incorporated during the training to boost performance further. To prevent over-correction, Sentence-level QE models were employed to select between MT and APE outputs. Both users (**HW-TSC_yjwsss** and **HW-TSC_zhaoxf4**) from this team made the same submissions for En-Ta, but different submissions for En-Hi.

**IT-Unbabel (T3; all):** IT-Unbabel submission leveraged xTower (Treviso et al., 2024), a model built on top of TowerLLM (Alves et al., 2024), which is designed to provide free-text explanations for translation errors to guide the generation of an improved translation. The system was trained on material that includes the xTower dataset (GPT-4 generated explanations for translation correction), TowerBlocks, and additional training datasets provided by the WMT24[19] organizers for English-Hindi and English-

Tamil, augmented with error span annotations from xCOMET (Guerreiro et al., 2023). A hybrid approach is used to dynamically select between the original translation and the corrected version produced by the xTower model using a quality estimation model.

## 6 Results

In this section, we present and discuss the results of our shared task. Please note that for all the three sub-tasks we used statistical significance testing with $p = 0.05$.

### 6.1 Task 1

As described in the Task 1 overview (§2.1.1), sentence-level submissions are evaluated against the true z-normalised sentence scores using Spearman's rank correlation coefficient $\rho$ along with the following secondary metrics: Pearson's correlation coefficient, $r$, and Kendall's $\tau$. Nonetheless, the final ranking between systems is calculated using the primary metric only (Spearman's $\rho$). Statistical significance was computed using William's test. The results are shown in Table 3.

Looking at the obtained scores, we observe an overall performance increase for the sentence-level scores compared to previous years for all language pairs (that have been previously tested) except for En-Ta, where we observe a small drop. We note, that while the domains and sources in the En-De MQM test-set are different, all DA test-sets are drawn from the same sources and observe similar score distributions to previous years, thus facilitating comparisons.

It should be noted that there is no clear winner across language pairs. Instead, different systems rank first for each language.

### 6.2 Task 2

For Task 2, the submissions are scored using the F1-score, computed at character level for the annotated error spans, as described in Section 2.2. Precision and Recall scores are also provided as complementary information to help contextualise the performance observed. Statistical significance was computed using randomisation tests (Yeh, 2000) with Bonferroni correction (Abdi, 2007) for each language pair. The results for Task 2 are described in Table 4.

This year, the fine-grained annotation task (Task 2) had a lower participation rate compared to the

---

[19]https://www2.statmt.org/wmt24/

| Model | Multi | Multidimensional Quality Metric (MQM) | | | Direct Assessment (DA) | | | |
|---|---|---|---|---|---|---|---|---|
| | | En-De | En-Es | En-Hi | En-Hi | En-Gu | En-Te | En-Ta |
| Unbabel | 0.553 | 0.512 † | 0.345 † | 0.412 | 0.714 | 0.703 † | 0.510 † | 0.675 † |
| Pister Labs | 0.452 | 0.513 † | 0.242 | 0.363 | 0.564 | 0.587 | 0.379 | 0.478 |
| HW-TSC | - | - | - | - | 0.719 † | 0.757 † | 0.482 † | 0.683† |
| TMU-HIT | - | - | - | - | 0.739 † | 0.713 | 0.482 | 0.603 |
| BASELINE | 0.520 | 0.514 † | 0.340 † | 0.441 † | 0.678 | 0.661 | 0.414 | 0.592 |

Table 3: Spearman correlation for the official submissions to WMT24 Quality Estimation **Task 1 Sentence-level**. Baseline systems are highlighted in grey. For each language pair, results marked with † correspond to the winning submissions, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959).

.

| Model | Multidimensional Quality Metric (MQM) | | | |
|---|---|---|---|---|
| | Multi | En-De | En-Es | En-Hi |
| BASELINE | 0.278 | 0.192 † | 0.161 † | 0.481 † |
| HW-TSC | 0.227 | 0.178 | 0.151 | 0.362 |

Table 4: F1-score for the official submissions to WMT24 Quality Estimation **Task 2 Error Span Detection**. Baseline systems are highlighted in grey. For each language pair, results marked with † correspond to the best system (not significantly outperformed by any other system) according to randomized paired t-test.



(a) HWTSC - Major  (b) HWTSC -Minor

(c) Baseline - Major  (d) Baseline - Minor

Figure 7: Confusion matrices for Task 2 English-German, comparing Minor and Major predictions between the Baseline system and the HWTSC one.

previous edition, and we can also see that the obtained scores remained particularly low, indicating that the task remains challenging and difficult to address.

Specifically, if we focus on confusion matrices shown in Figure 7 for the submission received, we can see that the Baseline is over-predicting *Major* error spans, which gives a slight advantage regarding the F1 score since it leads to higher recall. This finding is consistent with higher precision obtained by HW-TSC submission as seen in the Appendix C, Table 17. We provide the confusion matrices for all language pairs in Appendix E.

Despite this, it is important to note that the methods submitted for Task 1 still seem to benefit from a multi-task approach that considers word-level information. Taking both these observations into account and looking towards future editions, it might be useful to redesign the task, aiming either at a different span representation that would perhaps attempt a better normalisation over different span lengths or deviate from the character level representation. Another alternative view would be to encourage methods that use error spans to support or interpret sentence-level quality (Leiter et al., 2023) or concentrate only on specific error types.

### 6.3 Task 3

#### 6.3.1 Automatic Evaluation

Automatic Post-Editing evaluation results are shown in Table 5. The submitted runs are ranked based on the average TER (case-sensitive) computed using human post-edits of the MT segments as a reference, which is the APE task's primary evaluation metric. To provide a broader view of the systems' performance, BLEU, chrF, and COMET results computed using the same references are also reported. As can be seen from the table, all submissions for English-Hindi outperform the baseline by a significant margin across all metrics, with TER reductions that are always statistically significant. the baseline. The best system is able to improve trans-

|        |                   | TER   | BLEU  | CHRF  | COMET  |
|--------|-------------------|-------|-------|-------|--------|
| En-Hi  | IT-Unbabel        | **27.08** | 58.38 | 73.45 | 0.8646 |
|        | HW-TSC_yjwsss     | **30.37** | 54.50 | 71.06 | 0.8514 |
|        | HW-TSC_zhaoxf4    | **31.32** | 52.74 | 69.83 | 0.8517 |
|        | BASELINE (MT)     | 46.36 | 39.28 | 59.48 | 0.8084 |
| En-Ta  | HW-TSC            | 24.24 | 69.64 | 82.36 | 0.9186 |
|        | IT-Unbabel        | 24.54 | 70.05 | 82.30 | 0.9163 |
|        | BASELINE (MT)     | 24.71 | 70.16 | 81.80 | 0.9137 |

Table 5: Official results for the WMT24 Quality Estimation **Task 3 QE-informed APE** English-Hindi and English-Tamil shared task – average TER (↓), BLEU (↑), chrF (↑), COMET (↑). Statistical significance test is computed for the primary metric (TER) *wrt.* the baseline and the significant results are highlighted in bold. Baseline systems are highlighted in grey.

lation quality by nearly 20.0 TER points. However, for English-Tamil, we observe that while all submissions performed slightly better than the baseline in terms of absolute scores across all metrics except BLEU, none of the systems show statistically significant gains compared to the baseline. As discussed in Section 3.4, this can be attributed to the combined effect of less repetitive data (between 1.1-2.0) compared to English-Hindi (between 2.7-3.3) and a stronger baseline (24.7 vs 46.4 TER), leaving less room for improvement.

### 6.3.2 Analysis: Systems' Behaviour

**Modified, improved and deteriorated sentences.** To better understand the behaviour of each APE system, we now turn toward the changes made by each system to the test instances. To this end, Table 6 shows, for each submitted run, the number of modified, improved and deteriorated sentences, as well as the overall system's precision (*i.e.,* the proportion of improved sentences out of the total number of modified instances for which improvement/deterioration is observed). It's worth noting that, as in the previous rounds, the number of sentences modified by each system is higher than the sum of the improved and the deteriorated ones. This difference is represented by modified sentences for which the corrections do not yield any TER variations.

As can be seen from Table 6, for English-Hindi, all submissions perform aggressive post-editing, with the top submission modifying 96.5% of the translations, where most of the modifications lead to improving the translation quality with a precision score of 84.56%. In contrast, for English-Tamil, all submissions adopt a conservative approach, limit-



Figure 8: Distribution of edit operations (insertions, deletions, substitutions and shifts) performed by the three primary submissions to the WMT24 APE English-Hindi shared task.

ing edits to 3.8%-4.8% of the test set. This aligns with our previous observations on task difficulty, driven by the higher MT baseline and the skewed TER distribution, with samples concentrated in the near-perfect translation range. In this challenging scenario, all submissions are able to improve the majority of modified translations with a precision score between 54%-59%.

**Edit operations.** Similar to previous rounds, we analysed systems' behaviour also in terms of the distribution of edit operations (insertions, deletions, substitutions and shifts) done by each system. This fine-grained analysis of how systems corrected the test set instances is obtained by computing the TER between the original MT output and the output of each primary submission taken as reference. As shown in Figures 8 and 9, similar to last year, differences in systems' behaviour are minimal. All of them are characterised by a large number of deletions, followed by insertions, shifts and substitutions. For English-Tamil, we observe a relatively lower proportion of shifts and substitutions compared to English-Hindi. This might indicate that English-Tamil might have more diverse APE outputs, which might be more challenging to evaluate with reference-based automatic metrics.

## 7 Evaluation on challenge sets

We received two submissions that we could evaluate on challenge sets: Pister Lab's submission, based on prompting Llama 3.1, and Unbabel's, based on CometKiwi. In Figure 10, we report the percentage of samples where the hyp translation is

|  | Systems | Modified | Improved | Deteriorated | Prec. |
|---|---|---|---|---|---|
| En-Hi | IT-Unbabel | 965 (96.5%) | 756 (78.35%) | 138 (14.30%) | 84.56 |
|  | HW-TSC_yjwsss | 952 (95.2%) | 688 (72.27%) | 171 (17.96%) | 80.09 |
|  | HW-TSC_zhaoxf4 | 665 (66.5%) | 532 (80.00%) | 85 (12.78%) | 86.22 |
| En-Ta | HW-TSC | 48 (4.8%) | 25 (52.08%) | 18 (37.50%) | 58.14 |
|  | IT-Unbabel | 38 (3.8%) | 19 (50.00%) | 16 (42.11%) | 54.29 |

Table 6: Number (raw and proportion) of test sentences modified, improved and deteriorated by each run submitted to the APE 2024 English-Hindi and English-Tamil sub-task. The "Prec." column shows systems' precision as the ratio between the number of improved sentences and the number of modified instances for which improvement/deterioration is observed (*i.e.,* Improved + Deteriorated).



Figure 9: Distribution of edit operations (insertions, deletions, substitutions and shifts) performed by the three primary submissions to the WMT24 APE English-Tamil shared task.

scored higher, lower, or is tied to the con hypothesis.[20] Please refer to Section 2.2 for details on constructing these translation pairs for each phenomenon.

**Detached translations and omissions** Out of all the phenomena studied, these two constitute the most critical errors. It is thus highly encouraging that both models perform perfectly across the two language pairs in consistently scoring the correct hyp translation higher than the erroneous con translation.

**Currency and date formatting** This category reveals interesting differences between the two models. Llama 3.1 shows a high tie rate, indicating

---

[20]Inspired by the analysis in Kocmi et al. (2024), we consider a tie with CometKiwi when the absolute difference between the scores of the hyp and con hypotheses is lower or equal to 0.1 points. For the Llama-based submission, for its more coarse-grained scoring range (more akin to a categorical distribution), we consider a tie when both translations receive the same score.

it often does not distinguish between original and localized formats. This suggests a more neutral stance towards formatting choices. In contrast, CometKiwi is more sensitive to these formats, behaving less predictably. Although, in most cases, it either prefers the source format or is indifferent to the localized format, there are some cases, in particular for en-es translations, where it does prefer the localized format that does not lexically match that found in the source text.

**Idioms** Llama 3.1 predominantly shows ties or a slight preference for non-literal, idiomatic renderings (hyp) that accurately convey the meaning in the target language. In contrast, CometKiwi's behavior is more varied and, perhaps surprisingly, often favors literal translations (con) even when they may not preserve the source text's meaning in the target language. This tendency towards literalness can be quite problematic in the context of idioms and other figurative texts, where meaning often diverges from word-for-word translations. One potential way to alleviate these trends is to train neural metrics with more diverse data that includes idiomatic and figurative language to improve their robustness.

**Word order** Here, Llama 3.1 shows a high rate of ties, suggesting that, similarly to what we found for the currency and date formatting phenomenon, it does not distinguish between monotonic translations that closely follow the source sentence order and non-monotonic translations that rearrange words while preserving meaning. This suggests that Llama 3.1's scoring may be more tied to the overall meaning of the translation. In contrast, CometKiwi demonstrates more preference for monotonic translations (hyp) across both language pairs, particularly for en-de. As such, CometKiwi appears to be more sensitive to word order, poten-

Figure 10: Share of instances in challenge sets where participant systems ranked the hypothesis translation higher than (green), lower than (salmon), or equal to (grey) the contrast. Results on en-de (top) and en-es (bottom).

tially favouring translations that maintain a structure closer to the source text. As a learned metric, this behaviour might be attributed to CometKiwi's training data, which may have contained more monotonic translations (more common among classical encoder-decoder NMT models that constitute most of the translations that the model has seen during training) than paraphrastic or non-literal ones (more prevalent among the more novel LLM-based translation approaches (Raunak et al., 2023)).

**Gender subset** In most instances, both systems score the hypothesis with the correct gender inflection higher. However, we noticed that some cases have ties, which we consider as errors: the model does not capture the difference in gender forms and wrongly assigns equal scores to the hypothesis and the contrast. Expectedly, this phenomenon is more present in Pister Lab's scores, as Llama 3.1 tends to

assign more coarse-grained assessments. In analyzing sources with non-overlapping content, Llama 3.1 exhibits a higher frequency of errors for male sources in en-de translation while demonstrating increased error rates for female sources in en-es. Conversely, CometKiwi maintains a comparable error rate across genders in both language pairs, with an elevated error rate in en-es translation overall. When examining sources with identical content differentiated only by gender (categorized as "overlapping"), we observed higher errors for female sources across all configurations, except for CometKiwi's performance in en-es.

**Closing remarks** Our analysis of Llama 3.1 and CometKiwi on various challenge sets reveals distinct behaviours and potential areas for improvement. Both models excel at identifying critical errors like detached translations and omissions.

However, they differ in their handling of formatting, idioms, and word order, with Llama 3.1—perhaps for the more discrete nature of its quality assessments—often showing neutrality (manifested through a large number of ties) and CometKiwi demonstrating more varied preferences, some of which are problematic (e.g., preference towards literalness in the translation of idiomatic expressions). Gender-related evaluations suggest potential biases in both systems, mainly due to scoring masculine and feminine gender inflections equally despite only one being correct. When controlling for the source content, we notice more errors for the instances mentioning a feminine referent in specific contexts. These findings indicate that both models display gender-dependent behaviour in source processing, warranting further investigation into potential model biases.

## 8 Discussion

In the following, we discuss the main findings of this year's shared task based on the goals we had previously identified for it.

**Large language models in Quality Estimation** In this edition, we observed an increased use of LLMs, not only in order to generate pseudo-data for training or as a complementary system –which was the trend in the previous year– but rather as the primary model to address a task. Indeed, across tasks, it was possible to observe the performance of encoder-based models that follow the predictor-estimator architecture (Kim et al., 2017), as well as models that relied on large decoder-based approaches, where the emphasis was more on prompt engineering or instruction tuning. This is in line with recent works (Huang et al., 2023; Fernandes et al., 2023a; Kocmi and Federmann, 2023; Vu et al., 2024; Hada et al., 2024) that suggest that multilingual LLMs can be prompted to predict the quality of a translation, given some tuning or in-context learning.

Looking at the results for Tasks 1 and 2, however, we can see that the methods that rely on LLMs are still outperformed by predictor-estimator-based systems, especially when it comes to predicting sentence-level scores. One key disparity, in this case, relates to the fact that methods relying on scores generated by such models lack the granularity of predictor-estimator architectures that treat the QE task as a regression and, hence, can differentiate better between different translations and

quality levels. Instead, LLMs tend to default to a smaller range of values (as we can also see in the ties detected in the analysis of Section 7 and Figures 10a and 10b. However, we can see that the LLM-based methods are closing the gap in terms of performance when compared to the predictor-estimator-based model for Task 2, which involves error detection. More importantly, LLM-based approaches perform on par and even outperform other methods for Task 3, which focuses on translation correction (APE). Thus, it seems that in the MT evaluation and correction family of tasks, there is potential for both LLMs and "traditional" neural systems. Potentially, more hybrid methods, i.e. methods that employ sentence-level quality scores predicted from encoder-based models to inform LLM decisions on error detection and correction, would lead to improved performance and could take the lead in future editions for the shared task.

**Role of QE signals in APE** Both participants in Task 3 used QE information to perform APE in alignment with the task objectives. Their approaches share similarities, as they both involve a final QE-driven selection step to choose between the original MT output and the generated APE hypothesis. One participant (HW-TSC) exploited QE information only for this final selection step, while the other (IT-Unbabel) integrated the two technologies more tightly by generating APE outputs with an LLM informed by free-text explanations for translation errors, which can be considered as proxies for QE predictions. Overall, despite being obtained with different degrees of QE integration, the evaluation results reinforce previous findings regarding the effectiveness of combined QEAPE and approaches for enhancing MT output (Chatterjee et al., 2017; Deoghare et al., 2023).

## 9 Conclusions

This year's edition of the QE Shared Task introduced two key new elements besides fresh test sets: (1) A new task on QE-informed APE, motivating participants to consider the QE scores to improve the generated MT corrections and (2) an updated challenge set for En-De and En-Es language pairs to help analyse the behaviour and robustness of submitted models for different phenomena such as gender bias, idiomatic expressions, handling of numerical entities, hallucinations, and word order changes.

We found that overall QE performance is consis-

tently high across languages on the sentence level. Still, there is ample room for improvement regarding fine-grained error span detection. The addition of quality informed APE sub-task made it easier for participants to leverage their QE system for the APE task, achieving significant gains for *en-hi* and marginal (non-significant) gains on *en-ta* language pairs. In addition, we found that approaches that employ LLMs still have some way to go in competing on correlations with human scores at the sentence level but can provide competitive solutions for error span detection and QE-informed APE tasks.

In future iterations, we aim to redefine meaningful fine-grained QE tasks, targeting attainable error detection that can help detect critical errors, explain predicted quality, and better inform APE systems. Additionally, we intend to expand further the provided resources to aid the finer grained analysis of model behaviour, as it was discussed in Section 7.

## 10 Ethical Considerations

Post-editing, MQM, and DA annotations in this paper are carried out by professional translators. They are all paid at professional rates. In creating the gender subset, we drew examples from MT-GenEval (Currey et al., 2022), a corpus where gender is treated as a binary variable. We recognize that gender identities exist on a spectrum, going beyond just the masculine-feminine dichotomy. Our intention is to expand the evaluation of gender-related aspects to include more inclusive forms of machine translation.

Organisers from Unbabel and IT have submitted to this task without using prior access to test sets or any insider information.

## References

Hervé Abdi. 2007. The bonferroni and šidák corrections for multiple comparisons. *Encyclopedia of measurement and statistics*, 3:103–107.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88, Online. Association for Computational Linguistics.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Frederic Blain, Chrysoula Zerva, Ricardo Rei, Nuno M. Guerreiro, Diptesh Kanojia, José G. C. de Souza, Beatriz Silva, Tânia Vaz, Yan Jingxuan, Fatemeh Azadi, Constantin Orasan, and André Martins. 2023. Findings of the WMT 2023 shared task on quality estimation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 629–653, Singapore. Association for Computational Linguistics.

Aljoscha Burchardt. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.

Rajen Chatterjee, Christian Federmann, Matteo Negri, and Marco Turchi. 2019. Findings of the WMT 2019

shared task on automatic post-editing. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 11–28, Florence, Italy. Association for Computational Linguistics.

Rajen Chatterjee, Markus Freitag, Matteo Negri, and Marco Turchi. 2020. Findings of the WMT 2020 shared task on automatic post-editing. In *Proceedings of the Fifth Conference on Machine Translation*, pages 646–659, Online. Association for Computational Linguistics.

Rajen Chatterjee, Matteo Negri, Raphael Rubino, and Marco Turchi. 2018a. Findings of the WMT 2018 shared task on automatic post-editing. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 710–725, Belgium, Brussels. Association for Computational Linguistics.

Rajen Chatterjee, Matteo Negri, Marco Turchi, Frédéric Blain, and Lucia Specia. 2018b. Combining quality estimation and automatic post-editing to enhance machine translation output. In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 26–38.

Rajen Chatterjee, Matteo Negri, Marco Turchi, Marcello Federico, Lucia Specia, and Frédéric Blain. 2017. Guiding neural machine translation decoding with external knowledge. In *Proceedings of the Second Conference on Machine Translation, Volume 1: Research Papers*, Copenhagen, Denmark. Association for Computational Linguistics.

Rajen Chatterjee, Marion Weller, Matteo Negri, and Marco Turchi. 2015. Exploring the planet of the APEs: a comparative study of state-of-the-art methods for MT automatic post-editing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 156–161, Beijing, China. Association for Computational Linguistics.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Anna Currey, Maria Nadejde, Raghavendra Reddy Pappagari, Mia Mayer, Stanislas Lauly, Xing Niu, Benjamin Hsu, and Georgiana Dinu. 2022. MT-GenEval: A counterfactual and contextual dataset for evaluating gender accuracy in machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4287–4299, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Sourabh Deoghare, Diptesh Kanojia, Fred Blain, Tharindu Ranasinghe, and Pushpak Bhattacharyya.

2023. Quality estimation-assisted automatic post-editing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1686–1698, Singapore. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Taisei Enomoto, Hwichan Kim, Tosho Hirasawa, Yoshinari Nagai, Ayako Sato, Kyotaro Nakajima, and Mamoru Komachi. 2024. TMU-HIT at MLSP 2024: How well can GPT-4 tackle multilingual lexical simplification? In *Proceedings of the 19th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2024)*, pages 590–598, Mexico City, Mexico. Association for Computational Linguistics.

Alexander Richard Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. Qafacteval: Improved qa-based factual consistency evaluation for summarization. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2587–2601.

Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André Martins, Graham Neubig, Ankush Garg, Jonathan Clark, Markus Freitag, and Orhan Firat. 2023a. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1066–1083, Singapore. Association for Computational Linguistics.

Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André FT Martins, Graham Neubig, Ankush Garg, Jonathan H Clark, Markus Freitag, and Orhan Firat. 2023b. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. *arXiv preprint arXiv:2308.07286*.

Marina Fomicheva, Shuo Sun, Erick Fonseca, Chrysoula Zerva, Frédéric Blain, Vishrav Chaudhary, Francisco Guzmán, Nina Lopatina, Lucia Specia, and André F. T. Martins. 2022. MLQE-PE: A multilingual quality estimation and post-editing dataset. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4963–4974, Marseille, France. European Language Resources Association.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej

Bojar. 2021b. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774, Online. Association for Computational Linguistics.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.

Rishav Hada, Varun Gumma, Adrian Wynter, Harshita Diddee, Mohamed Ahmed, Monojit Choudhury, Kalika Bali, and Sunayana Sitaram. 2024. Are large language model-based evaluators the solution to scaling up multilingual evaluation? In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1051–1070.

Hui Huang, Shuangzhi Wu, Xinnian Liang, Bing Wang, Yanrui Shi, Peihao Wu, Muyun Yang, and Tiejun Zhao. 2023. Towards making the most of llm for translation quality estimation. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 375–386. Springer.

Raisa Islam and Owana Marzia Moushi. 2024. Gpt-4o: The cutting-edge advancement in multimodal llm. *Authorea Preprints*.

Hyun Kim, Jong-Hyeok Lee, and Seung-Hoon Na. 2017. Predictor-estimator using multilevel task learning with stack propagation for neural quality estimation. In *Proceedings of the Second Conference on Machine Translation*, pages 562–568, Copenhagen, Denmark. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. Gemba-mqm: Detecting translation quality error spans with gpt-4. In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775.

Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024. Navigating the metrics maze: Reconciling score magnitudes and accuracies.

Anoop Kunchukuttan, Pratik Mehta, and Pushpak Bhattacharyya. 2018. The IIT Bombay English-Hindi parallel corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Christoph Leiter, Juri Opitz, Daniel Deutsch, Yang Gao, Rotem Dror, and Steffen Eger. 2023. The Eval4NLP 2023 shared task on prompting large language models as explainable metrics. In *Proceedings of the 4th Workshop on Evaluation and Comparison of NLP Systems*, pages 117–138, Bali, Indonesia. Association for Computational Linguistics.

Yuang Li, Chang Su, Ming Zhu, Mengyao Piao, Xinglin Lyu, Min Zhang, and Hao Yang. 2023. Hw-tsc 2023 submission for the quality estimation shared task. In *Proceedings of the Eigth Conference on Machine Translation*, Singapore. Association for Computational Linguistics.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.

Arle Lommel, Aljoscha Burchardt, Maja Popović, Kim Harris, Eleftherios Avramidis, and Hans Uszkoreit. 2014. Using a new analytic measure for the annotation and analysis of mt errors on real data. In *Proceedings of the 17th Annual conference of the European Association for Machine Translation*, pages 165–172.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Vikas Raunak, Arul Menezes, Matt Post, and Hany Hassan Awadalla. 2023. Do gpts produce less literal translations?

Ricardo Rei, Nuno M Guerreiro, Daan van Stigt, Marcos Treviso, Luísa Coheur, José GC de Souza, André FT Martins, et al. 2023. Scaling up cometkiwi: Unbabel-ist 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848.

Ricardo Rei, Marcos Treviso, Nuno M Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José GC de Souza, Taisiya Glushkova, Duarte M Alves, Alon Lavie, et al. 2022. Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task. *WMT 2022*, page 634.

Ananya Sai B, Tanay Dixit, Vignesh Nagarajan, Anoop Kunchukuttan, Pratyush Kumar, Mitesh M. Khapra, and Raj Dabre. 2023. IndicMT eval: A dataset to meta-evaluate machine translation metrics for Indian languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14210–14228, Toronto, Canada. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Lucia Specia, Frédéric Blain, Marina Fomicheva, Chrysoula Zerva, Zhenhao Li, Vishrav Chaudhary, and André F. T. Martins. 2021. Findings of the WMT 2021 shared task on quality estimation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 684–725, Online. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Marcos Treviso, Nuno M Guerreiro, Sweta Agrawal, Ricardo Rei, José Pombal, Tania Vaz, Helena Wu, Beatriz Silva, Daan van Stigt, and André FT Martins. 2024. xtower: A multilingual llm for explaining and correcting translation errors. *arXiv preprint arXiv:2406.19482*.

Marco Turchi, Matteo Negri, and Marcello Federico. 2013. Coping with the subjectivity of human judgements in MT quality estimation. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 240–251, Sofia, Bulgaria. Association for Computational Linguistics.

Tu Vu, Kalpesh Krishna, Salaheddin Alzubi, Chris Tar, Manaal Faruqui, and Yun-Hsuan Sung. 2024. Foundational autoraters: Taming large language models for better automatic evaluation. *arXiv e-prints*, pages arXiv–2407.

Evan J. Williams. 1959. *Regression Analysis*, volume 14. Wiley, New York, USA.

Alexander Yeh. 2000. More Accurate Tests for the Statistical Significance of Result Differences. In *Coling-2000: the 18th Conference on Computational Linguistics*, pages 947–953, Saarbrücken, Germany.

Jiawei Yu, Xiaofeng Zhao, Min Zhang, Yanqing Zhao, Yuang Li, Chang Su, Xiaosong Qiao, Miaomiao Ma, and Hao Yang. 2024. Hw-tsc's participation in the wmt 2024 qeape task. *In Proceedings of the Ninth Conference on Machine Translation (WMT24)*.

Emmanouil Zaranis, Giuseppe Attanasio, Sweta Agrawal, and André F. T. Martins. 2024. Watching the watchers: Exposing gender disparities in machine translation quality estimation.

Chrysoula Zerva, Frédéric Blain, Ricardo Rei, Piyawat Lertvittayakumjorn, José G. C. de Souza, Steffen Eger, Diptesh Kanojia, Duarte Alves, Constantin Orăsan, Marina Fomicheva, André F. T. Martins, and Lucia Specia. 2022. Findings of the WMT 2022 shared task on quality estimation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 69–99, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

## A Hyper-parameters of pre-trained baseline models for Task 1 and Task 2 Quality Estimation

| Hyper-parameter | T1 Sentence-level<br>COMETKIWI-DA-22 | T2 Fine-grained<br>COMETKIWI-MULTITASK-22 |
|---|---|---|
| Encoder Model | XLM-RoBERTa (large) | XLM-RoBERTa (large) |
| Optimizer | Adam (default parameters) | Adam (default parameters) |
| n frozen epochs | 0.3 | 0.3 |
| Keep embeddings frozen | True | True |
| Learning rate | 3e-05 and 1e-05 | 3e-06 and 1e-05 |
| Batch size | 4 | 4 |
| Loss function | MSE and CE | MSE and CE |
| Dropout | 0.15 | 0.1 |
| FP precision | 32 | 32 |
| Feed-Forward hidden units | [2048, 1024] | [3072, 1024] |
| Word weights | [0.3, 0.7] | [0.1, 0.9] |
| Feed-Forward activation | Tanh | – |
| Language prefix | False | False |

Table 7: Hyper-parameters of both the CometKiwi models used as baselines for Task 1 Quality Estimation.

## B  Official Results of the WMT24 Quality Estimation Task 1 Sentence-level

Tables 8, 9, 10, 11, 12, 13, 14 and 15 show the results for all language pairs and the multilingual variants, ranking participating systems best to worst using Spearman correlation as primary key for each of these cases.

| Model | Spearman | Pearson | Kendall |
|---|---|---|---|
| Unbabel | 0.553 | 0.438 | 0.410 |
| BASELINE | 0.520 | 0.474 | 0.382 |
| Pister Labs | 0.452 | 0.378 | 0.354 |

Table 8: Official results of the WMT24 Quality Estimation Task 1 Sentence-level **Multilingual** (average over all language pairs). Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

| Model | Spearman | Pearson | Kendall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| BASELINE ● | 0.514 | 0.050 | 0.397 | 2,260,734,705 | 569,330,715 | 1 |
| Pister Labs ● | 0.513 | 0.114 | 0.455 | 1,400,000,000 | 70,000,000,000 | 1 |
| Unbabel ● | 0.512 | 0.037 | 0.393 | 42,868,104,221 | 10,716,932,147 | 6 |

Table 9: Official results of the WMT24 Quality Estimation Task 1 Sentence-level for **Engligh-German (MQM)**. Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

| Model | Spearman | Pearson | Kendall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| Unbabel ● | 0.345 | 0.116 | 0.257 | 42,868,104,221 | 10,716,932,147 | 6 |
| BASELINE ● | 0.340 | 0.197 | 0.253 | 2,260,734,705 | 569,330,715 | 1 |
| Pister Labs | 0.282 | 0.104 | 0.215 | 1,400,000,000 | 70,000,000,000 | 1 |

Table 10: Official results of the WMT24 Quality Estimation Task 1 Sentence-level for **English-Spanish (MQM)**. Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

| Model | Spearman | Pearson | Kendall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| BASELINE ● | 0.441 | 0.223 | 0.328 | 2,260,734,705 | 569,330,715 | 1 |
| Unbabel | 0.412 | 0.065 | 0.318 | 42,868,104,221 | 10,716,932,147 | 6 |
| Pister Labs | 0.363 | 0.142 | 0.300 | 1,400,000,000 | 70,000,000,000 | 1 |

Table 11: Official results of the WMT24 Quality Estimation Task 1 Sentence-level for **English-Hindi (MQM)**. Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

| Model | Spearman | Pearson | Kendall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| TMU-HIT ● | 0.739 | 0.760 | 0.547 | - | - | 1 |
| HW-TSC ● | 0.719 | 0.783 | 0.531 | 2,387,827,161 | 596,896,035 | 8 |
| Unbabel | 0.714 | 0.679 | 0.524 | 42,868,104,221 | 10,716,932,147 | 6 |
| BASELINE | 0.678 | 0.771 | 0.497 | 2,260,734,705 | 569,330,715 | 1 |
| Pister Labs | 0.564 | 0.536 | 0.443 | 1,400,000,000 | 70,000,000,000 | 1 |

Table 12: Official results of the WMT24 Quality Estimation Task 1 Sentence-level for **English-Hindi (DA)**. Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

| Model | **Spearman** | Pearson | Kendall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| TMU-HIT ● | 0.713 | 0.808 | 0.531 | - | - | 1 |
| Unbabel ● | 0.703 | 0.751 | 0.514 | 42,868,104,221 | 10,716,932,147 | 6 |
| HW-TSC | 0.686 | 0.757 | 0.500 | 2,387,827,161 | 596,896,035 | 8 |
| BASELINE | 0.661 | 0.776 | 0.486 | 2,260,734,705 | 569,330,715 | 1 |
| Pister Labs | 0.587 | 0.716 | 0.366 | 1,400,000,000 | 70,000,000,000 | 1 |

Table 13: Official results of the WMT24 Quality Estimation Task 1 Sentence-level **English-Gujarati (DA)**. Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

| Model | **Spearman** | Pearson | Kendall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| Unbabel ● | 0.510 | 0.719 | 0.363 | 42,868,104,221 | 10,716,932,147 | 6 |
| HW-TSC ● | 0.482 | 0.643 | 0.340 | 2,387,827,161 | 596,896,035 | 8 |
| TMU-HIT | 0.465 | 0.550 | 0.329 | - | - | 1 |
| BASELINE | 0.414 | 0.716 | 0.294 | 2,260,734,705 | 569,330,715 | 1 |
| Pister Labs | 0.379 | 0.535 | 0.304 | 1,400,000,000 | 70,000,000,000 | 1 |

Table 14: Official results of the WMT24 Quality Estimation Task 1 Sentence-level **English-Telugu (DA)**. Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

| Model | **Spearman** | Pearson | Kendall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| HW-TSC ● | 0.683 | 0.719 | 0.506 | 2,387,827,161 | 596,896,035 | 8 |
| Unbabel ● | 0.675 | 0.702 | 0.499 | 42,868,104,221 | 10,716,932,147 | 6 |
| TMU-HIT | 0.603 | 0.664 | 0.445 | - | - | 1 |
| BASELINE | 0.592 | 0.584 | 0.419 | 2,260,734,705 | 569,330,715 | 1 |
| Pister Labs | 0.478 | 0.503 | 0.366 | 1,400,000,000 | 70,000,000,000 | 1 |

Table 15: Official results of the WMT24 Quality Estimation Task 1 Sentence-level **English-Tamil (DA)**. Teams marked with "●" are the winners, as they are not significantly outperformed by any other system according to the Williams Significance Test (Williams, 1959). Baseline systems are highlighted in grey.

# C  Official Results of the WMT24 Quality Estimation Task 2 Fine grained Error Detection

Tables 16, 17, 18 and 19 show the results for all language pairs and the multilingual variant, ranking participating systems best to worst using F1-score as primary key for each of these cases.

| Model | **F1-score** | Precision | Recall |
|---|---|---|---|
| BASELINE | 0.278 | 0.220 | 0.427 |
| HW-TSC | 0.227 | 0.203 | 0.268 |

Table 16: Official results of the WMT24 Quality Estimation Task 2 Fine grained Error Detection **Multilingual** (average over all language pairs). The winning submission is indicated by a ●. Baseline systems are highlighted in grey.

| Model | **F1-score** | Precision | Recall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| BASELINE | 0.192 | 0.127 | 0.394 | 2,260,743,915 | 569,309,780 | 1 |
| HW-TSC | 0.178 | 0.175 | 0.181 | 2,409,244,995 | 2,280,000,000 | 1 |

Table 17: Official results of the WMT24 Quality Estimation Task 2 Fine grained Error Detection **English-German (MQM)**. The winning submission is indicated by a ●. Baseline systems are highlighted in grey.

| Model | **F1-score** | Precision | Recall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| BASELINE | 0.161 | 0.106 | 0.337 | 2,260,743,915 | 569,309,780 | 1 |
| HW-TSC | 0.151 | 0.106 | 0.261 | 2,409,244,995 | 2,280,000,000 | 1 |

Table 18: Official results of the WMT24 Quality Estimation Task 2 Fine grained Error Detection **English-Spanish (MQM)**. The winning submission is indicated by a ●. Baseline systems are highlighted in grey.

| Model | **F1-score** | Precision | Recall | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|
| BASELINE | 0.481 | 0.428 | 0.551 | 2,260,743,915 | 569,309,780 | 1 |
| HW-TSC | 0.362 | 0.329 | 0.401 | 2,409,244,995 | 2,280,000,000 | 1 |

Table 19: Official results of the WMT24 Quality Estimation Task 2 Fine grained Error Detection **English-Hindi (MQM)**. The winning submission is indicated by a ●. Baseline systems are highlighted in grey.

# D  Official Results of the WMT24 Quality Estimation Task 3 Quality-informed APE

Tables 20 and 21 show the results for all language pairs, ranking participating systems from best to worst using TER as the primary key for each of these cases.

| Model | TER | BLEU | ChrF | COMET | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|---|
| IT-Unbabel ● | 27.08 | 58.38 | 73.45 | 0.8646 | 28,991,029,248 | 7,000,000,000 | 1 |
| HW-TSC ● | 31.32 | 52.74 | 69.83 | 0.8517 | 1,265,490,783 | 99,388,416 | 1 |
| BASELINE | 46.36 | 39.28 | 59.48 | 0.8084 | - | - | - |

Table 20: Official results of the WMT24 Quality Estimation Task 3 Quality-informed APE **English-Hindi (DA)**. The winning submission is indicated by a ●. Baseline systems are highlighted in grey.

| Model | TER | BLEU | ChrF | COMET | Disk footprint (B) | # Model params | Ensemble |
|---|---|---|---|---|---|---|---|
| HW-TSC | 24.24 | 69.64 | 82.36 | 0.9186 | 1,265,490,783 | 99,388,416 | 1 |
| IT-Unbabel | 24.54 | 70.05 | 82.30 | 0.9163 | 28,991,029,248 | 7,000,000,000 | 1 |
| BASELINE | 24.71 | 70.16 | 81.80 | 0.9137 | - | - | - |

Table 21: Official results of the WMT24 Quality Estimation Task 3 Quality-informed APE **English-Tamil (DA)**. The winning submission is indicated by a ●. Baseline systems are highlighted in grey.

# E  Confusion Matrices for Task 2

We present below the confusion matrices for Major and Minor error span prediction between HW-TSC and the Baseline, for each language pair. We can see that overall HW-TSC targets precision, being more conservative in error span prediction, while the Baseline model greedily predicts major errors.

(a) HWTSC - Major

(b) HWTSC -Minor

(c) Baseline - Major

(d) Baseline - Minor

Figure 11: Confusion matrices for Task 2 English-German, comparing Minor and Major predictions between the Baseline system and the HWTSC one.

(a) HWTSC - Major

(b) HWTSC -Minor

(c) Baseline - Major

(d) Baseline - Minor

Figure 12: Confusion matrices for Task 2 English-Spanish, comparing Minor and Major predictions between the Baseline system and the HWTSC one.

(a) HWTSC - Major

(b) HWTSC -Minor

(c) Baseline - Major

(d) Baseline - Minor

Figure 13: Confusion matrices for Task 2 English-Hindi, comparing Minor and Major predictions between the Baseline system and the HWTSC one.

# Findings of the WMT 2024 Shared Task
# of the Open Language Data Initiative

**Laurie Burchell**[*]
University of Edinburgh

**Jean Maillard**[*]
Meta FAIR

**Antonios Anastasopoulos**
George Mason University

**Christian Federman**
Microsoft

**Philipp Koehn**
Johns Hopkins University

**Skyler Wang**
McGill University

Correspondence: info@oldi.org

## Abstract

We present the results of the WMT 2024 shared task of the Open Language Data Initiative. Participants were invited to contribute to the FLORES+ and MT Seed multilingual datasets, two foundational open resources that facilitate the organic expansion of language technology's reach. We accepted ten submissions covering 16 languages, which extended the range of languages included in the datasets and improved the quality of existing data.

## 1 Introduction

Machine translation research has advanced at breakneck speed in recent years (Kocmi et al., 2023). That said, progress made in translation quality has largely been directed at high-resource languages, leaving many languages behind. More recently, the focus has shifted towards under-served languages (also called low-resource) (Haddow et al., 2022). Foundational, high-coverage datasets have made it easier to develop and evaluate language technologies for a growing number of languages. Given the high impact of these components, extending such datasets becomes imperative.

The aim of the WMT 2024 shared task of the Open Language Data Initiative (OLDI) is to empower language communities to contribute such key datasets. In particular, we solicited contributions to the MT evaluation dataset FLORES+ and the MT Seed dataset. Additionally, we also solicited other high-quality, human-verified monolingual text datasets in under-resource languages. This builds on previous work to create these datasets and extend machine translation (MT) models and evaluation tools to more languages (Guzmán et al., 2019; Goyal et al., 2022; NLLB Team et al., 2024; Maillard et al., 2023).

We accepted ten submissions to the task, and the data contributed covered 16 languages. We re-

quired all contributions to be released under open licenses so that they can be useful to as many community members as possible. We make the data available online and encourage future work to build on these foundational datasets even further.[1]

## 2 Related Work

In recent years, there has been a growing recognition of the need for high-quality, representative datasets to broaden access to language technologies across a more diverse range of languages. Several initiatives have emerged to address this need.

In machine translation, the FLORES family of datasets (Guzmán et al., 2019; Goyal et al., 2022; NLLB Team et al., 2024) and NTREX-128 (Federmann et al., 2022) have provided the research community with massively multilingual, professionally translated benchmark data that is open source; while NLLB-Seed (Maillard et al., 2023; NLLB Team et al., 2024) played a similar role but focused on training data. Since the release of these resources, several authors have provided coverage for new languages (Gala et al., 2023; Doumbouya et al., 2023; Aepli et al., 2023) or even extended the datasets to the speech modality (Conneau et al., 2022).

Thanks to the availability of higher-quality data for an increasingly larger number of languages, recent language identification models have been able to expand coverage. Projects such as AfroLID (Adebara et al., 2022) and OpenLID (Burchell et al., 2023) improved upon pre-existing models by a careful curation and auditing of existing data sources; while LIMIT (Agarwal et al., 2023) further expanded data coverage and performance by creating and releasing a new high-quality corpus.

Several crowdsourced projects have proven invaluable as a source of knowledge for under-served

---

[*]Equal contribution

[1]https://huggingface.co/openlanguagedata

languages. The Tatoeba project,[2] not designed explicitly for language technologies but as a language learning aid, provides a large database of aligned multilingual sentences. Mozilla Common Voice (Ardila et al., 2020) has enabled communities to build open-source ASR corpora for their own language and counts over 160 languages to date. The Aya initiative (Singh et al., 2024) has created the largest instruction finetuning dataset for large language models.

## 3 Datasets: FLORES+ and MT Seed

### 3.1 FLORES+

One of the biggest challenges in extending effective natural language processing (NLP) to underserved languages is a lack of high-quality, high-coverage evaluation benchmarks. In particular, few benchmarks are suitable for evaluating multilingual translation, since this requires many-to-many alignment between different languages in the evaluation dataset.

The FLORES family of datasets was released to address this problem. While the first iteration of this dataset covered only three languages (Guzmán et al., 2019), following iterations increased coverage to 101 languages (FLORES-101, Goyal et al., 2022) and finally to over 200 languages as part of the "No Language Left Behind" project (FLORES-200, NLLB Team et al., 2024). The current iteration of this dataset set is managed by OLDI, and we refer to it as FLORES+ to disambiguate between the original datasets and this new actively developed version.

FLORES+ consists of sentences extracted from English Wikinews, Wikijunior, and Wikivoyage: 997 for the dev split and 1012 for the devtest split.[3] These were then professionally translated into each language (almost universally from English) and underwent quality assessment and adjustment as necessary. The fact that all sentences in all languages are translations of each other means that they can be used for any-way multilingual evaluation.

### 3.2 MT Seed

The MT Seed dataset (previously NLLB Seed) was created as a source of "starter data" for languages without publicly-available high-quality bitext in sufficient quantity for training NLP models (NLLB

Team et al., 2024, p.23). Previous work showed that employing the relatively small amount of high-quality data in MT Seed for training models had a significant impact on performance even when larger but lower quality corpora are used (Maillard et al., 2023). By extending MT Seed, OLDI aims to improve the quality of NLP applications for underserved languages by providing an initial source of reliable training data.

MT Seed consists of around 6000 sentences sampled from the Wikipedia articles listed in English Wikimedia's "List of articles every Wikipedia should have". These were professionally translated into each of the 38 languages covered by the first iteration of this dataset (39 if including English). Since this dataset is intended as a source of training data rather than evaluation, it did not undergo the quality assurance as the FLORES family of datasets.

## 4 Shared Task Definition

The goal of this shared task was to expand the open datasets managed by OLDI. Primarily, we solicited contributions to FLORES+ and MT Seed (described in Section 3), which could be either fixes to existing data or entirely new translations. It also accepted other high-quality, human-verified monolingual text datasets in under-resource languages.

### 4.1 Contributing to FLORES+ and MT Seed

To contribute to FLORES+ and MT Seed, we encouraged participants to translate from English into the target language so as to follow the original standard FLORES-200 workflow (NLLB Team et al., 2024, p.21). We required that translations were performed by qualified, native speakers of the target language and that translators acknowledged our translation guidelines (Appendix A). We strongly encouraged the verification of the data by at least one additional native speaker.

The acceptability of machine-translated content varied between the two datasets. Since the FLORES+ dataset is used to evaluate MT systems, new contributions must be entirely human-translated. Using or even referencing MT output was not allowed, including post-editing. However, post-edited MT content was allowed for contributions to MT Seed, provided all content was verified manually. This was done because MT Seed is intended for training rather than evaluation and, therefore, has less stringent translation requirements.

---

[2] https://www.tatoeba.org
[3] The separate blind test set, originally developed by Meta, is not managed by OLDI and is not part of FLORES+.

Participants were encouraged to provide experimental validation to demonstrate the quality of their submitted datasets. Due to the heterogeneous nature of submissions, we left the exact nature of the experimental validation up to the participants, though we gave some suggestions. For example, MT Seed data contributions could train a simple MT model and evaluate it on FLORES+.

All submissions were labeled with the same standardized language codes used throughout OLDI. These are made up of three parts, separated by underscores:

- An ISO 639-3 language code. Macrolanguage codes must not be used if a more specific code is possible: e.g., cmn, yue, wuu, etc., rather than zho.

- An ISO 15924 script code

- A Glottocode identifying the specific language variety.

For example, `apc_Arab_sout3123` indicates South Levantine Arabic written in the Arabic script.

All submissions were accompanied by a dataset card summarizing key facts about the data and how it was created. This is critical to foster informed and responsible use of the submitted data (Pushkarna et al., 2022). Submitted datasets were required to be released under the open CC BY-SA 4.0 license to match FLORES+ and MT Seed.

## 4.2 Contributing other monolingual data

Contributions of monolingual data had similar requirements to those for FLORES+ or MT Seed. The aim was to collect high-quality, human-verified monolingual text in multiple under-served languages for training NLP tools and systems. Synthetic data of any kind was not allowed. Parallel datasets were excluded from the scope of the shared task to not conflict with existing corpus-building efforts like Opus (Tiedemann, 2009).

For FLORES+ and MT Seed, submissions were encouraged to be manually verified by native speakers of the target language. All submissions needed to be accompanied by a data card and released under an open license (allowing free research use as a minimum).

## 5 Submissions

There were 24 expressions of interest in the shared task, and we ultimately accepted 10 papers. Table 1

summarizes the data submitted. We describe each submission in the following section.

**Abdulmumin et al. (2024)** contributed an improved version of the FLORES+ datasets for Hausa, Northern Sotho (Sepedi), Xitsonga, and isiZulu. They carried out error analysis of the datasets for the four languages and found problems such as poor translation of named entities, incorrect handling of morphological changes, a lack of consistency in vocabulary, and poor handling of borrowed terms. The Hausa dataset was particularly weak, with evidence that it was built upon Google Translate outputs. The participants corrected the translations following the guidelines in the shared task description and evaluated the alterations to the dataset using a range of metrics.

**Ahmed et al. (2024)** contributed a translation of MT Seed into the Bangla variety of Bangla/Bengali, an Indo-Aryan language that is the official language of Bangladesh and the state of West Bengal in India (as well as others). The dataset was translated by a native speaker with translation experience, per the OLDI translation guidelines. They validated the quality of their dataset by fine-tuning a range of pre-trained multilingual models on their generated translations and compared performance with the same pre-trained models fine-tuned on different but comparable datasets. They found that the models pre-trained on their translation of MT Seed showed the best performance after controlling for dataset size.

**Ali et al. (2024)** produced a translation of the FLORES+ dataset into the Central variety of Emakhuwa, a Bantu language spoken primarily in Mozambique. They verified their translation by using a second translator to revise the work of the first, followed by quality assessment involving three raters using a Direct Assessment pipeline. The participants conducted several experiments to benchmark current progress in Emakhuwa–Portuguese MT. They found that a lack of standardized orthography remains a challenge for Emakhuwa MT, though multiple reference translations can help with this issue.

**Cols (2024)** released Seed-CAT, an open-source web application specifically designed to assist human translators in translating MT Seed dataset files.[4] Using Seed-CAT, they produced a trans-

---

| Contributors | Type of contribution | Languages(s) |
|---|---|---|
| Abdulmumin et al. (2024) | FLORES+ (corrected) | Hausa, Northern Sotho (Sepedi), Xitsonga, isiZulu. |
| Ahmed et al. (2024) | MT Seed | Bangla/Bengali |
| Ali et al. (2024) | FLORES+ (new) | Emakhuwa |
| Cols (2024) | MT Seed (new) and CAT tool | Spanish (Latin American) |
| Ferrante (2024) | MT Seed (new) | Italian |
| Gordeev et al. (2024) | FLORES+ (new) | Erzya |
| Kuzhuget et al. (2024) | FLORES+ (new) | Tuvan |
| Mamasaidov and Shopulatov (2024) | FLORES+ devtest (new) | Karakalpak |
| Perez-Ortiz et al. (2024) | FLORES+ (new and corrected) | Aragonese, Aranese, Asturian, Valencian |
| Yu et al. (2024) | FLORES+ (new) | Wu Chinese |

Table 1: A summary of all accepted contributions to the WMT 2024 Shared Task of the Open Language Data Initiative.

lation of MT Seed into Latin American Spanish. To validate their dataset's quality, they trained an English–Spanish MT model using the MT Seed data and compared its performance to models trained to translate between English and three Italic languages using existing MT Seed data. They found similar performance, suggesting that quality was similar to existing data in MT Seed.

**Ferrante (2024)** contributed a translation of MT Seed into Italian, building on a previous translation by Haberland et al. (2024). For this submission, the existing post-edited machine translation was reviewed and amended by two native speakers. The dataset was verified by training an Italian–Ligurian MT system and finding comparable results to those of Haberland et al. (2024).

**Gordeev et al. (2024)** contributed a translation of FLORES+ into Erzya, a Finno-Ugric language spoken primarily in Russia. As part of their work, they created a set of neologisms to aid future translators working in the digital space. They used their FLORES+ translation to evaluate the quality of existing English–Erzya and Russian–Erzya MT systems and train new competitive models for translating these language pairs.

**Kuzhuget et al. (2024)** translated the FLORES+ dataset from Russian into the Central dialect of Tuvan, a Turkic language primarily spoken in the Republic of Tuva in South Central Siberia, Russia. The team of translators worked from guidelines prepared in Russian to ensure consistent and high-quality translation.

**Mamasaidov and Shopulatov (2024)** contributed a translation of FLORES+ devtest split into Karakalpak, a Turkic language primarily spoken in the Republic of Karakalpakstan, which is

an autonomous region within Uzbekistan. In addition, they also released a training dataset containing 100,000 sentence pairs for each of the language pairs: Uzbek–Karakalpak, Russian–Karakalpak, and English–Karakalpak. They carried out MT experiments using their datasets, releasing the trained models for further research.

**Perez-Ortiz et al. (2024)** contributed translations of FLORES+ into four Romance languages spoken in Spain: specifically new datasets for Aragonese, Aranese, and Valencian, and a corrected dataset for Asturian. The datasets were used as part of the evaluation of a shared task on MT from Spanish to low-resource languages of Spain (Sánchez-Martínez et al., 2024). Even though post-edited MT was used in the creation of these datasets, they were exceptionally accepted due to their use in a major shared task with the use of post-editing flagged in the metadata.

**Yu et al. (2024)** contributed a translation of FLORES+ into the Chongming dialect of Wu Chinese. The translation was done by two native speakers and checked by a third. Since Wu Chinese is typically colloquial while FLORES+ contains relatively formal text, the translators examined online written content and asked for community guidance about translations on fora to arrive at the best translations. To validate their dataset, the participants ran a three-way language identification task between Wu Chinese, Mandarin Chinese, and Yue Chinese. Their language identification model could distinguish between the three language varieties with high accuracy, though there was some confusion between Mandarin and Wu Chinese.

## 6 Discussion

Despite recent releases of state-of-the-art large-scale models (NLLB Team et al., 2024) and the growing attention directed at speech and sign language translations (Seamless Communication et al., 2023a,b; Rust et al., 2024), the work on text-based MT remains ongoing. This is particularly true for many of the world's under-served languages, which compete with their higher-resource counterparts for research attention. Without sustained interest and contributions to key evaluation and seed data sets, the delta between high and low-resource languages will continue to expand, exacerbating already prominent technical divides.

Covering 16 languages spanning five continents, the papers in this shared task present a rigorous effort to improve the quality and scope of such data sets. Taken collectively, the authors developed protocols and tools to both refine and introduce new languages to existing FLORES+ and MT Seed data sets. Beyond their technical attributes, the work presented here also aligns with one of OLDI's core commitments: to be community-centric. Every paper in this shared task involves engaging with speakers of the languages of interest, with many authors being native speakers themselves. The linguistics expertise and cultural nuances these researchers brought, alongside the personal convictions many may have, culminated in a body of work that is both scientifically and socially meaningful. It is our hope that the papers showcased in this shared task are the first of a long series designed to consolidate the building blocks needed to advance language technologies for under-served linguistics communities across the world.

## 7 Conclusion

We presented the results of the WMT 2024 OLDI shared task. We accepted ten submissions covering 16 languages, which extend the range of languages included in the foundational datasets FLORES+ and MT Seed. We thank all participants for their contributions and hope that this shared task encourages further efforts towards improved language technologies for more language varieties.

## References

Idris Abdulmumin, Sthembiso Mkhwanazi, Mahlatse S. Mbooi, Shamsuddeen Hassan Muhammad, Ibrahim Said Ahmad, Neo N. Putini, Miehleketo Mathebula, Matimba Shingange, Tajuddeen Gwadabe, and Vukosi Marivate. 2024. Correcting FLORES evaluation dataset for four African languages. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Ife Adebara, AbdelRahim Elmadany, Muhammad Abdul-Mageed, and Alcides Inciarte. 2022. AfroLID: A neural language identification tool for African languages. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1958–1981, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Noëmi Aepli, Chantal Amrhein, Florian Schottmann, and Rico Sennrich. 2023. A benchmark for evaluating machine translation metrics on dialects without standard orthography. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1045–1065, Singapore. Association for Computational Linguistics.

Milind Agarwal, Md Mahfuz Ibn Alam, and Antonios Anastasopoulos. 2023. LIMIT: Language identification, misidentification, and translation using hierarchical models in 350+ languages. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14496–14519, Singapore. Association for Computational Linguistics.

Firoz Ahmed, Nitin Venkateswaran, and Sarah Moeller. 2024. The Bangla/Bengali seed dataset submission to the WMT24 open language data initiative shared task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Felermino Dario Mario Ali, Henrique Lopes Cardoso, and Rui Sousa-Silva. 2024. Expanding FLORES+ benchmark for more low-resource settings: Portuguese-Emakhuwa machine translation evaluation. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4218–4222, Marseille, France. European Language Resources Association.

Laurie Burchell, Alexandra Birch, Nikolay Bogoychev, and Kenneth Heafield. 2023. An open dataset and model for language identification. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 865–879, Toronto, Canada. Association for Computational Linguistics.

Jose Cols. 2024. Spanish corpus and provenance with computer-aided translation for the WMT24 OLDI

shared task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2022. Fleurs: Few-shot learning evaluation of universal representations of speech. *Preprint*, arXiv:2205.12446.

Moussa Doumbouya, Baba Mamadi Diané, Solo Farabado Cissé, Djibrila Diané, Abdoulaye Sow, Séré Moussa Doumbouya, Daouda Bangoura, Fodé Moriba Bayo, Ibrahima Sory 2. Condé, Kalo Mory Diané, Chris Piech, and Christopher Manning. 2023. Machine translation for nko: Tools, corpora, and baseline results. In *Proceedings of the Eighth Conference on Machine Translation*, pages 312–343, Singapore. Association for Computational Linguistics.

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online. Association for Computational Linguistics.

Edoardo Ferrante. 2024. A high-quality seed dataset for Italian machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Isai Gordeev, Sergey Kuldin, and David Dale. 2024. Flores+ translation and machine translation evaluation for the Erzya language. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc'Aurelio Ranzato. 2019. The FLORES evaluation datasets for low-resource machine translation: Nepali–English and Sinhala–English. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages

6098–6111, Hong Kong, China. Association for Computational Linguistics.

Christopher R. Haberland, Jean Maillard, and Stefano Lusito. 2024. Italian-Ligurian machine translation in its cultural context. In *Proceedings of the 3rd Annual Meeting of the Special Interest Group on Under-resourced Languages @ LREC-COLING 2024*, pages 168–176, Torino, Italia. ELRA and ICCL.

Barry Haddow, Rachel Bawden, Antonio Valerio Miceli Barone, Jindřich Helcl, and Alexandra Birch. 2022. Survey of low-resource machine translation. *Computational Linguistics*, 48(3):673–732.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Ali Kuzhuget, Airana Mongush, and Nachyn-Enkhedorzhu Oorzhak. 2024. Enhancing Tuvan language resources through the FLORES dataset. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzman. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

Mukhammadsaid Mamasaidov and Abror Shopulatov. 2024. Open Language Data Initiative: Advancing low-resource machine translation for Karakalpak. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers,

115

Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2024. Scaling neural machine translation to 200 languages. *Nature*, 630(8018):841–846.

Juan Antonio Perez-Ortiz, Felipe Sánchez-Martínez, Víctor M. Sánchez-Cartagena, Miquel Esplà-Gomis, Aaron Galiano Jimenez, Antoni Oliver, Claudi Aventín-Boya, Alejandro Pardos, Cristina Valdés, Jusép Loís Sans Socasau, and Juan Pablo Martínez. 2024. Expanding the flores+ multilingual benchmark with translations for Aragonese, Aranese, Asturian, and Valencian. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. 2022. Data cards: Purposeful and transparent dataset documentation for responsible AI. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1776–1826.

Phillip Rust, Bowen Shi, Skyler Wang, Necati Cihan Camgoz, and Jean Maillard. 2024. Towards privacy-aware sign language translation at scale. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8624–8641, Bangkok, Thailand. Association for Computational Linguistics.

Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, Christopher Klaiber, Pengwei Li, Daniel Licht, Jean Maillard, Alice Rakotoarison, Kaushik Ram Sadagopan, Guillaume Wenzek, Ethan Ye, Bapi Akula, Peng-Jen Chen, Naji El Hachem, Brian Ellis, Gabriel Mejia Gonzalez, Justin Haaheim, Prangthip Hansanti, Russ Howes, Bernie Huang, Min-Jae Hwang, Hirofumi Inaguma, Somya Jain, Elahe Kalbassi, Amanda Kallet, Ilia Kulikov, Janice Lam, Daniel Li, Xutai Ma, Ruslan Mavlyutov, Benjamin Peloquin, Mohamed Ramadan, Abinesh Ramakrishnan, Anna Sun, Kevin Tran, Tuan Tran, Igor Tufanov, Vish Vogeti, Carleigh Wood, Yilin Yang, Bokai Yu, Pierre Andrews, Can Balioglu, Marta R. Costa-jussà, Onur Celebi, Maha Elbayad, Cynthia Gao, Francisco Guzmán, Justine Kao, Ann Lee, Alexandre Mourachko, Juan Pino, Sravya Popuri, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Paden Tomasello, Changhan Wang, Jeff Wang, and Skyler Wang. 2023a. Seamlessm4t-massively multilingual & multimodal machine translation.

Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, John Hoffman, Min-Jae Hwang, Hirofumi Inaguma, Christopher Klaiber, Ilia Kulikov, Pengwei Li, Daniel Licht, Jean Maillard, Ruslan Mavlyutov, Alice Rakotoarison, Kaushik Ram Sadagopan, Abinesh Ramakrishnan, Tuan Tran, Guillaume Wenzek, Yilin Yang, Ethan Ye, Ivan Evtimov, Pierre Fernandez, Cynthia Gao, Prangthip Hansanti, Elahe Kalbassi, Amanda

Kallet, Artyom Kozhevnikov, Gabriel Mejia Gonzalez, Robin San Roman, Christophe Touret, Corinne Wong, Carleigh Wood, Bokai Yu, Pierre Andrews, Can Balioglu, Peng-Jen Chen, Marta R. Costa-jussà, Maha Elbayad, Hongyu Gong, Francisco Guzmán, Kevin Heffernan, Somya Jain, Justine Kao, Ann Lee, Xutai Ma, Alex Mourachko, Benjamin Peloquin, Juan Pino, Sravya Popuri, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Anna Sun, Paden Tomasello, Changhan Wang, Jeff Wang, Skyler Wang, and Mary Williamson. 2023b. Seamless: Multilingual expressive and streaming speech translation. *Preprint*, arXiv:2312.05187.

Shivalika Singh, Freddie Vargus, Daniel D'souza, Börje Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura O'Mahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergun, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Chien, Sebastian Ruder, Surya Guthikonda, Emad Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. 2024. Aya dataset: An open-access collection for multilingual instruction tuning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11521–11567, Bangkok, Thailand. Association for Computational Linguistics.

Felipe Sánchez-Martínez, Juan Antonio Perez-Ortiz, Aaron Galiano Jimenez, and Antoni Oliver. 2024. Findings of the WMT 2024 shared task translation into low-resource languages of Spain: Blending rule-based and neural systems. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Jörg Tiedemann. 2009. News from OPUS - a collection of multilingual parallel corpora with tools and interfaces. In *Recent advances in natural language processing*, volume 5, pages 237–248.

Hongjian Yu, Yiming Shi, Zherui Zhou, and Christopher Haberland. 2024. Machine translation evaluation benchmark for Wu. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

## A  Translation Guidelines

These translation guidelines must be acknowledged by all translators who will be contributing data.

### Important note

Your translations will be used to help train or evaluate machine translation engines. For this reason, this project requires **human translation**.

- If you are translating data for evaluation purposes, such as for FLORES+, using or even referencing machine translation output is not allowed (this includes post-editing).

- **Note** that some machine translation services – including DeepL, Google Translate, and Chat-GPT – prohibit the use of their output for training other translation or AI models, so their use is not permitted.

### General Guidelines

1. You will be translating sentences coming from different sources. Please refer to the source document if available.

2. Do not convert any units of measurement. Translate them exactly as noted in the source content.

3. When translating, please maintain the same tone used in the source document. For example, encyclopedic content coming from sources like Wikipedia should be translated using a formal tone.

4. Provide fluent translations without deviating too much from the source structure. Only allow necessary changes.

5. Do not expand or replace information compared to what is present in the source documents. Do not add any explanatory or parenthetical information, definitions, etc.

6. Do not ignore any meaningful text present in the source.

7. In case of multiple possible translations, please pick the one that makes the most sense (e.g., for gender concordance, cultural fit in the target language, level of formality, etc.).

8. Translations must be faithful to the source in terms of pragmatics such as (if applicable) level of hedging/modality, sentiment and its intensity, negation, speech effects (disfluencies), etc.

9. For proper nouns and common abbreviations, please see the guidelines on Named Entities below.

10. Idiomatic expressions should not be translated word for word. Use an equivalent idiom if one exists. If no equivalent idiom exists, use an idiom of similar meaning. If no similar expressions exist in the target language, paraphrase the idiom such that the meaning is retained in the target language.

11. When a pronoun to be translated is ambiguous (for instance, when it could be interpreted as either him/her or he/she), opt for gender-neutral pronouns (such as them/they) if those exist in the target language. However, when a pronoun to be translated is clearly marked for gender, you should follow the source material and continue to mark for gender.

12. Foreign words and phrases used in the text should be kept in their original language when necessary to preserve the meaning of the sentence (e.g., if given as an example of a foreign word).

### Named entities

Named entities are people, places, organizations, etc., commonly referred to using a proper noun. This section provides guidance on how to handle named entities. Please review the following guidelines carefully:

1. If there is a commonly used term in the target language for the Named Entity:

   (a) If the most commonly used term is the same as in the source language, keep it as it is.

   (b) If the most commonly used term is a translation or a transliteration, use that.

2. If there is no commonly used term:

   (a) If possible, a transliteration of the original term should be used.

   (b) If a transliteration would not be commonly understood in the context, and the source term would be more acceptable, you may retain the original term.

# Results of the WAT/WMT 2024 Shared Task on Patent Translation

**Shohei Higashiyama**

National Institute of Information and Communications Technology, Japan
shohei.higashiyama@nict.go.jp

## Abstract

This paper presents the results of the patent translation shared task at the 11th Workshop on Asian Translation and 9th Conference on Machine Translation. Two teams participated in this task, and their submitted translation results for one or more of the six language directions were automatically and manually evaluated. The evaluation results demonstrate the strong performance of large language model-based systems from both participants.

## 1 Introduction

The patent translation task using the JPO Patent Corpus has been held under the Workshop on Asian Translation (WAT) in 2015–2023 (Nakazawa et al., 2023) and under the Conference on Machine Translation (WMT) this year.[1] Due to the high demand for patent translation, this task has attracted many participants particularly in the early WAT workshops: a total of 30 teams over the past 10 years as in Table 1.

This year, two teams participated in this task; one participant submitted translation results for two language directions, and the other for six out of six language directions, that is, Chinese↔Japanese, Korean↔Japanese, and English↔Japanese. Both teams employed large language model (LLM)-based systems, and the submitted translation results were evaluated by both automatic and human evaluation metrics. In this paper, we describe the evaluation dataset and procedure, and report the evaluation results for the submitted outputs.

## 2 Dataset

The JPO Patent Corpus (JPC)[2] was constructed by the Japan Patent Office (JPO) in collaboration with

| Year | # of teams |
|------|-----------|
| 2015 | 8 |
| 2016 | 6 |
| 2017 | 4 |
| 2018 | 2 |
| 2019 | 3 |
| 2020 | 2 |
| 2021 | 3 |
| 2022 | 0 |
| 2023 | 0 |
| 2024 | 2 |
| Total | 30 |

Table 1: The number of participant teams for the patent task over the years.

National Institute of Information and Communications Technology (NICT). The corpus consists of Chinese-Japanese (zh-ja), Korean-Japanese (ko-ja), and English-Japanese (en-ja) parallel sentences of patent descriptions. Most sentences were extracted from documents with one of four International Patent Classification sections: chemistry, electricity, mechanical engineering, and physics. As shown in Table 2, the dataset for each language pair consists of training, development, development-test, and multiple test sets. These datasets were constructed from patent families using automatic sentence alignment (Utiyama and Isahara, 2007), except for the test-N4 set where target sentences were manual translated from the source sentences.

A characteristic of the corpus is the use of fixed training and test datasets over the years, which allows for the comparison of new systems with past systems. The possible issue of data leakage is minimized: the data is provided only to applicants who have committed to participating in each annual workshop, and participants are required to delete the data after the workshop concludes.

---

[1] Similarly to other WAT shared tasks, this task is organized by WAT organizers but is held under WMT this year due to the collaboration between the workshop and conference.

[2] https://lotus.kuee.kyoto-u.ac.jp/WAT/patent/

| Set | # of Sentences | | | Published Years | Introduced Event |
|---|---|---|---|---|---|
| | zh-ja | ko-ja | en-ja | | |
| Train | 1,000,000 | 1,000,000 | 1,000,000 | 2011–2013 | WAT 2015–2016 |
| Dev | 2,000 | 2,000 | 2,000 | 2011–2013 | WAT 2015–2016 |
| DevTest | 2,000 | 2,000 | 2,000 | 2011–2013 | WAT 2015–2016 |
| Test-N1 | 2,000 | 2,000 | 2,000 | 2011–2013 | WAT 2015–2016 |
| Test-N2 | 3,000 | – | 3,000 | 2016–2017 | WAT 2018 |
| Test-N3 | 204 | 230 | 668 | 2016–2017 | WAT 2018 |
| Test-N4 | 5,000 | 5,000 | 5,000 | 2019–2020 | WAT 2022 |
| Test-2022 | 10,204 | 7,230 | 10,668 | 2011–2020 | WAT 2022 |

Table 2: Statistics of the JPO Corpus. The published years column represents the years for the source sentences. The introduced event column indicates the events for which each dataset was first introduced.

# 3 Evaluation Procedure

## 3.1 Automatic Evaluation

Task participants were required to submit translation results via the WAT Submission site.[3] For the results submitted with the "publish" checkbox selected, automatic evaluation scores were calculated and displayed in the WAT Evaluation site.[4] As the automatic evaluation metrics, we used BLEU (Papineni et al., 2002) with `multi-bleu.perl` in the Moses toolkit (Koehn et al., 2007) version 2.1.1[5] and RIBES (Isozaki et al., 2010) with `RIBES.py` version 1.02.4.[6]

Prior to calculating scores, reference sentences and output translation sentences were tokenized with the tokenization tools for each language: Juman 7.0 (Kurohashi et al., 1994), KyTea 0.4.6 (Neubig et al., 2011) with the full SVM model[7] and MeCab 0.996 (Kudo et al., 2004) with IPA dictionary 2.7.0[8] for Japanese, KyTea 0.4.6 with the full SVM Model (MSR model) and Stanford Word Segmenter (Tseng, 2005) version 2014-06-16 with the CTB and PKU models[9] for Chinese, mecab-ko[10] for Korean, and `tokenizer.perl`[11] in the Moses

| 5 | All important information is transmitted correctly. (100%) |
|---|---|
| 4 | Almost all important information is transmitted correctly. (80%–) |
| 3 | More than half of important information is transmitted correctly. (50%–) |
| 2 | Some of important information is transmitted correctly. (20%–) |
| 1 | Almost all important information is NOT transmitted correctly. (–20%) |

Table 3: Ratings and their descriptions in the JPO adequacy criterion.

toolkit for English. The detailed procedures are shown on the WAT Evaluation site.[12]

## 3.2 Human Evaluation

We conducted human evaluation for selected translation results based on the JPO adequacy evaluation criterion, which is originally defined by JPO to assess the quality of translated patent documents. For this evaluation, we used the test-N3 set for each language direction for the following reasons: (1) parallel sentences have been manually aligned (translations were manually created from the original sentences), and (2) both participants submitted results for this test set.

The evaluation was performed by two annotators (translation experts) for each system as follows. First, 200 sentences for evaluation were randomly selected from the test-N3 set in advance (the same 200 sentences were used for all systems). (2) The 200 pairs of the source sentences and translated sentences by the system were shown to each annotator, and the ratings between 1 and 5 were assigned to each sentence by the annotator as shown in Table 3.

---

[3] https://lotus.kuee.kyoto-u.ac.jp/WAT/submission/index.php

[4] https://lotus.kuee.kyoto-u.ac.jp/WAT/evaluation/index.html

[5] https://github.com/moses-smt/mosesdecoder/tree/RELEASE-2.1.1

[6] http://www.kecl.ntt.co.jp/icl/lirg/ribes/index.html

[7] http://www.phontron.com/kytea/model.html

[8] http://code.google.com/p/mecab/downloads/detail?name=mecab-ipadic-2.7.0-20070801.tar.gz

[9] http://nlp.stanford.edu/software/segmenter.shtml

[10] https://bitbucket.org/eunjeon/mecab-ko/

[11] https://github.com/moses-smt/mosesdecoder/tree/RELEASE-2.1.1/scripts/tokenizer/tokenizer.perl

[12] http://lotus.kuee.kyoto-u.ac.jp/WAT/evaluation/index.html

| Parameter | Value |
|---|---|
| encoder_type | brnn |
| brnn_merge | concat |
| src_seq_length | 150 |
| tgt_seq_length | 150 |
| src_vocab_size | 100,000 |
| tgt_vocab_size | 100,000 |
| src_words_min_frequency | 1 |
| tgt_words_min_frequency | 1 |

Table 4: The configuration used for the baseline model. For other parameters, tge default values were used.

## 4  Baseline System

The organizers built a baseline system, a recurrent neural network (RNN) encoder-decoder model with attention (Bahdanau et al., 2014) using Open-NMT (Klein et al., 2017) with the configuration shown in Table 4 and the same tokenizers for automatic evaluation explained in §3.1. This baseline system uses the old neural machine translation (NMT) model built for WAT 2018 and serves as a weak baseline for comparison. However, as shown in §6, many past participants have adopted Transformer-based systems, allowing for the performance comparison with Transformer models (Vaswani et al., 2017) for recent participants.

## 5  Participant Systems

Two teams participated in the patent translation task: GenAI (Yonsei University) and sakura (Rakuten Institute of Technology). The details on the submitted systems are as follows.

**sakura** used an LLM-based system fine-tuned with simple translation prompt on the JPC training set for the corresponding language pair. As their backbone model, they adopted RakutenAI-7B-chat,[13] which had been pretrained on English and Japanese texts.

**GenAI** used an LLM-based system fine-tuned on only 1,000 sentences from the JPC training set. Their backbone model is Mistral-Nemo-Instruct-2407 (12B),[14] which had been pretrained on multilingual texts. During both fine-tuning and testing, their system identified domain-specific terms in each input source sentence by matching them with

their bilingual terminology dictionary, and then generated the translation based on prompt that required the use of the specified term pairs.

## 6  Evaluation Results

### 6.1  Main Results

For the same reasons mentioned in §3.2, we only present the results for the test-N3 set; results for other test sets can be found at the WAT Evaluation site.[15] Table 5, 6, 7, 8, 9, and 10 show the performance of evaluated system for each language direction (systems with "∗r" indicate they used external resources). The tables present the automatic and human evaluation scores of the two participants' systems (one system per participant, selected based on the BLEU score), as well as the organizer's baseline and the best participant systems from previous years. The model type columns indicate whether the system employed statistical machine translation (SMT), RNN-based NMT, or Transformer (TF)-based NMT, and whether it corresponds to a decoder-only model (Dec) or an encoder-decoder model (EncDec). The BLEU/RIBES scores for the translation tasks into Japanese and Chinese represent the average BLEU/RIBES scores based on three different tokenizers.[16] The JPO adequacy scores (Adeq) represent the average of the scores assigned by two annotators.

We observed the following findings. (1) Unsurprisingly, both participants' systems as well as all previous best systems outperformed the baseline for all language directions in terms of automatic metrics. (2) The LLM-based systems by the two participants achieved strong results in terms of automatic metrics; GenAI's system outperformed the previous systems for ko→ja and ja→ko and sakura's system outperformed the previous systems for ja→ko and ja→en. However, the previous systems maintained the highest scores for zh→ja, ja→zh, and en→ja. (3) Both participants' systems achieved high adequacy scores of over 4. However, importantly, a system with a higher automatic evaluation score did not necessarily achieved a higher human evaluation score. Specifically, sakura's system yielded lower automatic evaluation scores than GenAI's system (e.g., BLEU of 52.77 vs. 67.10 for ja→ko and 68.00 vs. 70.60 for ja→ko), but

| System | Model Type | BLEU | RIBES | Adeq |
|---|---|---|---|---|
| GenAI best | TF Dec | 67.10 | 0.9225 | 4.66 |
| 2018 best | SMT | 54.63 | 0.9056 | – |
| 2019 best | TF EncDec | 54.42 | 0.9012 | – |
| 2020 best | TF EncDec | 53.77 | 0.9044 | – |
| 2021 best*ʳ | TF EncDec | 53.48 | 0.9014 | – |
| sakura best | TF Dec | 52.77 | 0.8982 | 4.67 |
| Baseline | RNN EncDec | 52.65 | 0.8975 | – |

Table 5: Results on the ko→ja test-N3 set.

| System | Model Type | BLEU | RIBES | Adeq |
|---|---|---|---|---|
| GenAI best | TF Dec | 70.60 | 0.9391 | 4.39 |
| sakura best | TF Dec | 68.00 | 0.9268 | 4.76 |
| 2021 best | TF EncDec | 66.25 | 0.9252 | – |
| 2019 best | TF EncDec | 65.74 | 0.9228 | – |
| 2020 best | TF EncDec | 64.30 | 0.9223 | – |
| Baseline | RNN EncDec | 62.43 | 0.9153 | – |

Table 6: Results on the ja→ko test-N3 set.

| System | Model Type | BLEU | RIBES | Adeq |
|---|---|---|---|---|
| 2020 best | TF EncDec | 40.51 | 0.7568 | – |
| 2019 best | TF EncDec | 24.96 | 0.7639 | – |
| 2018 best | TF EncDec | 24.87 | 0.7492 | – |
| 2021 best*ʳ | TF EncDec | 22.67 | 0.7716 | – |
| sakura best | TF Dec | 20.83 | 0.7615 | 4.24 |
| Baseline | RNN EncDec | 17.28 | 0.7322 | – |

Table 7: Results on the zh→ja test-N3 set.

| System | Model Type | BLEU | RIBES | Adeq |
|---|---|---|---|---|
| 2020 best | TF EncDec | 44.34 | 0.8340 | – |
| 2021 best*ʳ | TF EncDec | 31.09 | 0.8550 | – |
| 2019 best | TF EncDec | 29.82 | 0.8390 | – |
| sakura best | TF EncDec | 26.60 | 0.8245 | 4.33 |
| 2018 best | TF EncDec | 24.66 | 0.8261 | – |
| Baseline | RNN EncDec | 23.68 | 0.7886 | – |

Table 8: Results on the ja→zh test-N3 set.

| System | Model Type | BLEU | RIBES | Adeq |
|---|---|---|---|---|
| 2019 best*ʳ | TF Enc-Dec | 55.32 | 0.8827 | – |
| sakura best | TF Dec | 53.93 | 0.8803 | 4.44 |
| 2021 best*ʳ | TF Enc-Dec | 53.34 | 0.8753 | – |
| 2018 best*ʳ | SMT | 52.07 | 0.8643 | – |
| 2020 best | TF Enc-Dec | 50.95 | 0.8665 | – |
| Baseline | RNN Enc-Dec | 46.39 | 0.8438 | – |

Table 9: Results on the en→ja test-N3 set.

| System | Model Type | BLEU | RIBES | Adeq |
|---|---|---|---|---|
| sakura best | TF Dec | 43.20 | 0.8505 | 4.08 |
| 2019 best*ʳ | TF Enc-Dec | 41.37 | 0.8499 | – |
| 2021 best*ʳ | TF Enc-Dec | 40.73 | 0.8546 | – |
| 2020 best | TF Enc-Dec | 39.94 | 0.8413 | – |
| Baseline | RNN Enc-Dec | 35.01 | 0.8230 | – |

Table 10: Results on the ja→en test-N3 set.

achieved similar or better adequacy scores (4.67 vs. 4.66 for ja→ko and 4.76 vs. 4.39 for ja→ko). This result highlights the need for using a variety of evaluation metrics, such as neural-based metrics, which have been demonstrated to correlate well with human judgement (Freitag et al., 2023).

## 6.2 Detailed Human Evaluation Results

Table 11 shows the detailed results of the JPO adequacy evaluation for a total of eight participant systems, which were selected from among the same participant's systems based on the BLEU score. The "Adequacy Score" column represents the average of ratings assigned to 200 sentences by each annotator for the Annotator="A"/"B" rows and the average and standard deviation of the average score by the two annotators (A and B) for the Annotator="Both" row, which is shown as the adequacy score (Adeq) in Table 5–10.

We observed the following findings. First, most sentences were assigned scores over 4 (75% or more sentences for each translation result, except for sakura's ja-en result evaluated by Annotator B). This indicates that there were many high-quality translation overall, but more accurate systems have room for development, considering that the translations with a score lower than 5 account for more than 20–50% in most cases of annotator-level evaluation results.

Second, the difference of sentence-level scores between two annotators ("Diff Score") was 0 or 1 in most cases, and there were only nine sentences with the difference score of 2 over all translation results. As a result, the adequacy scores between two annotators were close in many cases, but relatively large standard deviation (close to or greater than 0.2) was observed in three cases, i.e., sakura ja-zh, GenAI ja-ko, and sakura ja-en results. In the latter cases, there were somewhat many mismatches; each translation result included over 100 sentences with a score difference of 1 from the two annotators and/or a few sentences with a score difference of 2.

For the nine sentences with a score difference of 2, we conducted a meta-review by a third evaluator, distinct from the two annotators. We found that which annotator provided the more appropriate rating varied depending on the example. In some examples, one annotator overlooked a mistranslation and assigned a higher rating. In other examples, there were no mistranslations, but one annotator still assigned a lower rating. Additionally, in cases

| Lang | Team | Data ID | Annotator | Adequacy Score (Avg. ± SD) | Distribution of Ratings | | | | | Diff Score | | |
|------|------|---------|-----------|------------|---|---|---|---|---|---|---|---|
| | | | | | 1 | 2 | 3 | 4 | 5 | 0 | 1 | 2 |
| zh-ja | sakura | 7302 | A | 4.24 | 4 | 4 | 24 | 76 | 92 | | | |
| | | | B | 4.24 | 2 | 6 | 26 | 74 | 92 | | | |
| | | | Both | 4.24 ± 0 | | | | | | 130 | 70 | 0 |
| ja-zh | sakura | 7257 | A | 4.50 | 2 | 5 | 17 | 43 | 133 | | | |
| | | | B | 4.15 | 7 | 10 | 30 | 52 | 101 | | | |
| | | | Both | 4.33 ± 0.18 | | | | | | 120 | 80 | 2 |
| ko-ja | sakura | 7311 | A | 4.79 | 1 | 1 | 7 | 21 | 170 | | | |
| | | | B | 4.55 | 2 | 0 | 9 | 65 | 124 | | | |
| | | | Both | 4.67 ± 0.12 | | | | | | 137 | 63 | 0 |
| ko-ja | GenAI | 7180 | A | 4.84 | 0 | 0 | 1 | 37 | 162 | | | |
| | | | B | 4.51 | 0 | 0 | 0 | 99 | 101 | | | |
| | | | Both | 4.66 ± 0.15 | | | | | | 124 | 76 | 0 |
| ja-ko | sakura | 7224 | A | 4.64 | 0 | 4 | 4 | 52 | 140 | | | |
| | | | B | 4.87 | 0 | 1 | 4 | 15 | 180 | | | |
| | | | Both | 4.76 ± 0.12 | | | | | | 148 | 52 | 0 |
| ja-ko | GenAI | 7267 | A | 4.16 | 0 | 7 | 38 | 71 | 84 | | | |
| | | | B | 4.61 | 0 | 0 | 9 | 60 | 131 | | | |
| | | | Both | 4.39 ± 0.23 | | | | | | 98 | 102 | 0 |
| en-ja | sakura | 7278 | A | 4.49 | 0 | 4 | 15 | 61 | 120 | | | |
| | | | B | 4.40 | 0 | 3 | 35 | 41 | 121 | | | |
| | | | Both | 4.44 ± 0.04 | | | | | | 123 | 73 | 0 |
| ja-en | sakura | 7309 | A | 3.83 | 2 | 22 | 59 | 43 | 74 | | | |
| | | | B | 4.33 | 1 | 5 | 26 | 64 | 104 | | | |
| | | | Both | 4.08 ± 0.25 | | | | | | 79 | 144 | 7 |

Table 11: Detailed results of the JPO adequacy evaluation for the test-N3 set. The "Distribution of Ratings" column shows the number of sentences with each rating of 1–5. The "Diff Score" represents the number of sentences with each difference score, which means the difference of ratings between two annotators.

where the translation contained garbled characters, one annotator assigned a lower rating.

# 7 Conclusion

This paper summarizes the results of the WAT/WMT 2024 shared task on patent translation. The patent translation task using the JPO Patent Corpus has been held for ten years, and this will be the last time.[17] We believe that extensive development efforts by task participants over the past 10 years have contributed to advance machine translation technologies for the patent domain.

# References

Dzmitry Bahdanau, KyungHyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv:1409.0473*.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, EMNLP '10, pages 944–952, Stroudsburg, PA, USA. Association for Computational Linguistics.

Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. OpenNMT: Open-source toolkit for neural machine translation. In *Proceedings of ACL 2017, System Demonstrations*, pages 67–72. Association for Computational Linguistics.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Annual Meeting of the Association for Computational Linguistics (ACL), demonstration session*.

---

[17]The JPO Patent Copurs will be provided to applicants via the ALAGIN forum (https://www.alagin.jp/index-e.html) for future research.

Taku Kudo, Kaoru Yamamoto, and Yuji Matsumoto. 2004. Applying conditional random fields to Japanese morphological analysis. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 230–237, Barcelona, Spain. Association for Computational Linguistics.

Sadao Kurohashi, Toshihisa Nakamura, Yuji Matsumoto, and Makoto Nagao. 1994. Improvements of Japanese morphological analyzer JUMAN. In *Proceedings of The International Workshop on Sharable Natural Language*, pages 22–28.

Toshiaki Nakazawa, Kazutaka Kinugawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida, Makoto Morishita, Ondřej Bojar, Akiko Eriguchi, Yusuke Oda, Chenhui Chu, and Sadao Kurohashi. 2023. Overview of the 10th workshop on Asian translation. In *Proceedings of the 10th Workshop on Asian Translation*, pages 1–28, Macau SAR, China. Asia-Pacific Association for Machine Translation.

Graham Neubig, Yosuke Nakata, and Shinsuke Mori. 2011. Pointwise prediction for robust, adaptable Japanese morphological analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: Short Papers - Volume 2*, HLT '11, pages 529–533, Stroudsburg, PA, USA. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *ACL*, pages 311–318.

Huihsin Tseng. 2005. A conditional random field word segmenter. In *In Fourth SIGHAN Workshop on Chinese Language Processing*.

Masao Utiyama and Hitoshi Isahara. 2007. A Japanese-English patent parallel corpus. In *Proceedings of Machine Translation Summit XI: Papers*, Copenhagen, Denmark.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

# Findings of the WMT 2024 Biomedical Translation Shared Task: Test Sets on Abstract Level

**Mariana Neves**[1] *   **Cristian Grozea**[2]   **Philippe Thomas**[3]   **Roland Roller**[3]
**Rachel Bawden**[4]   **Aurélie Névéol**[5]   **Steffen Castle**[3]
**Vanessa Bonato**[6]   **Giorgio Maria Di Nunzio**[6]   **Federica Vezzani**[6]
**Maika Vicente Navarro**[7]   **Lana Yeganova**[8]   **Antonio Jimeno Yepes**[9,10]

[1]German Centre for the Protection of Laboratory Animals (Bf3R),
German Federal Institute for Risk Assessment (BfR), Berlin, Germany
[2]Fraunhofer Institute FOKUS, Berlin, Germany
[3]German Research Center for Artificial Intelligence (DFKI), Berlin, Germany
[4]Inria, Paris, France
[5]Université Paris-Saclay, CNRS, LISN, Orsay, France
[6]Dept. of Linguistic and Literary Studies University of Padua, Italy
[7]Leica Biosystems, Australia
[8]NCBI/NLM/NIH, Bethesda, USA
[9]RMIT University, Australia
[10]Unstructured Technologies, USA

## Abstract

We present the results of the ninth edition of the Biomedical Translation Task at WMT'24. We released test sets for six language pairs, namely, French, German, Italian, Portuguese, Russian, and Spanish, from and into English. Each test set consists of 50 abstracts from PubMed. Differently from previous years, we did not split abstracts into sentences. We received submissions from five teams, and for almost all language directions. We used a baseline/comparison system based on Llama 3.1 and share the source code at `https://github.com/cgrozea/wmt24biomed-ref`.

## 1 Introduction

In this paper, we present a description and the findings of the ninth edition of the Biomedical Translation Task,[1] which took place at the ninth edition of the Conference for Machine Translation (WMT'24). The shared task aims to support advances in Machine Translation (MT) in the biomedical domain, especially for scientific literature. Previous editions of the shared task addressed up to seven language pairs and included the release of training and test sets (Bojar et al., 2016; Jimeno Yepes et al., 2017; Neves et al., 2018; Bawden et al., 2019, 2020; Yeganova et al., 2021; Neves et al., 2022, 2023). All previous data is available in the shared task repository.[2]

Similar to previous years, our test sets consist of biomedical abstracts, which have been included to PubMed[3] just before publishing the test set, to decrease the likelihood of data contamination. We prepared test sets for six languages from and into English, namely, French (fr2en, en2fr), German (de2en, en2de), Italian (it2en, en2it), Portuguese (pt2en, en2pt), Russian (ru2en, en2ru), and Spanish (es2en, en2es). The test sets consist of 50 abstract pairs for each of the 12 language directions above. Some of the test sets were also released as test suites in the General Task of WMT (Kocmi et al., 2024). After the release of the test sets, the participants had around two weeks to submit their automatic translations. For this year's shared task, the following features were introduced:

- The selection of the articles for the test sets was based on topics of interest to the task organizers (Section 2);

---

* The contributions of the authors are the following: MN prepared the MEDLINE test sets, performed manual validation, and organized the shared task; CG developed the baseline system; PT, RR, RB, AN, SC, VB, GMN, FV, MVN, LY performed manual validation; AJY performed manual validation and the automatic evaluation, as well as co-organized the shared task; All authors approved the final version of the manuscript. E-mail for contact: mariana.lara-neves@bfr.bund.de

[1] `http://www2.statmt.org/wmt24/biomedical-translation-task.html`

[2] `https://github.com/biomedical-translation-corpora/corpora`

[3] `https://pubmed.ncbi.nlm.nih.gov/`

- The test sets consist of paragraphs comprising the papers' title and the abstract, i.e. no sentence splitting and alignment were carried out (Section 2);

- Consequently, we only performed a manual evaluation on the abstract level (cf. Section 6);

- We used as a baseline/comparison a local large language model, Llama 3.1 (cf. Section 3);

- We performed the automatic evaluation also based on COMET (Rei et al., 2020), besides BLEU (Papineni et al., 2002).

## 2 Test sets

We downloaded the daily update files from PubMed[4] around mid-April for the preparation of the test sets. As usual, we first identify all articles that are available in English as well as one of the non-English languages that we address in the shared task. Subsequently, we selected 100 pairs of articles for each language pair, which were later split into two sets, i.e., from and into English.

This year, we aimed to prioritize three topics[5] in our test sets. While selecting the articles, we restricted each topic to around a third of the total. Still, this limit was frequently not reached because too few articles included any of the three selected topics. The three topics are listed below:

- Animals: D000818

- SARS-CoV-2: D000086402

- Pancreatic Neoplasms: D010190

Subsequently, the 100 selected articles for each language pair were split between the two test set directions. Test set statistics are shown in Table 1. No further processing was performed on the test sets, and these were released as a plain text file, one for each language pair, each with 50 lines, and one for each article. Each line is composed of the title and abstract of the article.

## 3 Baseline/Comparison system

While we used GPT 3.5 as a comparative model last year, we decided to use a self-hosted open-weight large language model this year. Several such models are available of various sizes, licenses, and performance levels in the MT task. Based on the previously accumulated hands-on experience in informally evaluating several open-weight models in multiple tasks, including translation, we selected one of the best performing models, namely Llama 3.1 (Dubey et al., 2024).[6]

The Llama models are open in the sense that their weights and supporting code are freely available, but the usage is limited by a relatively liberal license. In the case of the model used here, the precise licenses are "LLAMA 3.1 COMMUNITY LICENSE AGREEMENT" and "Llama 3.1 Acceptable Use Policy". The last one prohibits using the model to violate the law or the rights of others, to activities related to bodily harm, including military, to generate false information, and includes a clause making it compulsory to report "risky content generated by the model". This risky content can arise when used for medical texts in the form of mistranslated medical procedures.

To interact with the model, we used ollama,[7] through which the model can be queried (i.e. we can programmatically perform tasks with the selected LLM and retrieve the response to those tasks, e.g. from a program written in the Python programming language). In addition, ollama provides a command line interface that can be used to pull further models or to interact with a model in a text-based chat interface.

**Implementation decisions** We used "Meta Llama 3.1 70B Instruct"[8] (known in ollama as llama3.1:70b), which means the approximate number of parameters is $70 * 10^9$. Such an LLM is run fully accelerated by a GPU only when the parameters fit into the video RAM of the GPU. Since we used a Nvidia A6000 ADA, a 48 GB RAM GPU card, we used the quantization Q4_0 (4 bits per parameter). This makes the actual size of such a model 37.22 GiB and fits in the 48 GB VRAM of the GPU. With the other temporary data needed in the same memory during processing, the occupation of the VRAM went up to 41.2 GB (85%). To evaluate the impact of using the same model with a smaller card, we also tested a 24 GB VRAM card, Nvidia A5000. This raised the CPU usage to 28 cores (from 2) and processing was slower.

---

| topics | fr2en | de2en | it2en | pt2en | ru2en | es2en | en2fr | en2de | en2it | en2pt | en2ru | en2es |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SARS-CoV-2 | 15 | 18 | - | 5 | 11 | 17 | 9 | 15 | - | - | 13 | 16 |
| Pancr. Neopl. | 2 | 15 | 1 | - | 3 | 2 | - | 15 | 2 | 1 | 4 | 1 |
| Animals | 15 | 17 | 5 | 17 | 14 | 22 | 20 | 17 | 5 | 18 | 20 | 13 |
| other | 19 | - | 44 | 28 | 22 | 9 | 21 | 4 | 43 | 31 | 14 | 20 |

Table 1: Statistics of the topics in the test sets. The topic "other" refers to articles that do not contain any of the three selected topics. The sum of the values for one language pair might be higher than 50 because some articles contain more than one topic.

**Prompt used**  Choosing the right prompt is important for instruction-tuned LLMs and is still rather an art than a science. We started with the prompt "*You are a helpful assistant specialised in biomedical translation. You will be provided with a text in {src}, and your task is to translate it into {dest}.*" where *src* is the name of the source language and *dest* is the name of the target language.

Visual examination of one text entry (out of the 50 in the test set) per language pair showed the following undesirable behaviour in the MT output generated by the LLM, which we tried to fix by changes to the prompt:

- in one case some additional text, with the meaning "this is the translation into German", which was fixed by adding "*You will add nothing and comment nothing, just produce the accurate translation of the text in specialist language.*" to the prompt;

- additional formatting of the output text through the insertion of newlines, which was almost entirely fixed by adding "*Keep the formatting as close as possible to the source and especially do not insert any newline.*" to the prompt.

- the occasional replacement of digits by their names. We decided not to try to fix this.

After a complete run, we noticed that the LLM still failed to respect the original format of the source texts (it still sometimes produced multiple lines per source text). Visual inspection showed that in a few cases it still attempted to format the subsections of the translated test despite being asked to refrain from doing that. Therefore, explicit postprocessing was carried out to eliminate the line breaks from the LLM's outputs.

Some good features of the translated texts were also noticed, such as localized acronyms e.g. translating English *Real-time functional magnetic resonance imaging (fMRI)* to French *L'imagerie fonc-tionnelle par résonance magnétique (IRMf)*. Quite impressive was how well the translation retained the quantitative results in the fairly long source texts, while simultaneously applying number localization transformations, such as swapping the decimal point with the decimal comma.

**Run-time Statistics**  Measured duration in seconds with an A6000 in each case for 50 texts:

| | | | | | |
|---|---|---|---|---|---|
| en2de | 1232 | en2es | 1065 | en2fr | 1202 |
| de2en | 728 | es2en | 902 | fr2en | 859 |
| en2it | 1413 | en2pt | 1098 | en2ru | 1110 |
| it2en | 810 | pt2en | 748 | ru2en | 641 |

With an A5000, the speed was about 10 times slower. A GPU-free execution is also possible, but it can be too slow to be practical.

**Energy consumption, $CO_2$ emissions**  For the A6000 card, a total of $11,607$ seconds at about 1 kW (300W the GPU itself) equals an amount of 3.22 kWh and an equivalent $CO_2$ emission of 1.16 kg – at the average 360 g $CO_2$/kWh in Germany, equivalent to the emission of an ICE (internal combustion engine) car driven for about 9.5km. For the slower card, which totalled $131,898$ execution seconds, the figures are 36.64 kWh and therefore 13.2 kg $CO_2$.

## 4  Teams and systems

We followed similar dates to the WMT General Translation Shared Task, releasing the test sets on June 27th, 2024 and allowing submission until July 12th, 2024 (after an extension). We released all test sets both in our submission system (Google Form) and the OCELoT tool.[9] We also included our test sets for en2de, en2es, and en2ru as test suites in the General Task[10] in OCELoT. These were the only language pairs that overlapped with the ones from the General Task.

---

[9] https://ocelot-west-europe.azurewebsites.net/
[10] http://www2.statmt.org/wmt24/translation-task.html

| Team ID | Institution | Publication |
|---------|-------------|-------------|
| ADAPT | Dublin City University, Ireland | (Castaldo et al., 2024) |
| AIST | National Institute of Advanced Industrial Science and Technology, Japan | |
| DCU | Dublin City University, Ireland | |
| HW-TSC | Huawei Translation Service Center, China | |
| Unbabel | Unbabel, Portugal | |

Table 2: List of the participating teams and systems.

We received submissions from five teams that directly registered to our task. We list them in Table 2 and present details about their systems below.

**ADAPT** (Castaldo et al., 2024). For the submissions identified as "run1" for de2en, en2de, fr2en, and en2fr, the participants relied on NLLB-200's distilled 600M variant (NLLB Team et al., 2022), which was fine-tuned on around 10k parallel segments from in-domain training data in the respective language pair. Run2 for en2de, in addition to the above approach, included post-edition by LLM agents powered by GPT-4o.[11] Finally, for run3 for de2en, they relied on LLama-3-8B[12] fine-tuned on around 10k parallel sentences and few-shot prompting using fuzzy matches retrieved by similarity search from the training dataset.

**AIST.** For run1 of de2en, the team relied on a Mega model (Ma et al., 2023) trained from scratch and fine-tuned on parallel biomedical data from MEDLINE. For run2 for both en2de and de2en, they used a Mega model, an ensemble of four checkpoints trained from scratch and fine-tuned on the same data. For all submissions, they estimate the following sizes of training data used: 3M from in-domain, 5M from open domain, and 3M monolingual.

**DCU.** We do not have much information about the system behind the submissions for this team, except for a short description citing the Mistra-7B language model[13] for ru2en and fr2en.

**HW-TSC.** For all submissions to en2de and de2en, the team relied on a system based on Transformers that was trained from scratch on in-domain and open-domain parallel and monolingual data (Wu et al., 2023). It is not clear which changes were carried out for the distinct runs.

**Unbabel** The submissions for all language pairs consisted of a new version of the Tower LLM (Alves et al., 2024), either with Greedy (run1) or MBR (run2) decoding. The LLM has 70B parameters, was built on top of Llama3, and its continued pre-training phase used 25B tokens for 15 languages, followed by fine-tuning with instructions for all the languages in a variety of tasks, including MT.

## 5 Automatic evaluation

We ran automatic evaluation based on BLEU (Papineni et al., 2002) and COMET (Rei et al., 2020). We present the results for the submissions to the biomedical translation task using our form in Tables 3 (from English) and 4 (into English), as well as the ones from OCELoT for our task in Table 5 and for our test suites submitted to the General Task in Table 6. All scores were multiplied by 100.

### 5.1 Biomedical Task submission system

Among all submissions, including our baseline system, the highest BLEU score was 55.63 for pt2en (Unbabel run1) and the highest COMET score was of 89.71 for en2ru (Unbabel run2). The submissions that scored better were the ones from Unbabel and our baseline system, e.g., for en2de, en2fr, en2it, en2pt, and en2ru, with some few exceptions where another system also obtained a high score, e.g., AIST for en2de and DCU for en2ru. The submissions from Unbabel usually scored slightly higher than our baseline, with a few exceptions, e.g., en2pt, fr2en, and es2en.

We observed that the two types of metric score were rather equivalent and that submissions that scored high for BLEU also did so for COMET. However, some submissions had very different BLEU scores for similar COMET scores. For instance, the baseline system obtained the BLEU scores of 31.67 and 51.65 for en2de and en2pt, respectively, but around 87.00 for the COMET score in both cases. Overall, the scores from this year's

---

[11] https://platform.openai.com/docs/models/gpt-4o

[12] https://huggingface.co/meta-llama/Meta-Llama-3-8B

[13] mistralai/Mistral-7B-v0.1

| Team | Run | Metric | en2de | en2fr | en2it | en2pt | en2es | en2ru |
|---|---|---|---|---|---|---|---|---|
| ADAPT | 1 | BLEU | 25.03 | 29.92 | | | | |
| | | COMET | 84.31 | 78.14 | | | | |
| ADAPT | 2 | BLEU | *30.16 | | | | | |
| | | COMET | 85.30 | | | | | |
| AIST | 2 | BLEU | 33.80 | | | | | |
| | | COMET | 85.59 | | | | | |
| DCU | - | BLEU | 16.46 | | 29.12 | 38.97 | | 31.28 |
| | | COMET | 64.78 | | 80.39 | 74.17 | | 87.00 |
| HW-TSC | 1 | BLEU | *28.77 | | | | | |
| | | COMET | 82.92 | | | | | |
| HW-TSC | 2 | BLEU | 28.46 | | | | | |
| | | COMET | 82.83 | | | | | |
| HW-TSC | 3 | BLEU | 28.32 | | | | | |
| | | COMET | 83.14 | | | | | |
| Unbabel | 1 | BLEU | 34.22 | 53.54 | 34.84 | 50.35 | | 35.76 |
| | | COMET | 87.48 | 87.26 | 85.17 | 87.03 | | 88.97 |
| Unbabel | 2 | BLEU | *32.13 | *49.76 | *32.06 | *48.47 | | *32.35 |
| | | COMET | 88.09 | 87.60 | 86.04 | 87.55 | | 89.71 |
| Baseline | - | BLEU | 31.67 | 45.98 | 31.64 | 51.65 | 47.95 | 30.92 |
| | | COMET | 87.00 | 87.03 | 85.00 | 87.02 | 85.37 | 87.55 |

Table 3: BLEU and COMET scores for submissions to the Biomedical Task submission system, for translation from English. The runs marked with a star (*) were the ones selected for manual validation. For the submissions from Unbabel, runs "1" are the ones identified as "Greedy", and runs "2" are the ones for "MBR".

| Team | Run | Metric | de2en | fr2en | it2en | pt2en | es2en | ru2en |
|---|---|---|---|---|---|---|---|---|
| ADAPT | 1 | BLEU | *32.24 | 18.81 | | | | |
| | | COMET | 83.04 | 72.14 | | | | |
| ADAPT | 3 | BLEU | 36.93 | | | | | |
| | | COMET | 78.84 | | | | | |
| AIST | 1 | BLEU | 45.86 | | | | | |
| | | COMET | 84.65 | | | | | |
| AIST | 2 | BLEU | *45.92 | | | | | |
| | | COMET | 84.84 | | | | | |
| DCU | - | BLEU | 32.60 | 31.47 | 28.40 | 31.32 | 28.02 | 25.76 |
| | | COMET | 78.99 | 78.74 | 79.63 | 79.56 | 80.90 | 70.01 |
| HW-TSC | 1 | BLEU | *45.79 | | | | | |
| | | COMET | 83.98 | | | | | |
| HW-TSC | 2 | BLEU | 45.68 | | | | | |
| | | COMET | 83.86 | | | | | |
| HW-TSC | 3 | BLEU | 45.43 | | | | | |
| | | COMET | 84.08 | | | | | |
| Unbabel | 1 | BLEU | 49.05 | 53.29 | 38.91 | 55.63 | 51.32 | 47.28 |
| | | COMET | 86.67 | 86.05 | 85.32 | 85.11 | 86.99 | 83.82 |
| Unbabel | 2 | BLEU | *46.72 | *51.67 | *38.91 | *53.53 | *52.28 | *45.11 |
| | | COMET | 86.97 | 86.39 | 85.32 | 85.47 | 87.25 | 83.95 |
| Baseline | - | BLEU | 45.85 | 54.79 | 37.49 | 51.38 | 53.54 | 43.70 |
| | | COMET | 86.39 | 86.11 | 85.28 | 85.08 | 87.18 | 83.37 |

Table 4: BLEU and COMET scores for submissions to the Biomedical Task submission system, for translation into English. The runs marked with a star (*) were the ones selected for manual validation. For the submissions from Unbabel, runs "1" are the ones identified as "Greedy", and runs "2" are the ones for "MBR".

submissions are not directly comparable to the ones from the previous year since, for the first time, we ran an evaluation on the abstract level.

## 5.2 OCELoT Biomedical Translation task

Only one team (AIST) submitted to the biomedical task in OCELoT, but also for the same language pairs in our submission system and for our test

| Team | Run | Metric | en2de | de2en |
|------|-----|--------|-------|-------|
| AIST | 517 | BLEU | 28.30 | |
| | | COMET | 83.75 | |
| | 542 | BLEU | | 39.68 |
| | | COMET | | 82.55 |
| | 544 | BLEU | | 39.68 |
| | | COMET | | 82.55 |
| | 545 | BLEU | 28.30 | |
| | | COMET | 83.75 | |

Table 5: BLEU scores for submissions to OCELoT for the Biomedical Translation Task.

suites in the general task. While their results as shown in Table 5 were similar to the ones in Table 6, they were slightly inferior to the ones that the same team obtained for the runs to our submission system, e.g., for en2de, a BLEU score of 28.30 versus 33.80, and a COMET score of 85.59 versus 83.75.

### 5.3 OCELoT General Machine Translation task

We included test suites only for the language pairs in our task that overlap with the ones considered in the general task, namely, en2de, en2es, and en2ru. The scores for the submissions to the general task (cf. Table 6) varied much more than the ones submitted directly to the biomedical task (cf. Table 3), from very low to high, e.g., BLEU scores of 1.63 (certainly due to mistakes in the system) to 52.56. It is safe to assume that most systems were not trained especially for the biomedical domain. In spite of this, we observed some submissions with scores even higher than the ones for the biomedical task. Amongst the submissions to the general task, the highest scores for en2de were 38.07 BLEU (ONLINE-W) and 88.25 COMET (TranssionMT), as opposed to a BLEU score of 34.22 (Unbabel run1) and a COMET score 88.09 (Unbabel run2) in the biomedical task. For en2ru, the highest scores in the general task were 41.25 BLEU (Claude-3.5) and 89.88 COMET (Claude-3.5 and Unbabel-Tower70B), as opposed to 35.76 BLEU (Unbabel run1) and 89.71 COMET (Unbabel run2). Therefore, submissions from the same team (Unbabel) scored slightly higher in the general task than in the biomedical task.

## 6 Manual evaluation

Similar to previous years, we performed manual validation of a sample of the submissions for most of the language pairs. The number of abstracts that

we considered for each language was of either 10 or 20 depending on the availability of the human evaluators. We used the three-way function of the Appraise tool (Federmann, 2018), which includes the following elements:

- the abstract in the original language (e.g., English for en2fr);

- translation A: first translation in the target language (e.g. French for en2fr);

- translation B: second translation in the target language (e.g. French for en2fr).

The task consists of validating whether a translation is better than the other (i.e., A>B or A<B), or whether they are of similar quality (A=B). In cases where the evaluators notice that an error might have occurred, e.g., translation from another text or a translation shorter than it should be, it is possible to skip the validation of this particular pair.

For all language pairs, we considered the best run from each of the team that submitted directly to the biomedical task. The best run was the one identified by the participants during the submission process. Otherwise, we selected the best performing one. We evaluated pairs of either two translations from the teams, or one translation from a team and the reference translation. We present the results for submissions from English in Table 7 and for submission into English in Table 8.

We present below a summary of the mistakes that we observed during manual evaluation.

**en2fr** Translation quality was uneven, as suggested by the 20 point difference in BLEU scores obtained by the systems. While some translations were of very high quality, others exhibited serious issues including conveying meaning drastically different from the original sentence. In example 1, numerical values are erroneous and inconsistent with the corresponding percentages. In Example 2 the resulting translation is medically unacceptable.

(1) **en:** Of the 273 patients, <u>164</u> (60.1%) required invasive mechanical ventilation. <u>One hundred and forty-two</u> patients (52.0%) survived their hospital stay.
**fr\*:** Sur les 273 patients, <u>104</u> (60,1%) ont nécessité une ventilation mécanique invasive et <u>164</u> (52,0%) ont survécu à leur séjour à l'USI.
**fr:** Parmi les 273 patient·es, <u>164</u> (60,1 %) ont nécessité une ventilation mécanique invasive.

| Teams | en2de | | en2es | | en2ru | |
|---|---|---|---|---|---|---|
| | **BLEU** | **COMET** | **BLEU** | **COMET** | **BLEU** | **COMET** |
| AIST-AIRC | 28.28 | 84.85 | | | | |
| Aya23 | 30.77 | 87.11 | 49.49 | 85.32 | 31.90 | 86.69 |
| CUNI-DS | | | | | 27.93 | 86.96 |
| CUNI-NL | 20.06 | 83.38 | | | | |
| Claude-3.5 | 35.23 | 87.86 | 52.08 | 85.93 | 41.25 | 89.88 |
| CommandR-plus | 32.44 | 87.67 | 49.84 | 85.78 | 34.33 | 88.64 |
| CycleL | 1.32 | 38.35 | 3.00 | 45.17 | 0.32 | 34.65 |
| CycleL2 | 1.32 | 38.35 | | | 0.10 | 28.49 |
| Dubformer | 31.19 | 83.49 | 40.65 | 78.58 | 1.94 | 39.58 |
| GPT-4 | 35.80 | 87.93 | 51.53 | 85.85 | 34.00 | 88.45 |
| IKUN-C | 10.82 | 78.34 | 22.18 | 78.23 | 12.69 | 81.74 |
| IKUN | 11.07 | 79.14 | 12.67 | 74.02 | 13.28 | 82.98 |
| IOL_Research | 30.86 | 87.17 | 48.90 | 85.56 | 32.30 | 87.68 |
| Llama3-70B | 31.43 | 87.01 | 47.86 | 85.30 | 32.18 | 88.05 |
| MSLC | 25.17 | 82.24 | 46.30 | 84.27 | | |
| NVIDIA-NeMo | 15.91 | 80.21 | 30.00 | 79.32 | 20.37 | 83.28 |
| ONLINE-A | 36.09 | 87.34 | 52.56 | 85.62 | 40.20 | 89.23 |
| ONLINE-B | 36.48 | 88.21 | 51.56 | 85.13 | 40.23 | 88.73 |
| ONLINE-G | 34.86 | 87.08 | 50.98 | 85.34 | 37.22 | 89.44 |
| ONLINE-W | 38.07 | 88.04 | 52.47 | 85.78 | 39.77 | 89.52 |
| Occiglot | 6.33 | 70.19 | 31.93 | 78.52 | | |
| TSU-HITs | 1.63 | 37.00 | 17.23 | 60.20 | 2.80 | 52.36 |
| TranssionMT | 36.57 | 88.25 | 52.67 | 85.67 | 40.07 | 88.76 |
| Unbabel-Tower70B | 32.37 | 87.89 | 47.93 | 86.12 | 32.61 | 89.88 |
| Yandex | | | | | 35.09 | 89.81 |

Table 6: BLEU scores for submissions to OCELoT for the General Machine Translation Task.

| Languages | Systems | Abstracts | | | |
|---|---|---|---|---|---|
| | | **A>B** | **A=B** | **A<B** | **skipped** |
| **en2de** | AIST vs. ADAPT | 3 | 3 | 12 | 2 |
| | AIST vs. HW-TSC | 13 | 2 | 4 | 1 |
| | AIST vs. DCU | 10 | 3 | 4 | 3 |
| | AIST vs. reference | 2 | 7 | 10 | 1 |
| | AIST vs. Unbabel | 2 | 5 | 12 | 1 |
| | ADAPT vs. HW-TSC | 16 | 2 | 0 | 2 |
| | ADAPT vs. DCU | 10 | 5 | 1 | 4 |
| | ADAPT vs. reference | 0 | 8 | 10 | 2 |
| | ADAPT vs. Unbabel | 2 | 10 | 6 | 2 |
| | HW-TSC vs. DCU | 6 | 2 | 9 | 3 |
| | HW-TSC vs. reference | 0 | 0 | 19 | 1 |
| | HW-TSC vs. Unbabel | 0 | 1 | 18 | 1 |
| | DCU vs. reference | 0 | 3 | 14 | 3 |
| | DCU vs. Unbabel | 0 | 2 | 15 | 3 |
| | reference vs. Unbabel | 2 | 10 | 7 | 1 |
| **en2fr** | reference vs. Unbabel | 14 | 0 | 5 | 1 |
| | reference vs. ADAPT | 17 | 0 | 2 | 1 |
| | Unbabel vs. ADAPT | 18 | 0 | 1 | 1 |
| **en2it** | reference vs. DCU | 5 | 1 | 13 | 1 |
| | reference vs. Unbabel | 1 | 1 | 18 | 0 |
| | DCU vs. Unbabel | 4 | 6 | 9 | 1 |
| **en2pt** | DCU vs. Unbabel | 0 | 6 | 8 | 6 |
| | DCU vs. reference | 4 | 7 | 3 | 6 |
| | Unbabel vs. reference | 7 | 10 | 3 | 0 |
| **en2ru** | reference vs. Unbabel | 4 | 2 | 4 | 0 |
| | reference vs. DCU | 3 | 3 | 4 | 0 |
| | Unbabel vs. DCU | 7 | 2 | 1 | 0 |

Table 7: Pairwise manual evaluation results for the test set (from English).

| Languages | Systems | Abstracts | | | |
|---|---|---|---|---|---|
| | | A>B | A=B | A<B | skipped |
| **de2en** | DCU vs. AIST | 3 | 2 | 2 | 3 |
| | DCU vs. Unbabel | 2 | 2 | 3 | 3 |
| | DCU vs. reference | 2 | 2 | 3 | 3 |
| | DCU vs. HW-TSC | 5 | 0 | 2 | 3 |
| | DCU vs. ADAPT | 6 | 0 | 1 | 3 |
| | AIST vs. Unbabel | 1 | 0 | 9 | 0 |
| | AIST vs. reference | 3 | 2 | 5 | 0 |
| | AIST vs. HW-TSC | 4 | 3 | 3 | 0 |
| | AIST vs. ADAPT | 8 | 0 | 2 | 0 |
| | Unbabel vs. reference | 8 | 2 | 0 | 0 |
| | Unbabel vs. HW-TSC | 10 | 0 | 0 | 0 |
| | Unbabel vs. ADAPT | 10 | 0 | 0 | 0 |
| | reference vs. HW-TSC | 6 | 1 | 3 | 0 |
| | reference vs. ADAPT | 6 | 2 | 2 | 0 |
| | HW-TSC vs. ADAPT | 5 | 1 | 4 | 0 |
| **fr2en** | DCU vs. ADAPT | 6 | 0 | 4 | 0 |
| | DCU vs. reference | 1 | 2 | 7 | 0 |
| | DCU vs. Unbabel | 0 | 0 | 10 | 0 |
| | ADAPT vs. reference | 0 | 2 | 8 | 0 |
| | ADAPT vs. Unbabel | 0 | 0 | 10 | 0 |
| | reference vs. Unbabel | 1 | 3 | 6 | 0 |
| **it2en** | reference vs. Unbabel | 0 | 5 | 15 | 0 |
| | reference vs. DCU | 5 | 3 | 8 | 4 |
| | Unbabel vs. DCU | 11 | 4 | 1 | 4 |
| **es2en** | DCU vs. reference | 4 | 2 | 9 | 5 |
| | DCU vs. Unbabel | 3 | 4 | 8 | 5 |
| | reference vs. Unbabel | 5 | 6 | 9 | 0 |
| **ru2en** | reference vs. Unbabel | 2 | 2 | 5 | 1 |
| | reference vs. DCU | 4 | 0 | 4 | 2 |
| | Unbabel vs. DCU | 4 | 3 | 1 | 2 |

Table 8: Pairwise manual evaluation results for the test set (into English).

Cent quarante-deux personnes (52,0 %) ont survécu à leur séjour à l'hôpital.

(2) **en:** Deaths by mechanical asphyxiation constitute a social drama
**fr*:** La prévention constitue un drame social
**fr:** Les morts par asphyxies mécaniques constituent un drame social

In both cases, the translation errors likely result from mixing information contained in different parts of the original texts. Arguably, this is very concerning because users of such a translation system could conclude that the erroneous translations are correct by checking that the information is present in the original text. Other issues are more easily detected, such as the interruption of the translation by a loop repetition of a set of tokens (e.g., *une mobilité allant de 5,6% à 5,6% à 1211% à 1211% à 1211% à 1211% à 1211% à 1211% à 1211% à 1211%...*).

The choice of having full abstract translation instead of sentence-by-sentence translation this year seems to have both a positive impact on the overall consistency of translations (e.g., overall consistent use of terms and acronyms throughout a document) and a negative impact on the end of translation for some systems, where translation quality was decreasing as the text unfolded and sometimes just interrupted (with or without loop repetitions).

Specialized term translation was sometimes erroneous, in particular with terms referring to animal species (for example, translating *waterfowl* by *oiseaux d'eau* instead of *sauvagine*), which were more frequent this year due to the selection method for the test documents. Polysemous terms were also a source of erroneous translations (e.g., *hood* translated as *capot* – car context instead of *capuche*, which is correct in a clothing context).

In addition to the manual evaluation through appraise, a complementary assessment of the best system submission outputs was conducted, with a focus on *Acronyms* and *Lab Values*, consistently with the evaluations conducted in the two previous years. Overall, 31 out of 50 test documents contained

acronyms and none contained lab values. Acronym translations were considered correct when the system translation was identical to the reference translation or consisted of an attested acronym use in a similar context. Correct acronym translations (79%) included frequent acronyms such as USI (*Unité de Soins Intensifs* – Itensive Care Unit) or IC (*Intervalle de Confiance* – confidence interval). In other cases, acronyms were either untranslated (16%) or erroneous (5%). Some of the acronym translation strategies used by human translators and not by machine translation consist of explicitly stating that an English acronym is used, for example: *la santé mentale du nourrisson (IMH en anglais).* This is sometimes combined with a strategy of using the long form of a term in French, when an acronym was used in English. These strategies are often used with acronyms that stand for infrequent terms.

It is also interesting to notice that reference translations contain idiomatic linguistic traits not used in machine-translated text, such as inclusive writing (as seen in Example 1).

**en2pt**  All translations into Portuguese were of very good quality, except for some empty translations from one submission and the remains of the prompt used, which were included in the translations of the same submission. Therefore, the decision of whether one translation was better than the other was generally based on small details, often one single mistake.

Small mistakes that we found were the following: (a) lack of capitalization at the start of the sentence (e.g., "... profunda (TVP). o sangue ..."); (b) nominal concordance (e.g., "o febre pós-anestésica"); (c) missing words (e.g., "com uma [força] média de 526N"); (d) words that remained in English (e.g., "odds ratio") (e) typos ("registe" instead of "registre"); (f) and grammatical mistakes (e.g., "acompanhou [por] mais de 18 meses").

As in previous years, we found mistakes related to the non-translation of acronyms. For easier or more common terms, e.g., Artificial intelligence (AI), the translations were all correct, i.e., "inteligência artificial (IA)". However, mistakes were often found for other terms, as in Example 3 below in which only the translation pt$_3$ is correct and has the right acronym:

(3) **en:** Computer vision (CV)
 **pt$_1$:** visão por computador (CV)

**pt$_2$:** visão computacional (CV)
**pt$_3$:** visão computacional (VC)

Often we observed a copy of the English acronym for much more complex terms, as in Example 4:

(4) **en:** hydrogenated castor oil (HCO ethoxylates)
 **pt$_1$:** Óleo de castor hidrogenado polioxietileno (etoxilações de HCO)
 **pt$_2$:** óleo de rícino hidrogenado de polioxietileno (HCO-etoxilados)
 **pt$_3$:** hydrogenated castor oil (HCO ethoxylates)

However, we had some difficult examples in which the translation and acronym were correct, e.g. pt$_2$: in Example 5:

(5) **en:** hospital standardized mortality ratio (HSMR)
 **pt$_1$:**  taxa de mortalidade hospitalar padronizada (HSMR, na sigla em inglês)
 **pt$_2$:**  razão de mortalidade hospitalar padronizada (RMHP)

Finally, we observed many examples in which we favored some translation over others because they either sounded better or more correct, namely, translations pt$_2$: in Examples 6, 7, 8, and 9:

(6) **en:** was highly expressed in CTCs
 **pt$_1$:** foi altamente expresso em CTCs
 **pt$_2$:** tinha uma expressão elevada nas CTCs

(7) **en:**  A quasi-experimental study,  which compared
 **pt$_1$:** Estudo quase-experimental, que comparou
 **pt$_2$:**  Um estudo quase experimental,  que comparou

(8) **en:** Case signalment
 **pt$_1$:** Fatores de identificação do caso
 **pt$_2$:** O sinalamento do caso

(9) **en:** axis of the femoral neck
 **pt$_1$:** eixo do colo do fêmur
 **pt$_2$:** eixo do colo femoral

**fr2en** With the change in protocol this year (from sentence-level to paragraph-level translation and evaluation), there were several differences in the observed quality of translations.

Translation issues brought up in previous years remained present, namely the copying or wrong translation of acronyms and specialised terms (Example 10), the wrong translation of personal pronouns (e.g. *son* 'his/her/their' in Example 11) and errors linked to the ambiguity of source terms (e.g. *taille* 'height or waist' in Example 12).

(10) **fr:** la thérapie de substitution de la nicotine (TSN)
**en:** nicotine replacement companies (NTS)
**en\*:** nicotine replacement therapy (NRT)

(11) **fr:** …la capacité d'un individu à rechercher des soins … pour son animal de compagnie
**en:** an individual's ability to seek … care for their companion animal
**en\*:** an individual's ability to seek … care for his companion animal

(12) **fr:** la circonférence de la taille (CT)
**en:** waist circumference (WC)
**en\*:** circumference of height (CT)

However, the overall quality of the translations was visibly lower than in previous years, due to the use of LLMs and the translation of whole paragraphs rather than individual sentences. LLMs tended to exhibit more volatile behaviour, often copying the source document instead of translating, and also including the initial prompt in the output. The consequence of the longer documents to translate was mostly seen in skipping sentences within the documents or (more commonly) at the end of documents (i.e. translation finishing too early or repeating the final sentence multiple times). We also observed the merging of multiple sentences/clauses into a single one and the negative influence of previous sentences on later translations, resulting in the repetition of terms in inappropriate places and errors in the translation of numbers (both problems illustrated in Example 13).

(13) **fr:** Cent six médecins ont répondu au sondage et 12 ont participé à un entretien
**en:** One hundred and six physicians responded to the survey and 12 participated in an interview
**en\*:** One hundred and twelve respondents participated in the survey

The consequence of the appearance of these more serious errors (i.e. non-translation, missing parts of the translation etc.) meant that they often formed the basis of the evaluation rather than distinctions being based on errors more traditionally resulting from the translation of scientific texts (terminology, acronyms, etc.). Not evaluating on the sentence level meant that an improved translation on the sentence level was easily overridden by a more technical problem, such as a missing sentence at the end of the document. It could be useful in the following years to consider evaluation via error analysis to get more detailed insights into the strengths and weaknesses of different systems on a more granular level.

**es2en** Contrary to past years, the Spanish to English language pair had very few contributions, totalling 30 examples from two different MT models both compared between each other and against a reference human translation.

In the past, sentence-to-sentence translation has provided good results in terms of translation quality at sentence level. However, the trade-off was inconsistency in the usage of medical terminology and medical specific acronyms. This year however, the use of full abstracts for translation led to greater consistency in the translation of terminology and acronyms specific to medicine.

When working well, the MT output has a good quality, sometimes producing a result that was comparable to human translation in terms of quality, as shown in Table 8, where the MT system Unbabel had very good results compared against DCU and the reference translation.

However, the MT output still lacks the fluency of a human translation, as the systems had a tendency to replicate the structure of the original Spanish source text, resulting in translations that can be considered "literal translations". In many instances, the MT output would require copy editing and rewriting by a native English speaker to render the text more fluent and increase the overall quality of the output.

Despite the good quality level of some translations, the overall quality of the outputs for this year's challenge is very uneven, with some very good abstracts in English and some abstracts that were not translated or still contained Spanish words in them.

At least one of the system used LLMs to produce the output in English, with this prompt: "1. While being factual, accurate and not missing out any detail, translate the given Spanish text into the specified English language. Spanish Text:". The use of the prompt ensured the output did not miss information from the original source text, as has sometimes been the case in past years. Nevertheless, the LLM system was not very robust.

As shown in the example below, the LLM system sometimes did not translate the text in English as requested. The text remained in Spanish. That is considered a missing translation and is considered a major error.

(14) **en:** While being factual, accurate and not missing out any detail, translate the given Spanish text into the specified English language. La prevalencia de alergia alimentaria ha aumentado en algunas regiones del mundo, y con ello la incidencia, según la variabilidad geográfica, en el fenotipo y manifestaciones clínicas...

Another error the LLM system made was the inclusion of the prompt used to generate the translated output as part of the response. This add superfluous information to the English translation and breaks the readability and fluency of the text (see previous example).

As mentioned before, fluent translation was still an issue for the machine translation system, in particular for the DCU system. This system sometimes generated sentences that were clunky or ungrammatical in English.

(15) **es:** Se registraron 4 casos de morbilidad post punción (2 dolores epigástricos y 2 hematomas de pared abdominal

**en:** Were registered 4 cases of morbidity post puncture (2 pain epigastric and 2 hematomas of abdominal wall).

In conclusion, LLMs systems still seem to have an unreliable performance when it comes to machine translation, producing very good quality translations, missing translations or ungrammatical translations at the same time. A better out-of-box LLM or refine the prompting techniques might obtain better results with these systems.

It must be noted, however, that there were very few examples for the Spanish to English translation to reach an indisputable conclusion.

**en2de** Similar to previous years, a generally high level of translation quality was seen for English to German translation. The strongest models produced translations that not only conveyed the content well but also maintained consistency in terms of style and structure. However, certain systems exhibited notable flaws. In particular, one model consistently omitted portions of the text, often truncating the translation towards the end of the document or, at times, even mid-sentence. Another system struggled with basic capitalization, failing to begin sentences with an uppercase letter, which detracted from the overall readability of the output.

Numerical translations were also an issue, with *Eighty-nine* frequently mistranslated as either *Achtundachtzig* "eighty-eight" or *Achtundneunzig* "ninety-eight", revealing inaccuracies in number handling. The translation of abbreviations varied across systems, with some attempting to expand or translate them, occasionally resulting in errors. For example, the *European Commission (EC)* was incorrectly translated as *EG (Europäische Gemeinschaft)* instead of EU. Furthermore, specialized terminology presented additional challenges, with terms like *compulsory elective* rendered awkwardly as *obligatorische elektive Veranstaltung* rather than the more appropriate *Wahlpflichtkurs*.

Grammatical errors also persisted in some translations, indicating that while overall quality was high, there is still room for improvement in handling both sentence structure and more nuanced linguistic elements.

**de2en** Overall, results varied for the German-to-English translation task. While at least one system was able to provide a human-level translation for each source sample, there was generally also at least one translation that was either incomplete or difficult to understand.

The most serious mistakes included omission of whole sentences, or synthesis of text that was not present in the original. This was especially evident in cases where the sample text ended in an incomplete sentence, which caused some systems to generate a completion to the sentence. In the most egregious example of this phenomenon, an incomplete sentence at the end of a description of an animal's skin condition after an insect bite led to more than one translation mentioning euthanasia, when no such language was present in the source. In some instances, text would be translated to nonexistent words, e.g. translation of *porös* to

the nonexistent word *sporeous*. Other mistranslations included rendering *mittleren Werte* as *median* instead of *mean* values as was intended in the text.

The most frequently occurring mistakes were related to the capitalization of words at the beginning of sentences. Other formatting mistakes failed to take into account the structure of the text, omitting paragraph headings. These mistakes did not affect the overall intelligibility of the text.

All in all, the majority of the systems were able to provide a translation that, while not perfect, was understandable and correctly conveyed important information.

**en2it**  The quality of the translation was higher than in previous years, even more so than last year, which set a new threshold in the accuracy of the translation from English to Italian and vice versa. The quality of most of the abstract was almost identical and fluent in terms of the quality of language. The terminology and the syntax was of very high quality in both translation directions. There were rarely major issues with the choice of terms or the construction of the sentences.

One mistake was the addition of parts of the text that were not present in the original version. For example, the original version is "Among those diagnosed with COVID-19 during follow-ups between March 2020 and March 2021 [...]"

While the Italian translation: "MATERIALE E METODO: TRA marzo 2020 e maggio 2021, sono stati analizzati [...]"

Where there is the addition of "MATERIALE E METODO". There is also some minor issue with the punctuation (the semicolon between "rene" and "o dobbiamo farlo" should not be there) as well as uppercase letters ("TRA" instead of "tra").

There were two problems concerning the cause effect or correlation among pathologies. For example, in the original English version: "Chronic rhinosinusitis with nasal polyps is a common disease with still unclear pathophysiologic mechanisms." The "Chronic rhinosinusitis with nasal polyps" are one thing all together that is documented to be a common disease.

On the other hand, the Italian version: "La rinosinusite cronica e la poliposi nasale sono patologie frequenti" the "Rinusite cronica" ("Chronic rhinosinusitis) and "poliposi nasale" ("nasal polyps") are considered as two distinct pathologies.

The other example happens with the following sentence: "The airway epithelial barrier has been shown to be involved in different chronic disorders, including rhinitis, nasal polyposis and asthma" and its Italian translation: "La barriera epiteliale delle vie respiratorie sembra essere coinvolta in diverse patologie croniche come la rinite, la poliposi nasale e l'asma"

In this case, the translation gives a slightly different interpretation of the fact that, in the original version, "airway epithelial barrier has been shown to be [...]" as in "it has been demonstrated that", while the Italian "sembra essere coinvolta" ("seems to be involved") shoes a less strong connection between the entities (airway epithelial barrier and chronic disorders).

**it2en**  For the Italian to English translation direction, we observe an opposite problem compared to the English one that is removing a part of the text.

For example, in the original "Conclusione: sebbene non abbiamo riscontrato differenze significative tra i pazienti sottoposti a gastrectomia standard e quelli sottoposti a NACT prima della gastrectomia, [...]" we have "Conlusione:" as the initial part of this sentence.

In the English version, we have "Although we found no significant difference between the patients undergoing standard gastrectomy and those undergoing NACT before gastrectomy," Where "Conclusions" ("conclusione") is missing.

From Italian to English, there was a missing agreement in gender for the translation of the following sentence: "A total of 192 female feral cats were investigated for a large-scale trap-neuter-release program." One of the Italian translations overlooked the female gender with: "Un totale di 192 gatti selvatici sono stati studiati per un ampio programma di trappola, sterilizzazione e rilascio." Where "gatti" is the masculine plural of a cat which, in this case, is wrong.

Another type of wrong accordance was found in the translation of the following sentence: "La gangrena di Fournier è una fascite necrotizzante a rapida progressione che coinvolge il perineo, le regioni perianale e genitali e costituisce una vera emergenza chirurgica con un tasso di mortalità potenzialmente elevato" where the English version: "Fournier's gangrene is a rapidly progressing necrotizing fasciitis involving the perineal, perianal, or genital regions and constitutes a true surgical emergency with a potentially high mortality rate." considers the "perineal [...] region" instead of the "perineum" alone.

**en2ru and ru2en** This year, two systems, Unbabel and DCU, participated in the Biomedical Machine Translation task. Generally, the translations to and from English were of high quality. We did not encounter examples that were completely unacceptable, aside from a few cases where text boundaries were mapped incorrectly. Compared to previous years, we observed a general improvement in how the systems handled abbreviations, which is a notable challenge in biomedical translation.

This year translations were evaluated at the abstract level, and at times determining which translation was superior often came down to small details. In some instances, we preferred one translation over another purely due to stylistic differences. There were only a handful of cases where the systems diverged significantly in quality. Overall, Unbabel outperformed DCU, as reflected by manual evaluation (Table 7 and 8) and better BLEU and COMET scores (Tables 3 and 4).

## 7 Conclusions

We presented the results for this year's edition of the Biomedical Translation Task at WMT, in which we considered 12 language pairs. In this paper, we described the development of the test sets, the submissions we received, our baseline system, and the details about the automatic and manual evaluation. Different from previous years, we did not split and align the sentences, instead we had the test sets simply composed of the title and abstracts of the articles.

## Limitations

Concerning the quality of the extracted test sets, the passage from sentence to paragraph level is likely to require additional post-processing in future years. Whereas in previous years, sentence alignment resulted in additional validation of the extraction process, a number of errors were present in the test sets this year, resulting in more skipped evaluations. These included (i) missing or additional sentences in the reference translations with respect to the source texts, (ii) the truncation of certain sentences after special characters and subscript text, the inconsistent inclusion of headers (e.g. *Methods*, *Results*) in the abstracts and the non-capitalised of accented characters in the headers (e.g. French *RéSUMé* 'Abstract' instead of *RÉSUMÉ*), a consequence of the original source text, but which could be corrected in a post-processing step.

## References

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *Preprint*, arXiv:2402.17733.

Rachel Bawden, Kevin Bretonnel Cohen, Cristian Grozea, Antonio Jimeno Yepes, Madeleine Kittner, Martin Krallinger, Nancy Mah, Aurelie Neveol, Mariana Neves, Felipe Soares, Amy Siu, Karin Verspoor, and Maika Vicente Navarro. 2019. Findings of the WMT 2019 biomedical translation shared task: Evaluation for MEDLINE abstracts and biomedical terminologies. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 29–53, Florence, Italy. Association for Computational Linguistics.

Rachel Bawden, Giorgio Maria Di Nunzio, Cristian Grozea, Inigo Jauregi Unanue, Antonio Jimeno Yepes, Nancy Mah, David Martinez, Aurélie Névéol, Mariana Neves, Maite Oronoz, Olatz Perez-de Viñaspre, Massimo Piccardi, Roland Roller, Amy Siu, Philippe Thomas, Federica Vezzani, Maika Vicente Navarro, Dina Wiemann, and Lana Yeganova. 2020. Findings of the WMT 2020 biomedical translation shared task: Basque, Italian and Russian as new additional languages. In *Proceedings of the Fifth Conference on Machine Translation*, pages 660–687, Online. Association for Computational Linguistics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198, Berlin, Germany. Association for Computational Linguistics.

Antonio Castaldo, Maria Zafar, Prashanth Nayak, Rejwanul Haque, Andy Way, and Johanna Monti. 2024. The SETU-ADAPT Submission for WMT 24 Biomedical Shared Task. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*, Miami, USA. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Christian Federmann. 2018. Appraise evaluation framework for machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 86–88, Santa Fe, New Mexico. Association for Computational Linguistics.

Antonio Jimeno Yepes, Aurélie Névéol, Mariana Neves, Karin Verspoor, Ondřej Bojar, Arthur Boyer, Cristian Grozea, Barry Haddow, Madeleine Kittner, Yvonne Lichtblau, Pavel Pecina, Roland Roller, Rudolf Rosa, Amy Siu, Philippe Thomas, and Saskia Trescher. 2017. Findings of the WMT 2017 biomedical translation shared task. In *Proceedings of the Second Conference on Machine Translation*, pages 234–247, Copenhagen, Denmark. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda,

Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Findings of the WMT24 general machine translation shared task: the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. 2023. Mega: Moving average equipped gated attention. *Preprint*, arXiv:2209.10655.

Mariana Neves, Antonio Jimeno Yepes, Aurélie Névéol, Rachel Bawden, Giorgio Maria Di Nunzio, Roland Roller, Philippe Thomas, Federica Vezzani, Maika Vicente Navarro, Lana Yeganova, Dina Wiemann, and Cristian Grozea. 2023. Findings of the WMT 2023 biomedical translation shared task: Evaluation of ChatGPT 3.5 as a comparison system. In *Proceedings of the Eighth Conference on Machine Translation*, pages 43–54, Singapore. Association for Computational Linguistics.

Mariana Neves, Antonio Jimeno Yepes, Aurélie Névéol, Cristian Grozea, Amy Siu, Madeleine Kittner, and Karin Verspoor. 2018. Findings of the WMT 2018 biomedical translation shared task: Evaluation on Medline test sets. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 324–339, Belgium, Brussels. Association for Computational Linguistics.

Mariana Neves, Antonio Jimeno Yepes, Amy Siu, Roland Roller, Philippe Thomas, Maika Vicente Navarro, Lana Yeganova, Dina Wiemann, Giorgio Maria Di Nunzio, Federica Vezzani, Christel Gerardin, Rachel Bawden, Darryl Johan Estrada, Salvador Lima-lopez, Eulalia Farre-maduel, Martin Krallinger, Cristian Grozea, and Aurelie Neveol. 2022. Findings of the WMT 2022 biomedical translation shared task: Monolingual clinical case reports. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 694–723, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff

Wang. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Zhanglin Wu, Daimeng Wei, Zongyao Li, Zhengzhe Yu, Shaojun Li, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Yuhao Xie, Lizhi Lei, Hao Yang, and Yanfei Jiang. 2023. The path to continuous domain adaptation improvements by HW-TSC for the WMT23 biomedical translation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 271–274, Singapore. Association for Computational Linguistics.

Lana Yeganova, Dina Wiemann, Mariana Neves, Federica Vezzani, Amy Siu, Inigo Jauregi Unanue, Maite Oronoz, Nancy Mah, Aurélie Névéol, David Martinez, Rachel Bawden, Giorgio Maria Di Nunzio, Roland Roller, Philippe Thomas, Cristian Grozea, Olatz Perez-de Viñaspre, Maika Vicente Navarro, and Antonio Jimeno Yepes. 2021. Findings of the WMT 2021 biomedical translation shared task: Summaries of animal experiments as new test set. In *Proceedings of the Sixth Conference on Machine Translation*, pages 664–683, Online. Association for Computational Linguistics.

# MSLC24 Submissions to the General Machine Translation Task

**Samuel Larkin**        **Chi-kiu Lo 羅致翹**        **Rebecca Knowles**
Digital Technologies Research Centre
National Research Council Canada (NRC-CNRC)
{samuel.larkin,chikiu.lo,rebecca.knowles}@nrc-cnrc.gc.ca

## Abstract

The MSLC (Metric Score Landscape Challenge) submissions for English–German, English–Spanish, and Japanese–Chinese are constrained systems built using Transformer models for the purpose of better evaluating metric performance in the WMT24 Metrics Task. They are intended to be representative of the performance of systems that can be built relatively simply using constrained data and with minimal modifications to the translation training pipeline.

## 1 Introduction

Lo et al. (2023) introduced the Metric Score Landscape Challenge (MSLC) dataset for the WMT23 Metrics Task, with the goal of examining automatic MT evaluation metric performance across a wider range of quality. That work found unexpected behaviours in several MT metrics, by examining performance across a wide range of quality and by analyzing metric characteristics other than correlation. A major limitation of that work was that there was no human evaluation of the medium- to low-quality MT outputs that were included in the MSLC dataset. To resolve this disconnect between the high-quality WMT systems and the core MSLC systems, we submit the higher performing end of the MSLC systems to the WMT General MT task for human evaluation. The systems described here are not highly-competitive systems, and are useful primarily for their purpose in evaluating metrics.

We build MSLC models for three language pairs: English→German (eng→deu), English→Spanish (eng→spa), and Japanese→Chinese (jpn→zho). All models are sentence-level models that handle paragraph- or document-level translation by performing sentence splitting, translation, and then concatenating the translated sentences. They are built without any additional modifications to the Transformer architecture and without additional components like backtranslation, tagging, factors,

or domain-specific features (with one exception for preprocessing input in the Japanese→Chinese speech domain). The English→German model is the same model described in Lo et al. (2023). The English→Spanish model uses language identification for training data filtering. The Japanese→Chinese model incorporates additional postprocessing.

In the remainder of this system description paper, we describe the data used (Section 2), the preprocessing and postprocessing performed (Section 3), and the models trained (Section 4) for our submissions for the three language pairs. Using the human evaluations produced by the Metrics task, we use the MSLC systems as a case study of some risks of the new automatic metric-based pre-selection of systems for human annotation at the General MT task (Section 5).

## 2 Data

We retrieved the corpora using the provided tool mtdata==0.4.1 (Gowda, 2024) for eng→spa and jpn→zho and reused what we had downloaded (without the use of the tool) from the 2023 data download table for eng→deu.

### 2.1 English→German

We re-used the English→German model from Lo et al. (2023), and refer the reader to that paper for full details of the training data used. The *newstest2020* data was used for validation, and the training corpora were downloaded from the WMT 2023 General Machine Translation download table.[1]

---

[1] https://www2.statmt.org/wmt23/translation-task.html#download. Note that this includes News Commentary v18.1 rather than v16, which the download tool delivered. By email communication with the organizers, we confirmed that both versions were permitted for the constrained track.

## 2.2 English→Spanish

We used some of the available corpora for the General Machine Translation constrained track[2] and filtered based on language ID (due to large amounts of target-side English in some training corpora). We opted not to use *OPUS-multiccaligned-v1*, *ParaCrawl-paracrawl-9*, *Statmt-ccaligned-1* and *Statmt-commoncrawl_wmt13-1*, due to known issues of noise in web-crawled corpora; for more discussion see, i.a., Khayrallah and Koehn (2018); Lo et al. (2018); Kreutzer et al. (2022). The full set of corpora used is shown in Table 1.

As a first filtering step, we kept sentence pairs where sentences have less than or equal to 4000 characters and less or equal to 200 words. We then proceeded with a second filtering step. For each corpora, we used `lingua-language-detector==2.0.2` (M. Stahl, 2023) in two ways. First, we ran `lingua` in a constrained bilingual mode, limiting the available languages to only English and Spanish. Second, we ran it again but this time in an unconstrained mode where it had to guess the language using all of its supported languages. We then did the final filtering by dropping sentence pairs if any of the following were true:

1. the source English sentence wasn't detected as English by both modes of `lingua`

2. the target Spanish sentence wasn't detected as Spanish by both modes of `lingua`

3. both sentences were identical

While we did not perform ablation experiments to compare these steps for filtering by language ID, we note that this process of filtering was introduced due to the observation of English output observed (by manual inspection) in our preliminary systems. Introducing this filtering resulted in output that was qualitatively observed to contain much less English text.

Finally, with a restricted subset of the initially chosen corpora, we sampled 20,000,000 sentence pairs from the corpora listed in Table 1 using the implementation of reservoir sampling in Larkin (2024) with 2024 as the seed.

We used *Statmt-newstest-2012-eng-spa* as our *validation* set, as suggested by `mtdata.recipes.wmt24-constrained.yml`.

## 2.3 Japanese→Chinese

We fetched all `jpn→zho` corpora available for WMT24's General Machine Translation.[3] We sampled 2000 sentence pairs for *validation* and 2000 sentence pairs for *test* (unused) from *Facebook-wikimatrix-1*, *Neulab-tedtalks_train-1*, *OPUS-wikimedia-v20210402*, *Statmt-news_commentary-18.1*. The remaining sentence pairs and all sentence pairs listed in the corpora of the second part of Table 2 were included in *train*.

## 3 Preprocessing and Postprocessing

There are two main types of preprocessing performed: subword segmentation (Section 3.1), which is perfomed on both the training data and the test data, and sentence splitting (Section 3.2) which is performed only on the WMT test data (as our models are trained primarily as sentence-level systems and should thus be applied to sentences rather than the full paragraphs and documents supplied at test time). We also describe the postprocessing that we performed (Section 3.3).

### 3.1 Subword Segmentation (Train and Test)

For details on our subword segmentation approach for `eng→deu`, see Lo et al. (2023). Our subword segmentation approach for `eng→spa` and `jpn→zho` is described here. To segment the corpora, a separate bilingual tokenizer (`SentencePieceUnigramTokenizer`) for each language pair was trained using HuggingFace's tokenizers (Moi and Patry, 2022), library version `0.14.1`. For each language pair, the vocabulary size was set to 32k tokens. Each tokenizer performs:

- control character and white space normalizations through HuggingFace's `Nmt`[4]

- NFKC normalization using HuggingFace's `NFKC`[5]

- and also applies a few normalizations done by Portage (Larkin et al., 2022). Some of these may overlap with the other normalization steps; see Appendix A.

---

[2] `mtdata get-recipe -i wmt24-eng-spa -o wmt24-eng-spa –compress –no-merge`

[3] `mtdata get-recipe -i wmt24-jpn-zho -o wmt24-jpn-zho –compress –no-merge`

[4] https://huggingface.co/docs/tokenizers/api/normalizers#tokenizers.normalizers.Nmt and https://github.com/huggingface/tokenizers/blob/main/tokenizers/src/normalizers/unicode.rs#L44

[5] https://huggingface.co/docs/tokenizers/api/normalizers#tokenizers.normalizers.NFKC

| corpus | original | step1 | step2 | ratio (%) |
|---|---|---|---|---|
| *EU-dcep-1* | 3,710,534 | 3,708,524 | 2,570,271 | 69.3 |
| *Facebook-wikimatrix-1* | 6,452,177 | 6,448,669 | 4,854,605 | 75.2 |
| *LinguaTools-wikititles-2014* | 16,598,519 | 16,598,519 | 1,144,423 | 6.9 |
| *OPUS-dgt-v2019* | 5,127,624 | 5,126,271 | 3,432,757 | 66.9 |
| *OPUS-dgt-v4* | 3,168,368 | 3,167,629 | 2,138,218 | 67.5 |
| *OPUS-elrc_emea-v1* | 777,371 | 777,262 | 596,733 | 76.8 |
| *OPUS-eubookshop-v2* | 5,215,515 | 5,212,657 | 4,651,096 | 89.2 |
| *OPUS-europarl-v8* | 2,009,073 | 2,008,951 | 1,928,793 | 96.0 |
| *OPUS-europat-v3* | 51,352,279 | 51,352,021 | 48,077,464 | 93.6 |
| *OPUS-multiun-v1* | 11,350,967 | 11,339,127 | 9,864,021 | 86.9 |
| *OPUS-unpc-v1.0* | 25,227,001 | 25,209,933 | 19,437,858 | 77.1 |
| *OPUS-wikimatrix-v1* | 3,377,911 | 3,377,355 | 2,708,923 | 80.2 |
| *OPUS-wikimedia-v20210402* | 1,275,296 | 1,272,410 | 910,544 | 71.4 |
| *OPUS-wikipedia-v1.0* | 1,811,428 | 1,808,866 | 1,196,239 | 66.0 |
| *OPUS-xlent-v1.1* | 9,251,728 | 9,251,728 | 830,623 | 9.0 |
| *Statmt-news_commentary-18.1* | 500,180 | 500,173 | 481,628 | 96.3 |
| *Tilde-eesc-2017* | 2,531,892 | 2,531,718 | 2,209,249 | 87.3 |
| *Tilde-rapid-2016* | 684,260 | 684,202 | 599,462 | 87.6 |
| **total** | 150,422,123 | 150,376,015 | 107,632,907 | 71.6 |

Table 1: Number of sentence pairs left after each filtering step for English→Spanish. The ratio column indicates the percentage of sentences pairs left from the original corpora after been filtered.

| corpus | # sentence pairs |
|---|---|
| *Facebook-wikimatrix-1* | 1,325,674 |
| *Neulab-tedtalks_train-1* | 5,159 |
| *OPUS-wikimedia-v20210402* | 23,132 |
| *Statmt-news_commentary-18.1* | 1,625 |
| *KECL-paracrawl-2-zho* | 83,892 |
| *LinguaTools-wikititles-2014* | 1,661,283 |
| *OPUS-bible_uedin-v1* | 124,260 |
| *OPUS-ccmatrix-v1* | 12,403,136 |
| *OPUS-gnome-v1* | 50 |
| *OPUS-kde4-v2* | 118,258 |
| *OPUS-multiccaligned-v1* | 4,280,695 |
| *OPUS-openoffice-v3* | 68,952 |
| *OPUS-opensubtitles-v2018* | 1,091,295 |
| *OPUS-php-v1* | 12,214 |
| *OPUS-qed-v2.0a* | 18,098 |
| *OPUS-tanzil-v1* | 12,472 |
| *OPUS-ted2020-v1* | 15,982 |
| *OPUS-ubuntu-v14.10* | 226 |
| *OPUS-ubuntu-v14.10* | 34 |
| *OPUS-xlent-v1.1* | 1,396,116 |
| **total** | 21,316,879 |

Table 2: Number of sentence pairs in each jpn→zho corpus. Corpora in the first part (*Facebook-wikimatrix-1* to *Statmt-news_commentary-18.1*) were used to sample *validation* and *test*. All corpora, except for the sentence pairs in *validation* and *test* were use for *train*.

The Neural Machine Translation (NMT) vocabulary is also augmented with 25 generic tokens (unused in these experiments); this yields a final vocabulary of 32029 tokens.

To train the eng→spa tokenizer, we used all training corpora provided except for *Facebook-wikimatrix-1*, *LinguaTools-wikititles-2014*, *OPUS-multiccaligned-v1*, *OPUS-wikimatrix-v1*, *OPUS-wikimedia-v20210402*, *OPUS-wikipedia-v1.0*, *OPUS-xlent-v1.1*, *ParaCrawl-paracrawl-9*, *Statmt-ccaligned-1*.

We used all 40 corpora available to train the jpn→zho subtokenizer model.

### 3.2 Sentence Splitting (Test-Only)

This year's General News Task test segments consist of paragraphs. To match our system's training configuration, we first split the paragraphs and documents into sentences before performing subword segmentation and translation for all language pairs. We do this for both the official test set and the test suites. We used utokenize.pl from Larkin et al. (2022) to sentence split the English segments of eng→deu and eng→spa. Since utokenize.pl doesn't support Japanese, we used ersatz (Wicks and Post, 2021) for jpn→zho. The speech documents in jpn→zho contain some punctuation but, in some cases, utterances appear to be separated only by spaces. For this domain only, we first split sentences using ersatz then followed this with a heuristic of splitting on spaces. We kept track of each sentence's segment and document ID to later

enable us to reconstruct the translations into their corresponding segment.

After sentence splitting is complete, we apply the subword segmenters described in Section 3.1 and perform translation at the level of the sentence. Since we perform sentence splitting of the source, the original source segments (paragraphs and documents) have to be reconstructed. We take this sentence-level output and concatenate the sentences belonging to a given input segment back together; for English→German and English→Spanish, we insert a space between sentences, while for Japanese→Chinese we concatenate without spaces.

### 3.3 Postprocessing (Test-Only)

In two cases, we performed additional postprocessing to handle issues specific to a language pair and/or a domain (as our training and validation data is more news-focused).

#### 3.3.1 English→Spanish

Our eng→spa translations contained some <unk> that clearly aligned to an emoji in the source (likely due to our training data not having strong coverage of social media domains). As a custom postprocessing step for eng→spa, we replaced the first <unk> with the first emoji in the source, the second <unk> with the second emoji and so on. For <unk> that did not have an emoji, they were considered spurious and were simply removed. Any extra emojis that couldn't be matched to a <unk> were simply added at the end of that translation. This was done because we noticed that our system would produce a single <unk> for multiple consecutive emojis.

#### 3.3.2 Japanese→Chinese

We noted some recurrent deficiencies in our Chinese translations. To fix those, we applied the following postprocessing steps:

- remove spaces between two Chinese characters

- remove spaces surrounding Chinese punctuation ：；，。？！

- when a Chinese character is repeated three or more times in a row, replace this with a single instance of that character

- fold repeating quotation marks onto a single quotation mark

## 4 MT System

We train all NMT models using Sockeye version 3.1.31 (Hieber et al., 2022), commit 13c63be5, with PyTorch 1.13.1 (Paszke et al., 2019). Training was performed on 4 Tesla V100-SXM2-32GB GPUs. Table 3 lists the parameter settings in our experiments that differ from the Sockeye defaults.

We train the models until convergence which is defined as no improvement in BLEU (Papineni et al., 2002; Post, 2018) for 32 checkpoints (when a model reaches this definition of convergence, training stops). The jpn→zho model trained for 390 checkpoints yielding its best checkpoint at update 358 and a BLEU score of 34.3 as reported on OCELoT over the WMT General Test Set. The eng→spa model trained for 832 checkpoints yielding its best checkpoint at update 800 and a BLEU score of 17.6 as reported on OCELoT over the WMT General Test Set. The eng→deu model had a score of 20.1 as reported on OCELoT over the WMT General Test Set.

## 5 Risks of Automatic System Selection for Human Evaluation

We submitted these systems with the intent of having them evaluated by human annotators, based on the understanding that "All submitted systems will be scored and ranked by human judgement."[6] Unfortunately, the task included a larger number of submissions than anticipated (Kocmi et al., 2024), resulting in the decision to remove some systems from human evaluation, as per the note in the evaluation section of the task page: "In the unlikely event of an unprecedented number of system submissions that we couldn't evaluate, we may decide to preselect the best performing systems for human evaluation with automatic metrics (such as COMET), we will primarily remove closed systems from the evaluation. However, we believe this won't be applied and all primary systems will be evaluated by humans." Among these, our submitted eng→deu and jpn→zho systems were removed from human evaluation, leaving only the eng→spa system to receive human evaluation by the General Task evaluation process.

However, all three of our submitted systems were evaluated using MQM (Multidimensional Quality Metrics; Lommel et al., 2013) by the Met-

---

[6]https://www2.statmt.org/wmt24/translation-task.html, most recently accessed Sept. 24, 2024.

| Name | Value | Default |
|---:|---|---|
| **amp** | *True* | *False* |
| **grading clipping type** | *abs* | *None* |
| **max sequence length** | *200:200* | *95:95* |
| **attention heads** | *16:16* | *8:8* |
| **shared vocabulary** | *True* | *False* |
| **transformer FFN** | *4096:4096* | *2048:2048* |
| **transformer model size** | *1024:1024* | *512:512* |
| **weight tying** | *True* | *False* |
| **batch size** | *8192* | *4096* |
| **batch type** | *max-word* | *word* |
| **cache last best params** | *2* | *0* |
| **cache metric** | *BLEU* | *perplexity* |
| **checkpoint interval** | 10 | 4000 |
| **decode and evaluate** | *-1 (entire validation)* | *500* |
| **initial learning rate** | 0.06325 | 0.0002 |
| **learning rate scheduler type** | *inv-sqrt-decay* | *plateau-reduce* |
| **learning rate warmup** | *4000* | *0* |
| **max num checkpoint not improved** | *32* | *None* |
| **max num epochs** | *1000* | *None* |
| **metrics** | *perplexity & accuracy* | *undefined* |
| **optimized metric** | *BLEU* | *perplexity* |
| **optimizer Betas** | *0.9, 0.98* | *0.9, 0.999* |
| **update interval** | *2* | *1* |

Table 3: Differences between Sockeye's default parameters and our eng→spa/jpn→zho configuration.



(a) English→Spanish    (b) Japanese→Chinese    (c) English→German

Figure 1: MQM scores on the News portion of the General MT test data, produced by the Metrics Task over a subset of the submitted WMT systems. Error bars represent bootstrap resampling, 1000 times, for $p < 0.05$. In all cases, our MSLC system appears at the far left of the plots, which are ordered by mean segment-level MQM score.

rics Shared Task. This offers a rare opportunity to examine the risks of selecting a subset of systems for human evaluation by using automatic metrics. In Fig. 1, we observe that the human rankings produced by MQM differ enough from the predicted rankings that they arguably demonstrate exactly the two types of errors one might be concerned about making: including a poorer quality system in human evaluation and, worse, failing to include a system with substantial confidence interval overlap with a system that was included for evaluation. In the first case, our eng→spa system, which was included for evaluation, appears substantially worse than other systems evaluated by MQM (Fig. 1a); however, we do note that IKUN-C, which could conceivably bridge the gap, was not included for evaluation by the Metrics Task, so it is possible that this does not represent an error. Unfortunately, without either human evaluation containing both, it is unlikely we can reach a definitive answer. In the second case, our jpn→zho system was excluded from human evaluation by the General MT task but IKUN-C was included for General MT task evaluation. In Fig. 1b, we can see that there is substantial confidence interval overlap between the MQM scores for the MSLC jpn→zho system and the IKUN-C system. We note that there are stronger ways to more definitively make this comparison (e.g., to do pairwise significance tests), but we primarily provide these examples for discussion and consideration. Finally, the eng→deu appears to represent the successful intended result of this approach to filtering sytems (Fig. 1c).

This highlights the risks of the mismatches between automatic evaluation and human evaluation; it may be better to perform some sort of smaller-scale initial human evaluation to separate systems rather than doing so based on automatic metrics.

## 6 Conclusion

We have built simple Transformer NMT models, primarily for the purpose of the MSLC dataset at the Metrics Task. We submit them to the WMT General Task to enable human evaluation, which will be useful to better understand how metrics perform and compare to human evaluation on a wider range of MT output quality. Of the three submitted systems, only one was included for human evaluation in the shared task.

## Limitations

As described, we submit extremely simple models, with minimal additional modifications. As our focus for MSLC is on news data, we expend only minimal effort on additional domains. We submit only three language pairs. We would not recommend the use of these MT systems outside of their intended uses for metric evaluation in MSLC.

## Ethics Statement

We build constrained MT systems, using the permitted training data from WMT24. Since our goal in this work is to build systems to be used to evaluate metrics across a wider range of translation quality, we expect that these systems may have a number of problems, including but not limited to: producing errors in translation, producing output in dialects (or languages) other than the desired ones, or otherwise produced biased output. We do not recommend their use for purposes other than the intended purpose of MSLC; their limitations for that purpose are discussed in more depth in the corresponding Metrics Task submission.

## Acknowledgements

We thank the WMT General Task organizers for their clarifications regarding data for the constrained task. We thank the WMT Metrics Task organizers for including our systems in their MQM human evaluation, to enable us to use those results to better understand the performance of automatic metrics across a range of MT quality.

## References

Thamme Gowda. 2024. A tool that locates, downloads, and extracts machine translation corpora.

Felix Hieber, Michael Denkowski, Tobias Domhan, Barbara Darques Barros, Celina Dong Ye, Xing Niu, Cuong Hoang, Ke Tran, Benjamin Hsu, Maria Nadejde, Surafel Lakew, Prashant Mathur, Anna Currey, and Marcello Federico. 2022. Sockeye 3: Fast Neural Machine Translation with PyTorch. *arXiv*, abs/2207.05851.

Huda Khayrallah and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83, Melbourne, Australia. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Preliminary wmt24 ranking of general mt systems and llms.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.

Samuel Larkin. 2024. A Python Implementation of Reservoir Sampling. https://github.com/SamuelLarkin/reservoir_sampling.

Samuel Larkin, Eric Joanis, Darlene Stewart, Michel Simard, George Foster, Nicola Ueffing, and Aaron Tikuisis. 2022. Portage Text Processing. https://github.com/nrc-cnrc/PortageTextProcessing.

Chi-kiu Lo, Samuel Larkin, and Rebecca Knowles. 2023. Metric score landscape challenge (MSLC23): Understanding metrics' performance on a wider landscape of translation quality. In *Proceedings of the Eighth Conference on Machine Translation*, pages 776–799, Singapore. Association for Computational Linguistics.

Chi-kiu Lo, Michel Simard, Darlene Stewart, Samuel Larkin, Cyril Goutte, and Patrick Littell. 2018. Accurate semantic textual similarity for cleaning noisy parallel corpora using semantic machine translation evaluation metric: The NRC supervised submissions to the parallel corpus filtering task. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 908–916, Belgium, Brussels. Association for Computational Linguistics.

Arle Richard Lommel, Aljoscha Burchardt, and Hans Uszkoreit. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.

Peter M. Stahl. 2023. The most accurate natural language detection library for Python, suitable for short text and mixed-language text.

Anthony Moi and Nicolas Patry. 2022. HuggingFace's Tokenizers. https://github.com/huggingface/tokenizers.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Rachel Wicks and Matt Post. 2021. A unified approach to sentence segmentation of punctuated text in many languages. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3995–4007, Online. Association for Computational Linguistics.

## A  Portage's Normalization

Table 4 describes the normalization steps done by Portage.

## B  Software Snapshots

For the three additional pieces of software, namely mtdata (Gowda, 2024), lingua (M. Stahl, 2023), and reservoir_sampling (Larkin, 2024), snapshots from September 24, 2024 are available on WaybackMachine (http://web.archive.org/), should their current URLs become unavailable.

- lingua is available at https://github.com/pemistahl/lingua-py; its snapshot is available at https://web.archive.org/web/20240924170712/https:

| Textual Description | Code |
|---|---|
| Convert various non-breaking hyphens to $-$ | $[\backslash u001E\backslash u00AD\backslash u2011] \rightarrow -$ |
| Strip out the MS Word discretional hyphen | $\backslash x1F$ |
| Replace special purpose spaces by regular spaces | $[\backslash u2060\backslash uFEFF\backslash u00A0\backslash u2007\backslash u202F\backslash u2028\backslash u2029] \rightarrow \sqcup$ |
| Replace remaining control characters by spaces | $[\backslash x01 - \backslash x09\backslash x0B\backslash x0C\backslash x0E - \backslash x1F\backslash x7F] \rightarrow \sqcup$ |
| convert DOS newlines to Linux ones | $\backslash x0d$ |
| Collapse multiple spaces to a single space | $\backslash s+ \rightarrow \sqcup$ |

Table 4: Portage normalizations

```
//github.com/pemistahl/lingua-py/
archive/refs/tags/v2.0.2.tar.gz
```

- `reservoir_sampling` is available at `https://github.com/SamuelLarkin/reservoir_sampling`; its snapshot is available at `https://web.archive.org/web/20240924170941/https://github.com/SamuelLarkin/reservoir_sampling/archive/refs/tags/0.1.tar.gz`

- `mtdata` is available at `https://github.com/thammegowda/mtdata`; its snapshot is available at `https://web.archive.org/web/20240924171242/https://github.com/thammegowda/mtdata/archive/refs/tags/v0.4.1.tar.gz`

# IOL Research Machine Translation Systems for WMT24 General Machine Translation Shared Task

**Wenbo Zhang, Qiaobo Deng, Zeyu Yan,** and **Hongbao Mao**
Transn IOL Research, Wuhan, China

## Abstract

This paper illustrates the submission system of the IOL Research team for the WMT24 General Machine Translation shared task. We submitted translations for all translation directions in the general machine translation task. According to the official track categorization, our system qualifies as an open system due to the utilization of open-source resources in developing our machine translation model. With the growing prevalence of large language models (LLMs) as a conventional approach for managing diverse NLP tasks, we have developed our machine translation system by leveraging the capabilities of LLMs. Overall, we first performed continued pretraining using the open-source LLMs with tens of billions of parameters to enhance the model's multilingual capabilities. Subsequently, we employed open-source Large Language Models, equipped with hundreds of billions of parameters, to generate synthetic data. This data was then blended with a modest quantity of additional open-source data for precise supervised fine-tuning. In the final stage, we also used ensemble learning to improve translation quality. Based on the official automated evaluation metrics, our system excelled by securing the top position in 8 out of the total 11 translation directions, spanning both open and constrained system categories.

## 1 Introduction

In the current year's WMT General Translation shared task, our team, IOL Research, took part in all 11 translation tasks, which involved translating text between various language pairs such as Czech to Ukrainian (cs->uk), Japanese to Chinese (ja->zh), English to Chinese (en->zh), English to Czech (en->cs), English to German (en->de), English to Hindi (en->hi), English to Icelandic (en->is), English to Japanese (en->ja), English to Russian (en->ru), English to Spanish (en->es), and English to Ukrainian (en->uk). One notable difference in this year's task compared to previous

years is that participants were required to translate paragraph-level texts, with one paragraph equating to one line. This change has significantly increased the length of the text to be translated. While traditional neural machine translation systems (Vaswani et al., 2017) based on encoder-decoder structures may struggle with processing long texts due to the lack of enough document parallel data. However, the large language models (LLMs) do not necessitate a large amount of lengthy text data for fine-tuning, making them more effective in handling long texts. As a result, we meticulously trained an LLM with 20 billion parameters to successfully address all translation tasks in the competition.

Our main strategy is to explore using LLMs to build machine translation systems. This includes fine-tuning the translation task on foundational LLMs and leveraging advanced open-source instruction-tuned LLMs to generate high-quality translation data for further enhancement. Before supervised fine-tuning, we also performed continued pretraining, which has been proven to be very beneficial for translation tasks (Xu et al., 2023), because many open-source LLMs such as LLaMA (Touvron et al., 2023) are usually pretrained on English monolingual data, lacking the necessary knowledge of other languages required for translation tasks. Moreover, we experimented with ensemble learning, a technique known to be effective for neural machine translation models. We discovered that it provided some degree of assistance for machine translation tasks based on LLMs. In the end, our billion-parameter machine translation system achieved comparable performance to hundred billion parameter LLMs in high-resource languages and even outperformed them in certain low-resource languages.

The subsequent paper is designed as follows. We introduce the data source and processing strategy in Section 2; Section 3 describes the details of our training procedure; Section 4 presents the experi-

mental settings and results.

## 2 System Overview

### 2.1 Model Architecture

We selected the Qwen1.5 model (Bai et al., 2023) as our foundational model because of its outstanding performance and considerable multilingual capabilities. Specifically, we utilized the Qwen1.5-14B[1] as our starting point, which has 40 layers and 14 billion parameters. To enhance the model's capacity within our hardware constraints, we concatenated the first 32 layers with the last 32 layers, resulting in duplication of the middle 24 layers, following the approach used in SOLAR (Kim et al., 2023). This fusion led to a scaled-up model with 64 layers and 21 billion parameters. Since this approach alters the structure of the pretrained model, continual pretraining becomes a necessary step to recover its performance.

### 2.2 Continual Pretraining

Continual pretraining is an effective method to enhance the knowledge embedded within LLMs. This method has been extensively utilized to adapt LLMs from English to various other languages, as well as to augment the domain-specific knowledge inherent in these models. In the context of using LLMs for translation tasks, it has been substantiated that the continuous pretraining of LLMs with multilingual monolingual data, encompassing languages involved in all the translation directions, is crucial (Xu et al., 2023). This year's WMT24 general machine translation task includes 11 translation directions, involving 10 distinct languages. Therefore, our continued pretraining is carried out on monolingual data in these 10 languages.

We sampled the required multilingual monolingual data from the mC4 (Raffel et al., 2019) and OSCAR (Jansen et al., 2022) datasets, then proceeded to refine the chosen data. For refinement processes, we employed fastText (Joulin et al., 2017) for language identification, the minLSH algorithm for document deduplication, and KenLM (Heafield, 2011) tool for filtering the documents with high perplexity. Many studies (Lin et al., 2020; Yang et al., 2021) have shown that integrating bilingual data with monolingual data in the pretraining stage can help the model achieve better cross-lingual proficiency. Therefore, we also incorporated a portion of the CC-Aligned

parallel data (El-Kishky et al., 2019) into our continuous pretraining stage. This data includes language pairs such as English-Czech, English-Ukrainian, English-Japanese, English-Chinese, English-German, English-Hindi, English-Icelandic, English-Russian, and English-Spanish. Specifically, we randomly swapped the order of the two articles in the bilingual document, and then merged them into a new document as the pretraining document. The distribution of the number of documents in all languages in the pretraining dataset is shown in Table 1.

| Language | Rate(%) |
| --- | --- |
| en | 21.99 |
| ja | 15.02 |
| de | 12.48 |
| cs | 11.60 |
| es | 10.35 |
| zh | 9.32 |
| uk | 7.98 |
| ru | 7.2 |
| hi | 3.53 |
| is | 0.47 |

Table 1: The distribution of the number of documents in all languages in the pretraining dataset.

### 2.3 Supervised Fine-tuning

Through supervised Fine-tuning, we can unlock the capabilities of LLMs using only a minimal amount of aligned data. Many fine-tuning LLMs experiences (Zhou et al., 2024; Xia et al., 2024) have demonstrated that the quality and diversity of fine-tuning data are far more important than its quantity. In the context of translation tasks, high-quality parallel data is the ideal fine-tuning data for LLMs. However, obtaining such high-quality parallel data is challenging. Usually, we need to invest significant effort and undergo numerous steps to clean publicly available parallel data, aiming to achieve high-quality data. However, this process does not always guarantee the quality of filtered data due to its inherent complexity. On the other hand, start-of-the-art machine translation systems have shown competitive performance comparable to human translators. Consequently, we opted to employ LLMs to generate parallel data as the supervised fine-tuning data.

We used the c4ai-command-r-plus[2] and

---

[1]https://huggingface.co/Qwen/Qwen1.5-14B

[2]https://huggingface.co/CohereForAI/c4ai-command-r-

Qwen1.5-110B-Chat[3], these two instruction fine-tuned models, to generate synthetic parallel data for all languages, with the exception of Icelandic. Specifically, when the task requires generating Chinese content, our go-to model is the Qwen1.5-110B-Chat. However, for English content generation, we make a random selection between the Qwen1.5-110B-Chat and c4ai-command-r-plus models. For all other scenarios, we consistently utilize the c4ai-command-r-plus model. The selection of models in different languages is based on our evaluation of these two models in translation tasks. Please refer to Table 3 for specific comparison. Considering the lack of proficiency of both c4ai-command-r-plus and Qwen1.5-110B-Chat in generating Icelandic content, we adopted an alternative strategy. We leveraged our supervised fine-tuning model, which has been fine-tuned on synthetic data of all other languages, to produce the synthetic data for Icelandic. Therefore, our model only utilized Icelandic monolingual data for pre-training, and the Icelandic bilingual synthesis data was generated by unsupervised method.

We have tried two synthetic data generation methods commonly used in traditional neural machine translation systems, forward translation (Kim and Rush, 2016) and back translation (Sennrich et al., 2016). Forward translation refers to using the established translation model to translate real source language sentences into target language sentences, and then combining the translated target language sentences with the real source language sentences to form synthetic parallel sentence pairs. Back translation refers to translating real target language sentences back into the source language using another established reverse translation model, and then combining the real target language sentences with the translated source language sentences to form synthetic parallel sentence pairs. In the process of generating back translation data based on real target language data, we found that the real target language data has many problems such as incoherence, fluency deficits, and even grammatical errors. To address these problems, we utilized automatic post-editing technology. This approach involves taking the translated source language sentences and the real target language sentences as inputs, and subsequently producing su-

perior quality target language sentences. These improved sentences are then used to replace the real target language sentences in the back translation synthetic data. Lastly, we also utilized LLMs to filter all the generated synthetic data, including both forward and back translation data, to ensure higher quality fine-tuning data. All the prompts we use to generate synthetic data are shown in the table 2. For each language pair, after filtering, we retained around 100,000 FT and BT sentence pairs respectively.

In addition to synthetic data, we also incorporated document parallel data from News Commentary v18.1[4], which assists the model in translating long text, and instruction fine-tuning data TowerBlocks-v0.2 (Alves et al., 2024) to help the model follow more diverse instructions. The News Commentary v18.1 data we used includes sections ja-zh, en-zh, en-de, en-hi, en-ja, en-ru, en-es, en-cs, cs-ru, cs-de, cs-es, cs-hi, cs-ja, cs-zh, and ja-ru. We also excluded the data from TowerBlocks-v0.2 that includes FLoRes (Goyal et al., 2021), and the NTREX-128 (Federmann et al., 2022) sections, as we used these two datasets as our test sets to verify the performance of the model.

## 2.4 Ensemble Learning

The ensemble learning approach has demonstrated significant efficacy in a wide range of machine learning tasks. In machine translation tasks, ensemble learning completes the generation of the entire translation by using multiple different machine translation models to autoregressively vote for the probability distribution of the next word. However, for LLMs, this method implies a huge memory occupancy and computational resource consumption, so we use transductive ensemble learning (Wang et al., 2020) to replace this way of generating with multiple models simultaneously. Transductive ensemble learning first utilizes multiple different translation models to generate translations for the same test set separately, then aggregates all translations as fine-tuning data. The final translation is generated by one translation model after fine-tuning on this data. Ensemble learning conventionally entails training diverse models via different random initializations. However, this approach proves inefficient in our context, as we are mandated to employ the identical pre-trained model for supervised fine-tuning. Therefore, we used dif-

---

plus

| Task | Prompt |
|---|---|
| Forward and back translation | Translate the following text from SRC_LANG to TGT_LANG. SRC_CONTENT |
| Automatic post-editing | Given a source SRC_LANG sentence and its TGT_LANG translation, please modify and correct the TGT_LANG translation to get a more accurate and fluent TGT_LANG translation. Source (SRC_LANG): SRC_CONTENT Translation (TGT_LANG): TGT_CONTENT Corrected translation (TGT_LANG): |
| Synthetic data filtering | Source (SRC_LANG): SRC_CONTENT Translation (TGT_LANG): TGT_CONTENT Please check if the above translation is an accurate and fluent translation of its source text? Please only answer "yes" or "no" |

Table 2: All the prompts we use to generate synthetic data.

ferent fine-tuning data to train multiple models for ensemble learning. Different fine-tuning data is obtained by randomly sampling synthesized data from different parts.

## 3 Experiments

### 3.1 Experiment Settings

For continual pretraining phase, we trained the scaled-up model with 21 billion parameters on 8 NVIDIA H800 GPUs. For the optimization process, we employed the Adam optimizer (Kingma and Ba, 2014), with $\beta1 = 0.9, \beta2 = 0.99$. We adopted a learning rate scheduling strategy that remained constant after warmup phase, setting the number of warmup steps to 200, the maximum learning rate at 0.00001 and weight decay to 0.1. The batch size was set to 3.14 million tokens, the length of each sequence was set to 4096, and a total of 56 billion tokens have been trained.

For supervised fine-tuning phase, we fine-tuned the continual pretrained model on 16 NVIDIA H800 GPUs. We leveraged the Adam optimizer for the optimization process, setting $\beta1 = 0.9, \beta2 = 0.99$. We employed a cosine learning rate scheduling strategy, with a warmup ratio of 0.01, a peak learning rate at 0.000007, and a weight decay of 0.1. Configuring the batch size to 480 sentences, we trained the model for a single epoch encompassing approximately 1.5 million sentences.

When conducting transductive ensemble learning, we increased the batch size to 800 sentences, adopted a fixed learning rate, and reduced the learning rate to 0.000001. Similarly, we only fine-tune for one epoch on the ensemble data.

### 3.2 Results

The FLoRes (Goyal et al., 2021) and NTREX-128 (Federmann et al., 2022) test sets were utilized as our evaluation benchmarks. The performance of the machine translation system was assessed using SacreBLEUpost-2018-call and COMET (Rei et al., 2022)[5] metrics. We uesed vLLM (Kwon et al., 2023) to infer all LLMs. We chose c4ai-command-r-plus and Qwen1.5-110B-Chat as our baselines for comparison, and all results were obtained through zero-shot evaluation.

Test results on the FLoRes test set for all translation directions are shown in Table 3. We used greedy decoding and beam search with beam size = 5 to generate translations for our model, and provided the ensemble effect on this test set. It is clear that, just like traditional neural machine translation models, beam search performs better than greedy decoding in terms of BLUE and COMET scores across all translation directions. Ensemble learning has a steady improvement on BLEU scores, but the overall change in COMET scores is not significant. Compared with the two baseline systems CMD-R-P and Qwen1.5-L, our model achieved equivalent or better performance in the seven directions of cs→uk, en→zh, en→de, en→hi, en→is, en→uk, and en→cs. The performance outcomes presented in Table 4 are based on evaluations conducted using the NTREX-128 test set. These results mirror those observed in the FLoRes test set, indicating a consistent performance trend across both datasets.

---

[5]https://huggingface.co/Unbabel/wmt22-comet-da

| | | CMD-R-P | Qwen1.5-L | our model greedy decoding | our model beam search | our model ensemble learning |
|---|---|---|---|---|---|---|
| cs→uk | BLEU | 24.1 | 20.5 | 23.9 | 24.4 | 24.6 |
| | COMET | 90.47 | 87.96 | 90.18 | 90.41 | 90.47 |
| ja→zh | BLEU | 31.6 | 34.1 | 34.8 | 35.3 | 35.0 |
| | COMET | 87.91 | 88.10 | 87.99 | 88.11 | 87.98 |
| en→zh | BLEU | 39.9 | 44.0 | 46.9 | 47.5 | 47.6 |
| | COMET | 88.71 | 89.08 | 89.22 | 89.28 | 89.26 |
| en→de | BLEU | 41.1 | 33.9 | 40.5 | 41.1 | 41.6 |
| | COMET | 88.84 | 87.37 | 88.60 | 88.73 | 88.84 |
| en→hi | BLEU | 27.3 | 19.9 | 27.6 | 28.5 | 28.7 |
| | COMET | 80.47 | 75.01 | 79.99 | 80.75 | 80.67 |
| en→is | BLEU | 12.1 | 9.8 | 19.8 | 20.5 | 20.7 |
| | COMET | 71.41 | 63.82 | 82.77 | 83.66 | 84.02 |
| en→ja | BLEU | 49.8 | 42.2 | 49.4 | 50.1 | 50.4 |
| | COMET | 91.70 | 89.88 | 91.50 | 91.59 | 91.61 |
| en→ru | BLEU | 32.4 | 27.6 | 31.3 | 31.9 | 32.4 |
| | COMET | 90.70 | 87.98 | 90.09 | 90.32 | 90.28 |
| en→es | BLEU | 30.4 | 27.1 | 29.4 | 29.4 | 29.5 |
| | COMET | 87.29 | 86.64 | 87.01 | 87.06 | 86.98 |
| en→uk | BLEU | 30.4 | 24.6 | 31.2 | 32.0 | 32.2 |
| | COMET | 90.88 | 88.19 | 90.56 | 90.83 | 90.92 |
| en→cs | BLEU | 32.7 | 26.6 | 32.8 | 34.3 | 34.4 |
| | COMET | 92.09 | 90.04 | 91.78 | 92.15 | 92.13 |

Table 3: Test results on the FLoRes test set for all translation directions. CMD-R-P represents c4ai-command-r-plus, and Qwen1.5-L represents Qwen1.5-110B-Chat.

|  |  | CMD-R-P | Qwen1.5-L | our model greedy decoding | our model beam search |
|---|---|---|---|---|---|
| cs→uk | BLEU | 20.9 | 16.8 | 20.4 | 20.8 |
|  | COMET | 88.26 | 84.57 | 87.80 | 88.00 |
| ja→zh | BLEU | 25.6 | 28.7 | 28.7 | 29.0 |
|  | COMET | 84.42 | 84.84 | 84.83 | 84.85 |
| en→zh | BLEU | 31.7 | 36.6 | 39.0 | 39.5 |
|  | COMET | 85.60 | 86.41 | 86.76 | 86.83 |
| en→de | BLEU | 33.9 | 27.1 | 33.2 | 33.9 |
|  | COMET | 87.05 | 84.64 | 86.63 | 86.78 |
| en→hi | BLEU | 22.2 | 16.8 | 23.3 | 24.1 |
|  | COMET | 78.07 | 72.23 | 77.96 | 78.58 |
| en→is | BLEU | 14.8 | 11.2 | 23.4 | 24.1 |
|  | COMET | 70.14 | 62.32 | 82.75 | 83.67 |
| en→ja | BLEU | 41.3 | 35.0 | 41.3 | 42.4 |
|  | COMET | 89.51 | 87.40 | 89.37 | 89.45 |
| en→ru | BLEU | 29.9 | 23.8 | 30.4 | 31.5 |
|  | COMET | 88.13 | 84.47 | 87.47 | 87.88 |
| en→es | BLEU | 42.5 | 38.2 | 42.4 | 42.7 |
|  | COMET | 87.06 | 85.82 | 86.69 | 86.85 |
| en→uk | BLEU | 26.2 | 20.5 | 26.3 | 26.9 |
|  | COMET | 88.86 | 85.42 | 88.51 | 88.73 |
| en→cs | BLEU | 29.0 | 22.6 | 29.1 | 30.4 |
|  | COMET | 89.90 | 87.06 | 89.73 | 90.19 |

Table 4: Test results on the NTREX-128 test set for all translation directions. CMD-R-P represents c4ai-command-r-plus, and Qwen1.5-L represents Qwen1.5-110B-Chat.

|  | cs→uk | ja→zh | en→zh | en→de | en→hi | en→is | en→ja | en→ru | en→es | en→uk | en→cs |
|---|---|---|---|---|---|---|---|---|---|---|---|
| FT | 90.32 | 87.79 | 89.21 | 88.55 | 80.48 | 83.37 | 91.55 | 90.11 | 86.94 | 90.64 | 91.93 |
| BT | 90.43 | 88.18 | 89.24 | 88.77 | 81.45 | 84.12 | 91.64 | 90.43 | 87.20 | 90.99 | 92.38 |
| MIX | 90.32 | 88.12 | 89.26 | 88.63 | 80.76 | 84.08 | 91.50 | 90.19 | 87.06 | 90.63 | 91.93 |

Table 5: COMET scores of models fine-tuned on different data on the Flores test set. FT is fine-tuned on forward translation data. BT is fine-tuned on back translation data. MIX is fine-tuned on both forward and back translation data.

## 3.3 Forward Translation vs Back Translation

To determine the effectiveness of forward translation versus back translation, we separately fine-tuned the continual pretrained model using forward translation data, back translation data, and a combination of both. For each approach, we randomly chose 80,000 data samples per language translation direction. For the combined dataset, we selected 40,000 samples from both the forward translation and back translation pools. The results are presented in Table 5, all of which were generated using beam search. We can see that the back translation yields better performance, whereas mixed data does not result in significant improvement. Due to time constraints, we used mixed data in the WMT24 competition, this conclusion will guide us to further improve our model in the future.

## 4 Conclusion

In this paper, we present IOL Research's contributions to the WMT24 General Translation shared task, covering all translation aspects. Our approach utilizes LLMs to develop an effective translation system. Experimental results demonstrate that our model, which contains 21 billion parameters, achieves competitive results comparable to models with 100 billion parameters. According to the official automatic evaluation metrics (Kocmi et al., 2024), our system achieved 8 first places in 11 translation directions spanning both open and constrained system categories, including Czech to Ukrainian, English to German, English to Spanish, English to Hindi, English to Russian, English to Ukrainian, English to Chinese, and Japanese to Chinese.

## References

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu,

Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609.*

Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2019. Ccaligned: A massive collection of cross-lingual web-document pairs. *arXiv preprint arXiv:1911.06154.*

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. Ntrex-128–news test references for mt evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2021. The flores-101 evaluation benchmark for low-resource and multilingual machine translation.

Kenneth Heafield. 2011. Kenlm: Faster and smaller language model queries. In *Proceedings of the sixth workshop on statistical machine translation*, pages 187–197.

Tim Jansen, Yangling Tong, Victoria Zevallos, and Pedro Ortiz Suarez. 2022. Perplexed by Quality: A Perplexity-based Method for Adult and Harmful Content Detection in Multilingual Heterogeneous Web Data. *arXiv e-prints*, page arXiv:2212.10440.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431. Association for Computational Linguistics.

Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, et al. 2023. Solar 10.7 b: Scaling large language models with simple yet effective depth up-scaling. *arXiv preprint arXiv:2312.15166.*

Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. *arXiv preprint arXiv:1606.07947.*

DiederikP. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv: Learning,arXiv: Learning.*

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Findings of the WMT24 general machine translation shared task:

the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.

Zehui Lin, Xiao Pan, Mingxuan Wang, Xipeng Qiu, Jiangtao Feng, Hao Zhou, and Lei Li. 2020. Pre-training multilingual neural machine translation by leveraging alignment information. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2649–2663, Online. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*.

Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timoth'ee Lacroix, Baptiste Rozi'ere, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Yiren Wang, Lijun Wu, Yingce Xia, Tao Qin, ChengXiang Zhai, and Tie-Yan Liu. 2020. Transductive ensemble learning for neural machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6291–6298.

Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. Less: Selecting influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*.

Ziqing Yang, Wentao Ma, Yiming Cui, Jiani Ye, Wanxiang Che, and Shijin Wang. 2021. Bilingual alignment pre-training for zero-shot cross-lingual transfer. *arXiv preprint arXiv:2106.01732*.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

# Choose the Final Translation from NMT and LLM hypotheses Using MBR Decoding: HW-TSC's Submission to the WMT24 General MT Shared Task

**Zhanglin Wu, Daimeng Wei, Zongyao Li, Hengchao Shang, Jiaxin Guo,**
**Shaojun Li, Zhiqiang Rao, Yuanchang Luo, Ning Xie, Hao Yang**
Huawei Translation Service Center, Beijing, China
{wuzhanglin2,weidaimeng,lizongyao,shanghengchao,guojiaxin1,
lishaojun18,raozhiqiang,nicolas.xie,yanghao30}@huawei.com

## Abstract

This paper presents the submission of Huawei Translate Services Center (HW-TSC) to the WMT24 general machine translation (MT) shared task, where we participate in the English to Chinese (en→zh) language pair. Similar to previous years' work, we use training strategies such as regularized dropout, bidirectional training, data diversification, forward translation, back translation, alternated training, curriculum learning, and transductive ensemble learning to train the neural machine translation (NMT) model based on the deep Transformer-big architecture. The difference is that we also use continue pre-training, supervised fine-tuning, and contrastive preference optimization to train the large language model (LLM) based MT model. By using Minimum Bayesian risk (MBR) decoding to select the final translation from multiple hypotheses for NMT and LLM-based MT models, our submission receives competitive results in the final evaluation.

## 1 Introduction

Machine translation (MT) (Brown et al., 1990) predominantly utilizes transformer encoder-decoder architectures (Vaswani et al., 2017), which is evident in prominent models such as NLLB-200 (Costa-jussà et al., 2022), M2M100 (Fan et al., 2021), and MT5 (Xue et al., 2021). Significant research effort has been devoted to task-specific neural machine translation (NMT) models (Wei et al., 2022; Wu et al., 2023b) trained in a fully supervised manner with large volumes of parallel data. Their performance has been enhanced through techniques such as regularized dropout (Wu et al., 2021), bidirectional training (Ding et al., 2021), data diversification (Nguyen et al., 2020), forward translation (Abdulmumin, 2021), back translation (Sennrich et al., 2016), alternated training (Jiao et al., 2021), curriculum learning (Zhang et al., 2019), and transductive ensemble learning (Wang et al., 2020b).

The emergence of decoder-only large language models (LLMs) such as the GPT series (Wu et al., 2023a; Achiam et al., 2023), Mistral (Jiang et al., 2023), and LLaMA (Touvron et al., 2023a,b) shows remarkable efficacy in various NLP tasks, providing a fresh perspective on the MT task. Recent studies (Hendy et al., 2023; Jiao et al., 2023) indicate that larger LLMs such as GPT-3.5 (175B) and GPT-4 exhibit strong translation abilities. However, the performance of smaller-sized LLMs (7B or 13B) still falls short when compared to conventional NMT models (Zhu et al., 2024). Therefore, there are studies (Yang et al., 2023; Zeng et al., 2024) intend to enhance the translation performance for these smaller LLMs, but their improvements are relatively modest, primarily due to the predominant pre-training of LLMs on English-centric datasets, resulting in limited linguistic diversity. Addressing this limitation, Xu et al. (Xu et al., 2023) initially continue pre-training (CPT) LLaMA-2 (Touvron et al., 2023b) with extensive non-English monolingual data to enhance their multilingual abilities, and then perform supervised fine-tuning (SFT) with high-quality parallel data to instruct the model to generate translations. Nonetheless, the performance still lags behind leading translation models such as GPT-4 and WMT competition winners. Subsequently, Xu et al. (Xu et al., 2024) bridged this gap by further fine-tuning the LLM-based MT model using contrast preference optimization (CPO).

Ensembling (Zhou et al., 2002) has a long history in machine learning, being well known for leveraging multiple complementary systems to improve performance on a given task and provide good/robust generalization. Minimum Bayesian risk (MBR) (Finkelstein and Freitag, 2023; Farinhas et al., 2023) decoding has successfully improved translation quality using task-specific NMT models, and subsequently it has also been shown to be suitable for LLM-based MT models.

Figure 1: CPT, SFT and CPO data templates used for LLM-based MT training.

For the WMT24 general MT shared task, we participate in the en→zh language pair. Similar to previous years' work (Wei et al., 2021, 2022; Wu et al., 2023b), we use training strategies such as regularized dropout (Wu et al., 2021), bidirectional training (Ding et al., 2021), data diversification (Nguyen et al., 2020), forward translation (Abdulmumin, 2021), back translation (Sennrich et al., 2016), alternated training (Jiao et al., 2021), curriculum learning (Zhang et al., 2019), and transductive ensemble learning (Wang et al., 2020b) to train NMT models based on the deep transformer-big architecture. In addition, we use CPT, SFT and CPO methods to train LLM-based MT models. Finally, we use MBR decoding to select the final translation from multiple hypotheses of NMT and LLM-based MT models.

## 2 Data

### 2.1 Data Source

We obtain bilingual and monolingual data from ParaCrawl v9, News Commentary v18.1, Wiki Titles v3, UN Parallel Corpus V1.0, CCMT Corpus, WikiMatrix, News Crawl and Common Crawl data sources. The amount of data we used for training NMT and LLM-based MT models is shown in Table 1. It should be noted that in order to obtain better translation performance in the general domain, we mix the monolingual data from Common Crawl and News Crawl.

### 2.2 NMT Data Pre-processing

Our data pre-processing methods for NMT include:

| language pairs | bitext data | monolingual data |
|---|---|---|
| en→zh | 25M | en: 50M, zh: 50M |

Table 1: Bilingual and monolingual used for training NMT and LLM-based MT models.

- Remove duplicate sentences or sentence pairs.

- Convert full-width symbols to half-width.

- Use fasttext[1] (Joulin et al., 2016) to filter other language sentences.

- Use jieba[2] to pre-segment Chinese sentences.

- Use mosesdecoder[3] (Koehn et al., 2007) to normalize English punctuation.

- Filter out sentences with more than 150 words.

- Use fast-align (Dyer et al., 2013) to filter sentence pairs with poor alignment.

- Sentencepiece[4] (SPM) (Kudo and Richardson, 2018) is used to perform subword segmentation, and the vocabulary size is set to 32K.

Since there may be some semantically dissimilar sentence pairs in bilingual data, we use LaBSE[5]

---

[1] https://github.com/facebookresearch/fastText
[2] https://github.com/fxsjy/jieba
[3] https://github.com/moses-smt/mosesdecoder
[4] https://github.com/google/sentencepiece
[5] https://huggingface.co/sentence-transformers/LaBSE

([Feng et al., 2022](#)) to calculate the semantic similarity of each bilingual sentence pair, and exclude bilingual sentence pairs with a similarity score lower than 0.7 from our training corpus.

## 2.3 LLM-based MT Data Pre-processing

The training of the LLM-based MT model requires three stages: CPT, SFT and CPO. As shown in Figure 1, the training data templates of the LLM-based MT model in these three stages are different.

In the CPT stage, considering that most LLMs are trained on English-dominated data, we using Chinese and English monolinguals for CPT to improve LLM's proficiency in Chinese. To preserve the long-context modeling capability of LLM, we concatenate multiple sentences into a long text with no more than 4096 words, and preferentially concatenate sentences from the same document.

In the SFT stage, drawing inspiration from the recognized significance of data quality in other applications ([Zhou et al., 2024](#); [Maillard et al., 2023](#)), we fine-tune the model with high-quality parallel data. In order to obtain high-quality parallel data, we use cometkiwi model [6] ([Rei et al., 2022](#)) to calculate the score of bilingual data on the en→zh language pair, and then retain bilingual data with a cometkiwi score greater than 0.8.

In the CPO stage, to learn an objective that fosters superior translations and rejects inferior ones, access to labeled preference data is essential, yet such data is scarce in machine translation. The following describes our process of constructing the triplet preference data required for CPO training. First, we randomly sample 50,000 data from high-quality bilingual data. Then, we use the NMT model to obtain N-best (N=10) hypotheses based on beam search decoding, and then use the comet-da model[7] ([Rei et al., 2020](#)) to calculate the score of each hypothesis, select the hypothesis with the highest score as the preferred translation, and select the hypothesis with the lowest score as the dis-preferred translation.

## 3 NMT System

### 3.1 System Overview

Transformer is the state-of-the-art model structure in recent NMT evaluations. There are two

Figure 2: The overall training flow of NMT system.

parts of research to improve this kind: the first part uses wide networks (eg: Transformer-Big ([Vaswani et al., 2017](#))), and the other part uses deeper language representations (eg: Deep Transformer ([Wang et al., 2019](#))). For the WMT24 general MT shared task, we combine these two improvements, adopting the Deep Transformer-Big ([Wei et al., 2022](#); [Wu et al., 2023b](#)) model structure to train the NMT system. Deep Transformer-Big uses pre-layer normalization, features 25-layer encoder, 6-layer decoder, 16-heads self-attention, 1024-dimensional word embedding and 4096-dimensional FFN embedding.

Fig. 2 shows the overall training flow of NMT system. We use training strategies such as regularized dropout (R-Drop) ([Wu et al., 2021](#)), bidirectional training (BiT) ([Ding et al., 2021](#)), data diversification (DD) ([Nguyen et al., 2020](#)), forward translation FT) ([Abdulmumin, 2021](#)), back translation (BT) ([Sennrich et al., 2016](#)), alternated training (AT) ([Jiao et al., 2021](#)), curriculum learning (CL) ([Zhang et al., 2019](#)), and transductive ensemble learning (TEL) ([Wang et al., 2020b](#)) for training.

### 3.2 Regularized Dropout

Regularized Dropout (R-Drop)[8] ([Wu et al., 2021](#)) is a simple yet more effective alternative to regularize the training inconsistency induced by dropout ([Srivastava et al., 2014](#)). Concretely, in each mini-batch training, each data sample goes through the forward pass twice, and each pass is processed by a different sub model by randomly dropping out some hidden units. R-Drop forces the two distributions for the same data sample outputted by the two sub models to be consistent with each other, through minimizing the bidirectional Kullback-Leibler (KL) divergence ([Van Erven and Harremos, 2014](#)) between the two distributions. That is, R-Drop regularizes the outputs of two sub models ran-

domly sampled from dropout for each data sample in training. In this way, the inconsistency between the training and inference stage can be alleviated.

### 3.3 Bidirectional Training

Many studies have shown that pre-training can transfer the knowledge and data distribution, hence improving the model generalization. Bidirectional training (BiT) (Ding et al., 2021) is a simple and effective pre-training method for NMT. Bidirectional training is divided into two stages: (1) bidirectionally updates model parameters, and (2) tune the model. To achieve bidirectional updating, we only need to reconstruct the training samples from "src→tgt" to "src→tgt & tgt→src" without any complicated model modifications. Notably, BiT does not require additional parameters or training steps and only uses parallel data.

### 3.4 Data Diversification

Data Diversification (DD) (Nguyen et al., 2020) is a data augmentation method to boost NMT performance. It diversifies the training data by using the predictions of multiple forward and backward models and then merging them with the original dataset which the final NMT model is trained on. DD is applicable to all NMT models. It does not require extra monolingual data, nor does it add more parameters. To conserve training resources, we only use one forward model and one backward model to diversify the training data.

### 3.5 Forward Translation

Forward translation (FT) (Abdulmumin, 2021), also known as self-training, is one of the most commonly used data augmentation methods. FT has proven effective for improving NMT performance by augmenting model training with synthetic parallel data. Generally, FT is performed in three steps: (1) randomly sample a subset from the large-scale source monolingual data; (2) use a "teacher" NMT model to translate the subset data into the target language to construct the synthetic parallel data; (3) combine the synthetic and authentic parallel data to train a "student" NMT model.

### 3.6 Back Translation

An effective method to improve NMT with target monolingual data is to augment the parallel training data with back translation (BT) (Sennrich et al., 2016; Wei et al., 2023). There are many works expand the understanding of BT and investigates a number of methods to generate synthetic source sentences. Edunov et al. (2018) find that back translations obtained via sampling or noised beam outputs are more effective than back translations generated by beam or greedy search in most scenarios. Caswell et al. (2019) show that the main role of such noised beam outputs is not to diversify the source side, but simply to tell the model that the given source is synthetic. Therefore, they propose a simpler alternative strategy: Tagged BT. This method uses an extra token to mark back translated source sentences, which generally outperforms noised BT (Edunov et al., 2018). For better joint use with FT, we use sampling back translation (ST) (Edunov et al., 2018).

### 3.7 Alternated Training

While synthetic bilingual data have demonstrated their effectiveness in NMT, adding more synthetic data often deteriorates translation performance since the synthetic data inevitably contains noise and erroneous translations. Alternated training (AT) (Jiao et al., 2021) introduce authentic data as guidance to prevent the training of NMT models from being disturbed by noisy synthetic data. AT describes the synthetic and authentic data as two types of different approximations for the distribution of infinite authentic data, and its basic idea is to alternate synthetic and authentic data iteratively during training until the model converges.

### 3.8 Curriculum Learning

A practical curriculum learning (CL) (Zhang et al., 2019) method should address two main questions: how to rank the training examples, and how to modify the sampling procedure based on this ranking. For ranking, we choose to estimate the difficulty of training samples according to their domain feature (Wang et al., 2020a). The calculation formula of domain feature is as follows, where $\theta_{in}$ represents an in-domain NMT model, and $\theta_{out}$ represents a out-of-domain NMT model. One thing to note is that we treat domains including news, user-generated (social), conversational, and e-commerce domains as in-domain, and others as out-of-domain. Specifically, we use the WMT22 test set to fine-tune a baseline model, and then use the baseline model and the fine-tuned model as the out-of-domain model and the in-domain model respectively.

$$q(x, y) = \frac{\log P(y|x; \theta_{in}) - \log P(y|x; \theta_{out})}{|y|}$$

(1)

For sampling, we adopt a probabilistic CL strategy that leverages the concept of CL in a nondeterministic fashion without discarding the original standard training practice, such as bucketing and mini-batching.

### 3.9 Transductive Ensemble Learning

Ensemble learning (Garmash and Monz, 2016), which aggregates multiple diverse models for inference, is a common practice to improve the performance of machine learning models. However, it has been observed that the conventional ensemble methods only bring marginal improvement for NMT when individual models are strong or there are a large number of individual models. Transductive Ensemble Learning (TEL) (Zhang et al., 2019) studies how to effectively aggregate multiple NMT models under the transductive setting where the source sentences of the test set are known. TEL uses all individual models to translate the source test set into the target language space and then fine-tune a strong model on the translated synthetic data, which significantly boosts strong individual models and benefits a lot from more individual models.

## 4 LLM-based MT System

### 4.1 System Overview

There is recently a surge in research interests in Transformer-based LLMs, such as ChatGPT (Wu et al., 2023a), GPT-4 (Achiam et al., 2023), and LLaMA (Touvron et al., 2023a,b). Benefiting from the giant model size and oceans of training data, LLMs can understand better the language structures and semantic meanings behind raw text, thereby showing excellent performance in a wide range of natural language processing (NLP) tasks. Although the training methodology of LLMs is simple, high computational requirements have limited the development of LLMs to a few players. In order to avoid training LLM from scratch, we chose to conduct research work on the open source Llama2-13b[9] (Touvron et al., 2023b) model. Llama2-13b is an autoregressive language model using an optimized transformer architecture that is pre-trained on 2 trillion tokens of data from publicly available

Figure 3: The training flow of LLM-based MT system.

sources. As shown in Figure 3, we train Llama2-13b into a powerful LLM-based MT model through three-stage training of CPT, SFT and CPO.

### 4.2 Continue Pre-training

LLMs like LLaMA are pre-trained on English-dominated corpora. This potentially explains their inadequate translation performance which necessitates cross-lingual capabilities. To ameliorate this, our first stage is to perform continue pre-training (CPT) on LLM with Chinese and English monolingual data to improve proficiency in Chinese and prevent forgetting of English knowledge. Previous studies also offer some clues that monolingual data help in translation. For instance, guo et al. (Guo et al., 2024) proposed a three-stage training method, which proved that using CPT can improve the performance of MT task in the SFT stage. Note that we use full fine-tuning at this stage.

### 4.3 Supervised Fine-tuning

LLMs have shown remarkable performance on a wide range of NLP tasks by leveraging in-context learning (Brown et al., 2020). However, this approach exhibits several drawbacks: performance is highly dependent on the quality of examples (Vilar et al., 2023), outputs are plagued by overgeneration (Bawden and Yvon, 2023), and inference cost are greatly increased by processing all input pairs. When parallel data is available, LLMs can perform supervised fine-tuning (SFT) on translation instructions (Li et al., 2024). Drawing inspiration from the recognized significance of data quality in other applications (Zhou et al., 2024),we use the cometkiwi model (Rei et al., 2022) to filter out large amounts of high-quality parallel data. Here, we use efficient lightweight low-rank adaptation (LoRA) fine-

Final Target Language Hypothesis

科学事实来自实验

MBR decoding with COMET

N-best Hypotheses             N-best Hypotheses

科学事实源于实验
科学事实源自实验
科学事实来自于实验
科学事实来源于实验
科学事实是从实验中得出的

科学事实源自于实验
科学事实来自实验
科学事实是由实验得出的
科学事实是通过实验得出的
科学事实产生于实验

Beam search             Temperature and nucleus sampling

NMT System             LLM-based MT System

Source Language Text

Scientific facts result from experiments

Figure 4: Choose the Final Translation from NMT and LLM hypotheses Using MBR Decoding.

tuning, where we apply LoRA to all modules of feed-forward network.

### 4.4 Contrastive Preference Optimization

Contrastive Preference Optimization (CPO) (Xu et al., 2024) aims to mitigate two fundamental shortcomings of SFT. First, SFT's methodology of minimizing the discrepancy between predicted outputs and gold-standard references inherently caps model performance at the quality level of the training data. This limitation is significant, as even human-written data, traditionally considered high-quality, is not immune to quality issues. Secondly, SFT lacks a mechanism to prevent the model from rejecting mistakes in translations. While strong translation models can produce high-quality translations, they occasionally exhibit minor errors, such as omitting parts of the translation. Preventing the production of these near-perfect but ultimately flawed translation is essential. To overcome these issues, we introduce CPO to train the LLM-based MT model using specially curated triplet preference data. Here, we construct a high-quality preference data for the WMT24 general MT task, and like the SFT stage, only update the weights of the added LoRA parameters.

### 4.5 Minimum Bayes Risk Decoding

Minimum Bayesian Risk (MBR) (Kumar and Byrne, 2004; Eikema and Aziz, 2020) decoding

aims to find the output that maximizes the expected utility function, which measures the similarity between the hypothesis and the reference. For MT, this could be an automated evaluation metric such as COMET (Rei et al., 2020). Garcia et al. (Garcia et al., 2023) train their own language models, sample multiple hypotheses and choose a final translation using MBR decoding, which has been shown to improve the translation capabilities of task-specific models (Fernandes et al., 2022). Subsequently, Farinhas et al. (Farinhas et al., 2023) find that MBR is also suitable for LLM-based MT. They provide a comprehensive study on ensembling translation hypotheses, proving that MBR decoding is a very effective method and can improve translation quality using a small number of samples. As shown in Figure 4, we simultaneously collect the N-best translations generated by the NMT system based on beam search and the N-best translations generated by the LLM-based MT system based on temperature and nucleus sampling (with t=0.8 and p=0.95), and then use MBR Decoding selects the final translation.

## 5 Experiment

### 5.1 Setup

We use the open-source fairseq (Ott et al., 2019) to train NMT models, and then use SacreBLEU

160

(Post, 2018)[10] and wmt20-comet-da model (Rei et al., 2020) to measure system performance. The main parameters are as follows: each model is trained using 8 GPUs, batch size is 6144, parameter update frequency is 2, and learning rate is 5e-4. The number of warmup steps is 4000, and model is saved every 1000 steps. The architecture we used is described in section 3.1. We adopt dropout, and the rate varies across different training phases. R-Drop is used in model training, and we set $\lambda$ to 5.

We use Llama2-13B as the backbone model of our LLM-based MT system. In our three-stage training process, the first stage uses full fine-tuning, and the last two stages use LoRA fine-tuning. If LoRA is used, lora_rank is 32, lora_alpha is 64, lora_dropout is 0.05, and lora_modules are "q_proj", "v_proj", "k_proj", "o_proj", "gate_proj", "down_proj", "up_proj". Furthermore, in the first and third stages, we use open-source ALMA [11] for training, while in the second stage, we use open-source llama-recipes [12] for training. The parameters during training are the default configurations of the corresponding codes.

## 5.2 Results

Tables 2 shows the evaluation results of en→zh NMT systems and LLM-based MT systems on WMT23 general test sets. On NMT systems, we use BiT and R-Drop to build a strong baseline, then use DD, FT and ST for data enhancement, and use AT and CL for more efficient training, and finally use TEL to ensemble multiple models ability. On LLM-based MT systems, we use CPT and SFT to build a strong baseline, and use CPO for further optimization. To ensemble two different types of translation systems, we use MBR decoding to select the final translation, which has been shown to be better than MBR decoding of a single translation system in terms of COMET scores.

## 5.3 Pre-processing and Post-processing

On the WMT24 general test set, we observe that there are some emoticons and URLs in the source text. To prevent the model from translating them incorrectly, we replace the emoticons and URLs with "Do Not Translate" (DNT) labels in pre-processing, and then restore the DNT labels back in post-processing. By doing so, we can reduce some translation errors for emoticons and URLs.

---

[10] https://github.com/mjpost/sacrebleu
[11] https://github.com/fe1ixxu/ALMA
[12] https://github.com/meta-llama/llama-recipes

| WMT23 general test set | BLEU | COMET |
|---|---|---|
| NMT baseline (BiT & R-Drop) | 54.24 | 0.6289 |
| + DD, FT & ST | 56.33 | 0.6580 |
| + AT | 57.03 | 0.6648 |
| + CL | 58.58 | 0.6830 |
| + TEL | **59.34** | 0.6928 |
| + NMT MBR | 58.88 | 0.7178 |
| LLM-based MT baseline (CPT & SFT) | 52.18 | 0.6553 |
| + CPO | 53.09 | 0.6907 |
| + LLM-based MT MBR | 52.16 | 0.7102 |
| + NMT & LLM-based MT MBR | 56.41 | **0.7234** |

Table 2: BLEU and COMET scores of en→zh NMT systems and LLM-based MT systems.

## 6 Conclusion

This paper presents the submission of HW-TSC to the WMT24 general MT Task. On the one hand, we use training strategies such as R-Drop, BiT, DD, FT, BT, AT, CL, and TEL to train the NMT system based on the deep Transformer-big architecture. On the other hand, we use CPT, SFT, and CPO to train the LLM-based MT system. Finally, we use MBR decoding to select the final translation result from the hypotheses generated by these two systems. By using these enhancement strategies, our submission achieved a competitive result in the final evaluation. Relevant experimental results also demonstrate the effectiveness of our strategies.

## References

Idris Abdulmumin. 2021. Enhanced back-translation for low resource neural machine translation using self-training. In *Information and Communication Technology and Applications: Third International Conference, ICTA 2020, Minna, Nigeria, November 24–27, 2020, Revised Selected Papers*, volume 1350, page 355. Springer Nature.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of bloom. *arXiv preprint arXiv:2303.01911*.

Peter F Brown, John Cocke, Stephen A Della Pietra, Vincent J Della Pietra, Frederick Jelinek, John Lafferty, Robert L Mercer, and Paul S Roossin. 1990. A statistical approach to machine translation. *Computational linguistics*, 16(2):79–85.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind

Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Isaac Caswell, Ciprian Chelba, and David Grangier. 2019. Tagged back-translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 53–63.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Liang Ding, Di Wu, and Dacheng Tao. 2021. Improving neural machine translation by bidirectional training. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3278–3284.

Chris Dyer, Victor Chahuneau, and Noah A Smith. 2013. A simple, fast, and effective reparameterization of ibm model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648.

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, page 489. Association for Computational Linguistics.

Bryan Eikema and Wilker Aziz. 2020. Is map decoding all you need? the inadequacy of the mode in neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021. Beyond english-centric multilingual machine translation. *Journal of Machine Learning Research*, 22(107):1–48.

António Farinhas, José de Souza, and André FT Martins. 2023. An empirical study of translation hypothesis ensembling with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11956–11970.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic bert sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891.

Patrick Fernandes, António Farinhas, Ricardo Rei, José GC de Souza, Perez Ogayo, Graham Neubig, and André FT Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings*

*of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412.

Mara Finkelstein and Markus Freitag. 2023. Mbr and qe finetuning: Training-time distillation of the best and most expensive decoding methods. *arXiv preprint arXiv:2309.10966*.

Xavier Garcia, Yamini Bansal, Colin Cherry, George Foster, Maxim Krikun, Melvin Johnson, and Orhan Firat. 2023. The unreasonable effectiveness of few-shot learning for machine translation. In *International Conference on Machine Learning*, pages 10867–10878. PMLR.

Ekaterina Garmash and Christof Monz. 2016. Ensemble learning for multi-source neural machine translation. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1409–1418.

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 639–649.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Rui Jiao, Zonghan Yang, Maosong Sun, and Yang Liu. 2021. Alternated training with synthetic and authentic data for neural machine translation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1828–1834.

Wenxiang Jiao, Jen-tse Huang, Wenxuan Wang, Zhiwei He, Tian Liang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Parrot: Translating during chat using large language models tuned with human translation and feedback. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15009–15020.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hérve Jégou, and Tomas Mikolov. 2016. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open

source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *EMNLP 2018*, page 66.

Shankar Kumar and Bill Byrne. 2004. Minimum bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 169–176.

Jiahuan Li, Hao Zhou, Shujian Huang, Shanbo Cheng, and Jiajun Chen. 2024. Eliciting the translation ability of large language models via multilingual finetuning with translation instructions. *Transactions of the Association for Computational Linguistics*, 11:576–592.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzman. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756.

Xuan-Phi Nguyen, Shafiq Joty, Wu Kui, and Ai Ti Aw. 2020. Data diversification: a simple strategy for neural machine translation. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pages 10018–10029.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702.

Ricardo Rei, Marcos Treviso, Nuno M Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José GC de Souza, Taisiya Glushkova, Duarte M Alves, Alon Lavie, et al. 2022. Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task. *WMT 2022*, page 634.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *54th Annual Meeting of the Association for Computational Linguistics*, pages 86–96. Association for Computational Linguistics (ACL).

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Tim Van Erven and Peter Harremos. 2014. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting palm for translation: Assessing strategies and performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15406–15427.

Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao. 2019. Learning deep transformer models for machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.

Wei Wang, Ye Tian, Jiquan Ngiam, Yinfei Yang, Isaac Caswell, and Zarana Parekh. 2020a. Learning a multi-domain curriculum for neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7711–7723.

Yiren Wang, Lijun Wu, Yingce Xia, Tao Qin, ChengXiang Zhai, and Tie-Yan Liu. 2020b. Transductive ensemble learning for neural machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6291–6298.

Daimeng Wei, Zongyao Li, Zhanglin Wu, Zhengzhe Yu, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Minghan Wang, Lizhi Lei, Min Zhang, et al. 2021. Hwtsc's participation in the wmt 2021 news translation shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 225–231.

Daimeng Wei, Zhiqiang Rao, Zhanglin Wu, Shaojun Li, Yuanchang Luo, Yuhao Xie, Xiaoyu Chen, Hengchao Shang, Zongyao Li, Zhengzhe Yu, et al. 2022. Hwtsc's submissions to the wmt 2022 general machine translation shared task. In *Proceedings of the Seventh Conference on Machine Translation*, pages 403–410.

Daimeng Wei, Zhanglin Wu, Hengchao Shang, Zongyao Li, Minghan Wang, Jiaxin Guo, Xiaoyu Chen, Zhengzhe Yu, and Hao Yang. 2023. Text style transfer back-translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pages 7944–7959, Toronto, Canada.

Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, Tie-Yan Liu, et al. 2021. R-drop: Regularized dropout for neural networks. *Advances in Neural Information Processing Systems*, 34:10890–10905.

Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. 2023a. A brief overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal of Automatica Sinica*, 10(5):1122–1136.

Zhanglin Wu, Daimeng Wei, Zongyao Li, Zhengzhe Yu, Shaojun Li, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Yuhao Xie, Lizhi Lei, et al. 2023b. The path to continuous domain adaptation improvements by hw-tsc for the wmt23 biomedical translation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 271–274.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *arXiv preprint arXiv:2401.08417*.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498.

Wen Yang, Chong Li, Jiajun Zhang, and Chengqing Zong. 2023. Bigtranslate: Augmenting large language models with multilingual translation capability over 100 languages. *arXiv preprint arXiv:2305.18098*.

Jiali Zeng, Fandong Meng, Yongjing Yin, and Jie Zhou. 2024. Teaching large language models to translate with comparison. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 17, pages 19488–19496.

Xuan Zhang, Pamela Shapiro, Gaurav Kumar, Paul McNamee, Marine Carpuat, and Kevin Duh. 2019. Curriculum learning for domain adaptation in neural machine translation. In *Proceedings of NAACL-HLT*, pages 1903–1915.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

Zhi-Hua Zhou, Jianxin Wu, and Wei Tang. 2002. Ensembling neural networks: many could be better than all. *Artificial intelligence*, 137(1-2):239–263.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2765–2781.

# CycleGN: a Cycle Consistent approach for Neural Machine Translation

**Sören Dréano**
ML-Labs
Dublin City University
`soren.dreano2@mail.dcu.ie`

**Derek Molloy**
School of Electronic Engineering
Dublin City University
`derek.molloy@dcu.ie`

**Noel Murphy**
School of Electronic Engineering
Dublin City University
`noel.murphy@dcu.ie`

## Abstract

CycleGN is a Neural Machine Translation framework relying on the Transformer architecture. Its approach is similar to a Discriminator-less CycleGAN, specifically tailored for non-parallel text datasets.

The foundational concept of our research posits that in an ideal scenario, retro-translations of generated translations should revert to the original source sentences. Consequently, a pair of models can be trained using a Cycle Consistency Loss only, with one model translating in one direction and the second model in the opposite direction.

As a contribution to the WMT24 challenge, this study explores the efficacy of the CycleGN architectural framework in learning translation tasks across two language pairs, English-Chinese and German-English, under two distinct non-parallel dataset conditions: permuted and non-intersecting. Our findings demonstrate the robust adaptability of CycleGN in learning translation tasks, irrespective of the language pair.

## 1 Introduction

The introduction of the Transformer architecture (Vaswani et al., 2017) marked a significant advancement in the field of Machine Translation, witnessing widespread adoption since its inception. Although self-attention mechanisms were not novel and had been investigated in prior studies (Bahdanau et al., 2016), the Transformer model demonstrated its formidable capabilities within Natural Language Processing (NLP). Characterized by its parallelized structure, the Transformer architecture facilitated computational efficiency, enabling the incorporation of a larger number of parameters. This enhancement has been exemplified in NLP systems like Charles University Block-Backtranslation-Improved Transformer Translation (cubbitt) (Popel et al., 2020), which have surpassed

the performance levels of human professionals in certain contexts.

Neural Machine Translation (NMT) datasets necessitate substantial text corpora, structured as aligned pairs. This alignment implies the requirement for sentences with equivalent meaning to be present in a minimum of two distinct languages, enabling the initiation of model training to forge linguistic linkages. Ongoing initiatives, including OPUS (Tiedemann and Thottingal, 2020) and Tatoeba (Tiedemann, 2012), are committed to facilitating public access to these datasets. Parallel datasets comprise a small subset of the volume of data in monolingual datasets.

Despite the widespread availability of large parallel corpora for numerous language pairs, the capacity to employ solely monolingual datasets would substantially expand the pool of training data. This approach is particularly beneficial for languages with scarce parallel text corpora.

Regardless of the remarkable efficacy exhibited by Large Language Models (LLM) in NMT without the necessity of exclusive training on parallel data (Zhu et al., 2023), their considerable magnitude renders them costly in terms of both training and operation. This economic burden consequently restricts their widespread availability.

Back-translation (Sennrich et al., 2016) is a technique leveraging a trained MT (Machine Translation) model to translate sentences from a monolingual dataset to produce corresponding pairs, thereby synthetically augmenting the training data. Our research is founded on the premise that the process of translating a sentence from a source language to a target language, followed by its retro-translation from the target language back to the source language, allows for the measurement of the disparity between the original and the machine-retro-translated sentences. This disparity serves as a metric to assess the efficacy of the models and facilitates the backpropagation of gradients within

165

the networks. Notably, this methodology has been previously implemented in the realm of Image-to-Image Translation, as evidenced in the renowned CycleGAN framework from Zhu et al. (2017).

## 2 Previous work

The TextCycleGAN model (Lorandi et al., 2023), while not utilizing the Transformer architecture nor operating within the MT field, introduced an innovative strategy for text style transfer. This approach employed a CycleGAN on the Yelp dataset to facilitate the learning of mappings between positive and negative textual styles, notably in the absence of paired examples.

Shen et al. (2017) exemplified the feasibility of training two encoder-decoder networks in an unsupervised manner that enables the sharing of a latent space, thereby permitting style transfer. Lample et al. (2018), adopting a similar technique within the MT context, substantiated that the use of parallel datasets is not a prerequisite for effective translation.

## 3 Definitions

Machine Translation models are most commonly trained using "parallel" datasets, which are structured collections of text pairs. Each pair comprises a segment of text in a source language and its translation in the target language. By providing direct translations, models learn correspondences between text units to map the source language to the target language.

A non-parallel dataset on the other hand does not consist in pairs of text segments, consequently the source and target sentences do not share any explicit correspondence. Such a dataset can be creating by combining any two monolingual datasets of two distinct languages and adjusting for the number of samples. In the context of this research, two sub-categories of non-parallel datasets are introduced.

### 3.1 Permuted dataset

A "permuted" dataset is defined as a parallel dataset wherein the sentences of one language have been systematically rearranged. Consequently, this results in a non-parallel corpus where it is guaranteed that each sentence has a corresponding translation located at an unspecified index within the dataset. The authors postulate that when employing sufficiently large monolingual datasets, which are not

derived from permuted parallel corpora, it is likely that most sentences will possess an accurate translation "somewhere" within the dataset.

### 3.2 Non-intersecting dataset

A "non-intersecting" dataset is a non-parallel dataset for which it is guaranteed that no sentence has an exact translation. A non-intersecting dataset is derived from a meticulously curated parallel dataset devoid of duplicate entries. Two unique sets of natural integers are produced, each functioning as an index list of phrases to retain for each respective language.

## 4 Datasets

The datasets employed in this study are the English-German and Chinese-English language pairs from the WMT23 challenge (Kocmi et al., 2023). The data released for the WMT23 General MT task can be freely used for research purposes. Due to the current implementation's high computational demands, the models were not trained for the entirety of an epoch. Specifically, only 10% of the English-German dataset was used, while about half of the Chinese-English dataset in the non-intersecting condition.

| Type | English-German | Chinese-English |
|---|---|---|
| Permuted | 27,801,496 | 27,801,496 |
| Non-intersecting | 27,801,496 | 17,676,442 |
| Original dataset | 295,805,439 | 35,452,884 |

Table 1: Number of sentences used during training depending on the dataset type

## 5 Training

For greater clarity, the mathematical notations from the original CycleGAN work will be employed in the present study. Given two languages $\mathcal{X}$ and $\mathcal{Y}$ with appropriate datasets, the objective is to obtain two NMT models $\mathcal{G} : \mathcal{X} \mapsto \mathcal{Y}$ and $\mathcal{F} : \mathcal{Y} \mapsto \mathcal{X}$ such that if the translations are perfect, $\mathcal{G}(\mathcal{F}(y)) = y$ and $\mathcal{F}(\mathcal{G}(x)) = x$, with $x \in \mathcal{X}$ and for $y \in \mathcal{Y}$.

By using the Cross-Entropy Loss (CEL) (Zhang and Sabuncu, 2018) in the role of the Cycle Consistency Loss (CCL), we can determine the distance between the original sentence and its double translation in order to compute the gradients.

As in the original CycleGAN work, our current study also implements an Identity Loss (IL), which also relies on the CEL, to help with the training stability. As $\mathcal{G}$ consists in a mapping $\mathcal{X} \mapsto \mathcal{Y}$, if

given an input $y \in \mathcal{Y}$, the input should remain unchanged such that $\mathcal{G}(y) = y$. The same loss is applied to $\mathcal{F}$ between $\mathcal{F}(x)$ and $x$, as displayed in Figure 1.

## 5.1 Model architecture

The architecture used for both models, $\mathcal{G}$ and $\mathcal{F}$, is the Marian framework (Junczys-Dowmunt et al., 2018) implemented by Huggingface's Transformers library (Wolf et al., 2020), which is licensed under the Apache Licence. While most parameters follow the default configuration, Table 2 references the changes that were made in order to reduce the computational cost of the architecture.

| Parameter | Huggingface | Current work |
|---|---|---|
| Vocabulary size | 58,101 | 32,000 |
| Encoder layers | 12 | 6 |
| Decoder layers | 12 | 6 |
| Encoder attention heads | 16 | 8 |
| Decoder attention heads | 16 | 8 |
| Encoder feed-forward | 4096 | 2048 |
| Decoder feed-forward | 4096 | 2048 |
| Position embeddings | 1024 | 128 |
| Activation function | GELU | ReLU |

Table 2: Non-default parameters in the configuration of Marian Transformer models

## 5.2 Vocabulary organization

NMT models usually employ either a unified tokenizer or two distinct tokenizers. In the case of a single tokenizer, it is trained using sentences from both the source and target distributions, avoiding any duplicates. This approach facilitates the sharing of the encoder and decoder embedding layers, thereby diminishing computational demands and enhancing model accuracy (Press and Wolf, 2017).

Conversely, the alternative approach entails training one tokenizer on the source distribution and another one on the target distribution. While this method restricts the possibility of tying embeddings, it can potentially double the vocabulary size without increasing the dimensions of the embeddings. The overall vocabulary size of the model in this scenario, is the cumulative total of the two individual vocabularies, barring shared tokens like punctuation symbols.

While contemporary Transformer models like Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2019) and Generative Pre-trained Transformers (GPT) (Radford et al., 2018) typically utilize a single tokenizer, this study

introduces a novel vocabulary methodology that amalgamates the aforementioned approaches. This method involves training two tokenizers, each for a respective language and with half the vocabulary size. Subsequently, the identifiers of one tokenizer are adjusted to prevent overlap, yielding a result analogous to a single tokenizer that includes duplicates across languages. It is important to note that special tokens such as $< eos >$ (End of Sentence) and $< pad >$ (Padding) are shared and not duplicated. This strategy is designed to simplify model analysis during development, albeit at the expense of a reduced vocabulary.

## 5.3 Obtaining labels

In the training process of a Transformer model, it is imperative to have prior knowledge of the labels, as the decoder predicts tokens sequentially. Each token prediction, barring the initial one, is contingent upon all preceding predictions. By possessing prior knowledge of the reference translation, it becomes feasible to contrast each predicted token against the ground truth, enabling the calculation of the loss at every step.

Nevertheless, in the case of non-parallels datasets, the labels are by definition not known in advance. It is therefore not possible to calculate the loss after each predicted token. Furthermore, the act of selecting the most probable token for each prediction constitutes a non-differentiable operation, thus precluding the possibility of backpropagation once the sentence is fully generated.

Naturally, in inference mode, Transformers are able to generate sentences without labels. Thus, the first step is to generate the pseudo-labels $\hat{x}$ and $\hat{y}$, where $\hat{x}$ is used as the label of $y$ and $\hat{y}$ as the label of $x$. Even though this step cannot be used to compute the gradients, it is crucial for the entire process.

$\hat{\hat{x}}$ is computed from from $\mathcal{F}(\hat{y})$ with $x$ as the label, and $\hat{\hat{y}}$ is computed from $\mathcal{G}(\hat{x})$ with $y$ as the label. The CCL is applied between $\hat{\hat{x}}$ and $x$, and between $\hat{\hat{y}}$ and $y$ to compute the gradients and update the weights of $\mathcal{G}$ and $\mathcal{F}$.

## 5.4 A Discriminator-less GAN

The CycleGAN methodology, as indicated by its nomenclature, is predicated on the Generative Adversarial Network (GAN) framework, initially introduced in Goodfellow et al. (2014). This paradigm involves the training of a Generator model in conjunction with another model, termed

Figure 1: CycleGN training process

the Discriminator. The Discriminator is specifically trained to distinguish between authentic samples drawn from the dataset and synthetic samples produced by the Generator. In the CycleGAN training process, the Discriminators intervene after the generation of $\hat{x}$ and $\hat{y}$, helping the training of the Generators. However, as mentioned in Section 5.3, there can be no gradient computation during the generation of $\hat{x}$ and $\hat{y}$ in a Transformer and as such, Discriminators cannot be used in the present work. This is why CycleGN is not an "Adversarial" approach, hence the name.

# 6 Pre-training

During the development of CycleGN, a critical issue became clear, which prevented the model's ability to converge and learn effectively. As described in Section 5.3, the first step of the CycleGN framework is to generate $\hat{x}$ and $\hat{y}$. During the first initialisations, these pseudo-labels will be generated randomly and will depend only on the initialization of the weights of $\mathcal{G}$ and $\mathcal{F}$. However, the models consistently converge towards a trivial solution wherein by merely reproducing the input, they satistisfy the loss function criteria without achieving any meaningful learning or transformation of the data.

## 6.1 Absence of intermediate evaluation

As there is no Discriminator to ensure that $\hat{x}$ belongs to $\mathcal{X}$ and $\hat{y}$ belongs to $\mathcal{Y}$, $\mathcal{G}$ and $\mathcal{F}$ converge towards $x = \hat{y} = \hat{\hat{x}}$ and $y = \hat{x} = \hat{\hat{y}}$, as this approach achieves an optimal outcome on the CCL function, registering a value of zero, as schematised in Figure 2.



Figure 2: In the absence of a Discriminator $y \in \mathcal{Y}$ and pre-training is not employed, the CycleGN architecture will converge towards a state where no translation happens and still perfectly satisfy the CCL function

## 6.2 Moving away from the easiest path

Masked Language Modeling (MLM) is a pre-training strategy implemented in BERT, where a specified proportion of the input tokens are sub-

stituted with a unique $< mask >$ token. The objective of the neural network under this paradigm is to accurately reconstruct the original sentence from this degraded input. This process enables the model to discern intricate relationships between words and to develop a profound representation of the language. This pre-training has revealed excellent performances in diverse NLP application such as sentiment analysis (Alaparthi and Mishra, 2021), text classification (Sun et al., 2020), Named Entity Recognition (NER) (Souza et al., 2020) (Chang et al., 2021) (Akhtyamova, 2020) and paraphrase detection (Khairova et al., 2022).

As MLM does not require any labels, as the labels are generated from the dataset, it is perfectly adapted to the CycleGN approach. A single model $\mathcal{H}$ is trained on the non-parallel dataset to reconstruct both languages, with 15% of the input tokens masked. This model $\mathcal{H}$ has the exact same architecture as $\mathcal{G}$ and $\mathcal{F}$. When training the CycleGN, rather than randomly initializing $\mathcal{G}$ and $\mathcal{F}$, the parameters from $\mathcal{H}$ are directly transferred to both $\mathcal{G}$ and $\mathcal{F}$. Indeed, as $\mathcal{H}$ learns to reconstruct both language $\mathcal{X}$ and $\mathcal{Y}$, it can be used to initialize both networks. Figure 3 shows the training process of $\mathcal{H}$.



Figure 3: Masked Language Modeling training process

## 7 Training stability

It is crucial for the CycleGN framework that the two models exhibit approximately equivalent levels of performance. Given the interdependent nature of these models, where the output of one serves as the input for the other, maintaining consistency between them during training is imperative. Without a strategy in place to prevent the performance of the models from diverging, it is possible for one model to gain the "upper hand" over the other.

### 7.1 Divergence between the Generators

Figure 4 presents the evolution of the CCL of an early prototype of CycleGN and it can clearly be seen that one of the two generators, $\mathcal{F}$, ends up performing much better than its counterpart $\mathcal{G}$, which blocks any future training.



Figure 4: Evolution of the Cross-Entropy Loss during the training of an early prototype on the permuted German-English dataset

### 7.2 Gradient Clipping

Gradient clipping is a technique utilized in the training of Deep Learning (DL) models, to address the problem of "exploding" gradients. This issue occurs when gradients escalate to excessively high values during training, leading to numerical instability and impeding the model's convergence to an optimal solution.

Gradient clipping can be implemented through two primary methods: norm clipping and value clipping. Norm clipping involves establishing a threshold on the overall magnitude of the gradient vector. On the other hand, value clipping involves individually adjusting elements of the gradient vector that exceed the specified threshold.

By clipping the gradients by norm, with a threshold of 1.0, as advised by the Huggingface library, the training stabilizes and the divergence between $\mathcal{G}$ and $\mathcal{F}$ disappears.

Figure 5 demonstrates how the addition of gradient clipping helps with training stability during the training of the permuted German-English model.

### 7.3 Batch size

The original CycleGAN research mentions using a batch size of 1, and while they did not state the reason in the research paper, one of the authors explained it in a GitHub issue (Junyanz, 2017) as a lack of GPU memory.

Rajput et al. (2021) examined the impact of batch size within the CycleGAN architecture, observing a significant decline in performance the more the batch size is increased. This deterioration was evident both through the example images presented in

Figure 5: Evolution of the Cross-Entropy Loss during the training of the permuted German-English models

that study and through the calculated cosine dissimilarity, indicating inferior model performance with larger batch sizes. However, quality was achieved at the expense of computational efficiency, as the training duration to achieve 200 epochs was 8 hours with a batch size of 1, but this was reduced to just 2 hours with a batch size of 64.

In the context of this research, however, the trade-off between quality improvement and computing resource, as observed in the aforementioned study, does not hold true. Utilizing a batch size of 1 in the CycleGN experiments hindered any form of convergence. Consequently, a batch size of 32 was selected, as it represents the maximum capacity that could be accommodated within the available 24GB of GPU memory of the NVidia 4090 used for this work.

### 7.4 One large epoch or multiple smaller ones?

The CycleGAN framework is recognized for its computational expense due to several inherent factors. Primarily, as CycleGAN operates on the principle of cycle consistency, it necessitates the training of two GANs simultaneously – one for each direction of the transformation. This structure requires substantial computational resources, as each GAN consists of both a Generator and a Discriminator.

The resource-intensiveness of the CycleGAN process, thus limits the size of the dataset that can be used in a reasonable time. This necessitated a decision between training for a single epoch on a large dataset, or training for multiple epochs on a smaller corpus arose.

The CycleGN framework was compared on the permuted German-English dataset under four different conditions:

1. One epoch containing 1% of the dataset

2. Five epochs containing 0.2% of the dataset

3. One epoch containing 2% of the dataset

4. Five epochs containing 0.4% of the dataset

The Crosslingual Optimised Metric for Evaluation of Translation (COMET) score (Rei et al., 2020) was selected as our comparison criterion, as this metric has proven to be one of the most effective in recent WMT competitions, according to Kocmi et al. (2022), due to its strong correlation with human judgment, aligning well with our goal of mirroring human evaluative standards. Multiple COMET models have been made available and the default "wmt22-comet-da" model was chosen. The average scores obtained on 10,000 test sentences that were not part of the model training set are presented in Table 3.

| Condition | English->German | German->English |
|---|---|---|
| 1 | 0.2727 | 0.2715 |
| 2 | 0.2411 | 0.2635 |
| 3 | 0.2741 | 0.2665 |
| 4 | 0.2258 | 0.2658 |

Table 3: COMET scores of CycleGN models depending on the permuted German-English dataset condition

Models exposed to a larger portion of the total dataset demonstrate superior performance compared to those limited to a smaller, repetitive subset, especially when the dataset encompasses over half a million to a million sentences. The authors extrapolate this result to larger datasets and thus chose to train the CycleGN models for a single epoch on the largest dataset possible.

## 8 Results

Even if tracking the CCL is an inexpensive manner to estimate the progress of the training of the CycleGN architecture, a low loss value can also hide an absence of translation, as mentioned in Section 6.1. This is why an evaluation metric such as COMET is crucial to assess the progression of the CycleGN framework.

### 8.1 Evolution of COMET score during training

To measure the performances of CycleGN, every 1,000th batch the CCL was averaged and 1,000 sentences from the test set were translated to compute the COMET score.

Figures 6, 7, 8 and 9 demonstrate that the actual quality of translation, as measured by the COMET metric, increases with time. Figures 6 through 9 illustrate a progressive enhancement in the translation quality over time, as quantified by the COMET metric. This enhancement is observed respectively in the permuted and non-intersecting German-English models (Figures 6 and 7), as well as in the permuted and non-intersecting English-Chinese models (Figures 8 and 9). Figures 6 and 7 exhibit a sudden drop in the increase of accuracy, which is acknowledged by the authors. This anomaly will be thoroughly examined and discussed in a subsequent academic study.



Figure 7: Evolution of the COMET score during the training of the non-intersecting German-English models



Figure 8: Evolution of the COMET score during the training of the permuted Chinese-English models



Figure 6: Evolution of the COMET score during the training of the permuted German-English models

### 8.2 COMET Scores post-training completion

After the end of the training, a test set of 10,000 sentences per language were translated and the COMET scores are displayed in Table 4. In order to give a point of comparison, architecture-matched models using the original parallel datasets were trained. As in the case of the CycleGN training, these parallel models were only trained for a single epoch on the exact same number of sentences as the permuted models were.

The authors expected the COMET score of the CycleGN to be inferior to architecture-matched models trained using parallel corpora, as information is by definition lost during the permutation of



Figure 9: Evolution of the COMET score during the training of the non-intersecting Chinese-English models

the parallel datasets. However, the authors argue that the differences between the scores is likely smaller with larger datasets.

|                  | English → German | German → English |
|------------------|:----------------:|:----------------:|
| Permuted         | 0.505            | 0.537            |
| Non-intersecting | 0.556            | 0.579            |
| Parallel         | 0.780            | 0.775            |

Table 4: COMET score of the German-English models

|                  | English → Chinese | Chinese → English |
|------------------|:-----------------:|:-----------------:|
| Permuted         | 0.425             | 0.537             |
| Non-intersecting | 0.382             | 0.448             |
| Parallel         | 0.000             | 0.749             |

Table 5: COMET score of the Chinese-English models

## 9 Future Work

Further investigations will benefit from the incorporation of a more extensive dataset and an exploration of larger model architectures.

### 9.1 Larget dataset

The current work has been trained on a small dataset compared to MT standards. Future work should try to see how convergence progresses with more iterations. Further computational optimizations are probably necessary to shorten the training time required.

### 9.2 Larger models

The current architecture relies on a total of 158,769,152 parameters, which is only about a third of the size of the default in the Huggingface library. Although Tables 4 and 5 demonstrate that the current number of parameters, when trained using a parallel dataset, is capable of producing better translations than when exposed to permuted and non-intersecting datasets, an increase in both the number of epochs and the size of the dataset should be prioritized, larger models being common in NMT.

## 10 Source Code

The source code of CycleGN is available at https://github.com/SorenDreano/CycleGN.

## Limitations

The investigation acknowledges certain inherent limitations which may impact the generalizability and applicability of the findings.

### Language diversity

Another issue that arises from the computing cost of CycleGN is the lack in language diversity. Indeed, our current work only used the English-German and Chinese-English language pairs. Consequently, it cannot be certain that the approach presented can be applied to other languages and all alphabets. This is why CycleGN is taking part in WMT24, to explore the framework's performance on a wide range of language pairs.

### Training limitations

Since training a CycleGN model is particularly costly, there is a trade-off between training models on all language pairs, or choosing a subset of these pairs to train fewer models with more iterations and on a greater number of examples. In order to demonstrate the effectiveness of CycleGN on a wide range of language pairs, the first choice was made, i.e. to train models on all pairs, even if this means obtaining inferior results.

### Unused models

Unlike the previous edition (Kocmi, 2023), where most language pairs were bidirectional, i.e. the evaluations were to and from, the 2024 General Translation task is unidirectional. This means that for each language pair, it is sufficient to train a model that translates from the source to the target.

This is not a change that is favourable to CycleGN, since it is a bidirectional training architecture. Indeed, its cyclical nature means that one model must be trained from one language to another, and another model must complete the cycle, i.e. from this second language to the first. In other words, half the time spent training CycleGN is spent training a model which only serves to train the first, but which will never be evaluated in the contest.

This change has been accompanied by an increase in the number of language pairs, from 6 bidirectional and 2 unidirectional in 2023 to 11 unidirectional in 2024.

### Monolingual datasets

During the WMT challenge, teams are provided with monolingual datasets. Although this dataset format is perfectly suited to CycleGN training, they have been discarded for two reasons. The first is that for the majority of language pairs, the parallel datasets supplied have been truncated in order to reduce training time. The second is related to the construction of permuted and non-intersecting datasets, since it is preferable to build them from non-parallel datasets, as detailed in Section 3.

## Reduced dataset sizes

The datasets were truncated to obtain a maximum of 27,801,496 sentences for training and 100,000 sentences for the development set. The final size of the datasets used and the number of epochs is shown in Table 6 for permuted models and Table 7 for non-intersecting models. While the permuted models have all been trained, this was not the case for the non-intersecting models, due to lack of time.

## Training time

To make it possible to train so many models, several machines were used, with different technical characteristics, in particular different GPUs. However, by estimating the training time according to the number of sentences in the dataset and the GPU used, the total training time for all the models trained on the WMT24 datasets represents approximately 3,700 hours on an NVidia 4090.

## Ethics Statement

This study, focusing on the training of NMT models using non-parallel datasets, adheres to the highest ethical standards in research. We recognize the critical importance of ethical considerations in computational linguistics and machine learning, especially as they pertain to data sourcing, model development, and potential impacts on various linguistic communities.

Our research utilizes publicly available, non-parallel linguistic datasets. We ensure that all data is sourced following legal and ethical guidelines, respecting intellectual property rights and privacy concerns.

In our commitment to scientific integrity, we maintain transparency in our research methodologies, model development, and findings. We aim to make our results reproducible and accessible to the scientific community, contributing positively to the field of machine translation.

## Acknowledgements

## References

Liliya Akhtyamova. 2020. Named entity recognition in spanish biomedical literature: Short review and bert model. In *2020 26th Conference of Open Innovations Association (FRUCT)*, pages 1–7.

Shivaji Alaparthi and Manit Mishra. 2021. Bert: a sentiment analysis odyssey. *Journal of Marketing Analytics*, 9(2):118–126.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2016. Neural machine translation by jointly learning to align and translate.

Yuan Chang, Lei Kong, Kejia Jia, and Qinglei Meng. 2021. Chinese named entity recognition method based on bert. In *2021 IEEE International Conference on Data Science and Computer Application (ICDSCA)*, pages 294–299.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial networks.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Junyanz. 2017. Question: Batch size · issue 27 · junyanz/pytorch-cyclegan-and-pix2pix.

Nina Khairova, Anastasiia Shapovalova, Orken Mamyrbayev, Nataliia Sharonova, and Kuralay. 2022. Using bert model to identify sentences paraphrase in the news corpus. In *CEUR Workshop Proceedings, volume 3171*, pages 38–48.

Tom Kocmi. 2023. Shared task: General machine translation.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

| Language pair | Parallel sentences in WMT24 dataset | Sentences used | Number of epochs |
|---|---|---|---|
| Czech-Ukrainian | 10,757,756 | 10,657,756 | 1 |
| English-Chinese | 55,216,751 | 27,801,496 | 1 |
| English-Czech | 56,288,239 | 27,801,496 | 1 |
| English-German | 295,805,439 | 27,801,496 | 1 |
| English-Hindi | 315,070 | 314,070 | 10 |
| English-Icelandic | 23,434,361 | 23,334,361 | 1 |
| English-Japanese | 33,875,119 | 27,801,496 | 1 |
| English-Russian | 75,961,169 | 27,801,496 | 1 |
| English-Spanish | 626,076,911 | 27,801,496 | 1 |
| English-Ukrainian | 16,062,359 | 15,962,359 | 1 |
| Japanese-Chinese | 22,642,571 | 22,542,571 | 1 |

Table 6: Comparison between the number of sentences available in the WMT24 dataset and the number of sentences used to train the permuted models depending on the language pair

| Language pair | Parallel sentences in WMT24 dataset | Sentences used | Number of epochs |
|---|---|---|---|
| English-Chinese | 55,216,751 | 17,676,442 | 1 |
| English-Czech | 56,288,239 | 27,801,496 | 1 |
| English-German | 295,805,439 | 27,801,496 | 1 |
| English-Russian | 75,961,169 | 27,801,496 | 1 |

Table 7: Comparison between the number of sentences available in the WMT24 dataset and the number of sentences used to train the non-intersecting models depending on the language pair

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018. Unsupervised machine translation using monolingual corpora only.

Michela Lorandi, Maram A.Mohamed, and Kevin McGuinness. 2023. Adapting the CycleGAN Architecture for Text Style Transfer. *Irish Machine Vision and Image Processing Conference*.

Martin Popel, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature Communications*, 11(4381):1–15.

Ofir Press and Lior Wolf. 2017. Using the output embedding to improve language models.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. *OpenAI*.

Pranjal Singh Rajput, Kanya Satis, Sonnya Dellarosa, Wenxuan Huang, and Obinna Agba. 2021. cgans for cartoon to real-life images.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data.

Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment.

Fábio Souza, Rodrigo Nogueira, and Roberto Lotufo. 2020. Portuguese named entity recognition using bert-crf.

Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2020. How to fine-tune bert for text classification?

Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-MT — Building open translation services for the World. In *Proceedings of the 22nd Annual Conferenec of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In *Proceedings of the Eight International Conference on Language Resources and*

*Evaluation (LREC'12)*, Istanbul, Turkey. European Language Resources Association (ELRA).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface's transformers: State-of-the-art natural language processing.

Zhilu Zhang and Mert R. Sabuncu. 2018. Generalized cross entropy loss for training deep neural networks with noisy labels.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Multilingual machine translation with large language models: Empirical results and analysis.

# UvA-MT's Participation in the WMT24 General Translation Shared Task

**Shaomu Tan      Di Wu      David Stap      Seth Aycock      Christof Monz**

Language Technology Lab

University of Amsterdam

`{s.tan, d.wu, d.stap, s.aycock, c.monz}@uva.nl`

## Abstract

Fine-tuning Large Language Models (FT-LLMs) with parallel data has emerged as a promising paradigm in recent machine translation research. In this paper, we explore the effectiveness of FT-LLMs and compare them to traditional encoder-decoder Neural Machine Translation (NMT) systems under the WMT24 general MT shared task for English to Chinese direction. We implement several techniques, including Quality Estimation (QE) data filtering, supervised fine-tuning, and post-editing that integrate NMT systems with LLMs.

We demonstrate that fine-tuning LLaMA2 on a high-quality but relatively small bitext dataset (100K) yields COMET results comparable to much smaller encoder-decoder NMT systems trained on over 22 million bitexts. However, this approach largely underperforms on surface-level metrics like BLEU and ChrF. We further control the data quality using the COMET-based quality estimation method. Our experiments show that 1) filtering low COMET scores largely improves encoder-decoder systems, but 2) no clear gains are observed for LLMs when further refining the fine-tuning set. Finally, we show that combining NMT systems with LLMs via post-editing generally yields the best performance for the WMT24 official test set.

## 1 Introduction

Generative Large Language Models (LLMs) have demonstrated significant capabilities across various English-centric NLP tasks (Zhang et al., 2022; Touvron et al., 2023a,b). However, they often underperform in multilingual contexts, particularly with low-resource languages (Hendy et al., 2023; Stap and Araabi, 2023; Wang et al., 2023). To enhance the multilingual proficiency of LLMs, recent studies have explored several strategies, including vocabulary expansion (Lin et al., 2022; Liang et al., 2023; Yang et al., 2023), continual training on multilingual data (Le Scao et al., 2023; Dubey et al.,

2024; Xu et al., 2024a), and instruction tuning (Zhu et al., 2023; Alves et al., 2024; Stap et al., 2024). These approaches have collectively improved LLM performance on a variety of multilingual tasks, such as understanding (Lai et al., 2023), reasoning (Ponti et al., 2020; Shi et al., 2022), summarization (Hasan et al., 2021; Bhattacharjee et al., 2023), and machine translation (Kocmi et al., 2023).

Fine-tuning Large Language Models (FT-LLMs) with parallel data largely enhances translation capabilities, but such approach relies heavily on high-quality parallel data. For instance, prior research often uses development and test datasets like WMT and Flores (Alves et al., 2023; Xu et al., 2024a; Li et al., 2024) for the training, limiting the scalability to a broader range of languages. In this paper, we explore the feasibility of mining high-quality bi-texts from open-source corpora like OPUS. We utilize COMET (Rei et al., 2020), an automated Quality Estimation (QE) tool, to score sentences in the WMT-24 Constraint track. Unlike Peter et al. (2023), who found that selecting the highest quality sentences using COMET improves translation quality, our findings show that while this QE-based data filtering does not provide clear benefits for LLMs when refining fine-tuning datasets, it significantly enhances the performance of NMT systems when applied to filter training samples with low COMET scores.

Recent studies show that LLMs fine-tuned with MT data can rival state-of-the-art NMT models like NLLB (Costa-jussà et al., 2022). However, such comparisons may be unfair, as NMT models like NLLB typically support a broader range of languages. For example, ALMA-13b (Xu et al., 2024a) outperforms NLLB-54b (Costa-jussà et al., 2022) despite targeting only eight language pairs versus 200. Additionally, expanding languages in multilingual models often causes interference that degrades performance (Tan et al., 2024; Shaham et al., 2023). In this paper, we focus exclusively on

the English-to-Chinese translation direction[1], investigating how FT-LLMs compare to NMT models trained from scratch using the same parallel data source. Specifically, we use the full WMT-24 constraint track data to train an encoder-decoder NMT model, and we fine-tune LLaMA2-7B on a selected high-quality subset of up to 300K sentences. we found that, despite fine-tuned LLama2-7B being 17 times larger, it yields comparable COMET scores and worse scores for BLEU and ChrF.

While small NMT systems are resource-efficient in production, LLMs in practice, generate less literal translations (Vilar et al., 2023). In this paper, we integrate NMT and LLM systems by prompting LLMs to post-edit (PED) NMT outputs. Additionally, we implement a QE-guided PED system that selects the final outputs based on the higher QE score, as determined by COMET, between NMT and post-edited outputs. Our experiments show that the QE-guided PED system delivers the best performance on the WMT24 en-zh official test set, improving ChrF up to +3.7 over pure NMT outputs and +2.1 than direct translations by LLMs. Surprisingly, this approach brings negative performance gains on the Flores-devtest and Ntrex.

## 2  Data Preprocessing

In this section, we provide an overview of the data sources and the cleaning strategy. We use all the available data from the constrained track of the WMT-24 shared task for all three directions in which we participate, including English→Chinese, English→Japanese, and Japanese→Chinese. Following Wu et al. (2023), we perform a thorough preprocessing phase involving three key steps to enhance the data quality, as outlined below.

- Character-level Cleaning

  - Deescaping special characters in XML.
  - Removing non-printable characters.
  - Segmenting Chinese sentences with Jieba[2] and tokenizing Japanese data using KyTea (Neubig et al., 2011).

- Sentence-level Cleaning

  - Filtering out sentences longer than 256 tokens.

  - Eliminating sentences where over 75% of the words on both the source and target sides are identical.
  - Removing sentences with a source-to-target token ratio exceeding 3.0.
  - Eliminating duplicated sentences.

- Language-level Cleaning

  - Removing off-target sentences using the FastText language identification tool (Joulin et al., 2016).
  - Excluding sentences exhibiting one-to-many or many-to-one mappings, for example, a single source sentence having multiple different target sentences.

In specific, we use the Moses toolkit[3](Koehn et al., 2007) for all procedures in cleaning step 1 and use FastText (Joulin et al., 2016) for the language identification step. As shown in Table 1 (Cleaned), we removed 29%, 22%, and 45% of the data for en→zh, en→ja, and ja→zh directions.

| Directions | Raw | Cleaned | QE-filtered |
|---|---|---|---|
| en→zh | 55,346,004 | 39,354,051 | 22,606,804 |
| en→ja | 33,875,162 | 26,415,631 | 14,507,351 |
| ja→zh | 22,642,553 | 12,560,471 | 6,679,265 |

Table 1: Number of parallel sentences for three datasets.

## 3  Systems

### 3.1  NMT Systems

**MMT baseline**  In this section, we describe the backbone architecture and adjustments made to our baseline systems. We train a multilingual-Transformer-large (mT-large) model for all three en→zh, en→ja, ja→zh directions. The mT-large is a 12-layer Transformer (Vaswani et al., 2017) architecture with specific modifications, including pre-norm for both the encoder and decoder, and layer-norm for embedding. To enhance stability and performance, we tie the parameters of encoder embedding, decoder embedding, and decoder output. We also introduce dropout and attention dropout with a probability of 0.1, along with label smoothing at a rate of 0.1. In addition, to specify the translation directions, we prepend the source language tags in the source, and target language tags in the target side, e.g.: en2zh.

---

Similar to the approach described by Vaswani et al. (2017), we employ the Adam optimizer with a learning rate of 5e-4, implementing an inverse square root learning rate schedule with 4,000 warmup steps. We set the maximum number of tokens to 10,240, with gradient accumulation every 21 steps to facilitate large-batch training in Tang et al. (2021). We train all of our systems with 4 NVIDIA A6000 Gpus, and to expedite the training process, we conducted all experiments using half-precision training (FP16). Additionally, we save checkpoints every 2000 steps and implement early stopping based on perplexity, with a patience of 5 epochs.

**Quality-Estimation Filtering.** Due to data scarcity in the machine translation community, a large amount of Machine Translation data is mined from web-crawled data such as CCAligned (El-Kishky et al., 2020). Nonetheless, recent research found that there are many misaligned data exist in such web-crawled datasets, which impair performance when training models on it (Khayrallah and Koehn, 2018; Ranathunga et al., 2024). In addition, incorrect language and non-linguistic contents could affect the model in generating off-target or hallucinated outputs (Kreutzer et al., 2022). Similarly, recent studies on instruction fine-tuning of LLMs have shown that increasing data quality is more effective than data quantity (Du et al., 2023; Pan et al., 2024; Zhou et al., 2024), especially in inducing instruction-related capabilities (Xia et al., 2024). Additionally, Peter et al. (2023) shows that using QE metrics is not as effective at detecting translation noises like untranslated sentences, but is much better at identifying more fine-grained problems in the data, like small translation or grammatical errors.

Motivated by that, we investigate the feasibility of extracting high-quality parallel data using an automated Quality Estimation (QE) tool. We utilize the COMETKiwi model and apply this data-filtering phase to the cleaned data that we discussed in Section 2. Figure 1 presents the COMET score distributions for three directions. We found that for both English→Chinese and English→Japanese, the distributions are quite similar, that is, nearly half of the data falls into the poor quality range (0-80% Comet scores). For Japanese→Chinese, approximately half dataset ranges from 0% to 65% of COMET score. According to this observation, we filtered out parallel data that has smaller than

80% Comet scores for both English→Chinese and English→Japanese, and set the threshold at 65% for Japanese→Chinese. As a result, we show the number of parallel sentences after Quality Estimation filtering in Table 1.

**Directional Fine-tuning.** Lastly, to encourage the MMT model to gradually narrow down the data distribution to focus on task-specific data, we further fine-tune the MMT model on direction-specific data. Note that the direction-specific data, i.e., En → Zh, En → Ja, and Ja → Zh are the same data that included in the MMT baseline training data.

### 3.2 LLM Systems

We use LLaMA2-7B as the backbone because it is permitted for the constraint track of WMT24. We reuse the framework of ALMA (Xu et al., 2023) to conduct fine-tuning, however, we discard their first stage of monolingual continue training.

We set the training batch as 32 and accumulated 4-step gradients. The learning rate is set as 2e-5. The model was trained for one epoch using bf16 precision. The beam size is set as 5 for inference.

For the fine-tuning dataset, we further apply the quality estimation method described in Section 3.1 to filter out data with a QE score below a certain threshold. Then, we sample a certain number of bitext from the filter dataset. For example, in Table 4, the number of samples with a score above 89 is 53k, all of which are used for fine-tuning. Additionally, we sample data with scores higher than 87 at various levels, such as 53k, 100k, and 300k. We fine-tune LLaMA2 with different kinds of data to show the impact of data qualities.

### 3.3 NMT+LLM Systems

Previous studies have shown that leveraging Large Language Models (LLMs) to post-edit the outputs of supervised Neural Machine Translation (NMT) models can reduce translationese and enhance translation quality (Chen et al., 2023). This strategy has proven effective with LLMs such as ChatGPT (Chen et al., 2023), GPT-4 (Raunak et al., 2023), PaLM (Xu et al., 2024b), and LLaMA-2 (Ki and Carpuat, 2024). Specifically, post-editing utilizes LLMs either to refine the outputs of supervised NMT models or to perform "Self-Refinement" on their own outputs. Furthermore, Ki and Carpuat (2024) demonstrate that tuning LLMs with error-annotated translations can further enhance performance.

Figure 1: Comet score distributions for WMT-24 constraint training data on en→zh, en→ja, and ja→zh directions.

In this paper, we explore the effectiveness of Post-Editing (PED) in improving translation quality for the English-to-Chinese direction. We focus on a training-free PED approach due to computational constraints, utilizing pre-trained open LLMs to edit the outputs of our supervised NMT models. Given the limited Chinese capability of LLaMA2, we employ Tower-LLMs (Alves et al., 2024) (Tower-Instruct 7B and 13B), which have been continuously pre-trained on monolingual corpora including Chinese. Additionally, we implement a Quality Estimation-guided Post-Editing (QE-based PED) approach, where the NMT outputs and post-edited outputs are selected based on the higher QE score using COMETKiwi (wmt22-cometkiwi-da).

## 4 Experimental Setups

### 4.1 Systems

In this section, we briefly describe the systems we implemented. It is important to note that some of our implementations were focused only on the English-to-Chinese direction, specifically for FT-LLaMA2, Tower-Instruct, the PED system, and the QE-based PED system.

**mT-large.** A multilingual Transformer-large model trained in many-to-many directions using the "Cleaned" data (see Table 1 and Section 3.1 for details). It consists of 12 layers with 16 attention heads, $d = 1,024$, and $d_{ff} = 4,096$.

**mT-large + QE.** This model shares the same architecture and hyper-parameter settings as the *mT-large* model but is trained using the "QE-filtered" data outlined in Table 1.

**mT-large + QE + FT.** The *mT-large + QE* model was further fine-tuned on direction-specific data.

**FT-LLaMA2.** We use supervised fine-tuning to fine-tune LLaMA2. Detailed settings can be found in Section 3.2.

**Tower-Instruct.** We directly evaluate the performance of the Tower-Instruct models for comparison with our systems.

**Self-Refined PED.** We prompt the Tower-Instruct model to post-edit the translations they originally generated.

**PED system.** We prompt Tower-Instruct models to post-edit the outputs generated by our supervised NMT system (*mT-large + QE + FT*).

**QE-guided PED system.** We determined the final outputs by selecting between the NMT outputs and the post-edited outputs, based on the higher QE score as determined by COMETKiwi.

### 4.2 Data

For training, we utilize both the "Cleaned" and "QE-filtered" datasets, see details in section 2. For evaluation, we employ previous WMT validation and test sets as our validation set, and Flores, Ntrex as our test set.

### 4.3 Implementation and Evaluation

For our Neural Machine Translation (NMT) systems, we utilize the Fairseq toolkit (Ott et al., 2019) for both training and inference. For Large Language Model systems, we employ the Transformers toolkit for training and inference. To evaluate our models, we report detokenized SacreBLEU[4], ChrF++(Popović, 2017), and COMET (Rei et al., 2020) (wmt22-comet-da) scores.

---

[4]nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1

179

| ID | Methods | #Param | FLORES-Devtest | | | NTREX | | | WMT-24 Official | |
|----|---------|--------|------|------|-------|------|------|-------|------|------|
| | | | BLEU | ChrF | COMET | BLEU | ChrF | COMET | BLEU | ChrF |
| | **English→Japanese (NMT Systems Only)** | | | | | | | | | |
| ① | mT-large | 419M | 35.9 | 39.0 | 89.48 | 26.9 | 33.8 | 85.79 | 29.8 | 26.2 |
| ② | ① + QE | 419M | 36.6 | 39.8 | 90.00 | 27.3 | 34.5 | 86.95 | 34.2 | 29.6 |
| ③ | ② + FT | 419M | 37.1 | 40.3 | 90.24 | 28.3 | 35.1 | 87.18 | 34.7 | 30.1 |
| | **Japanese→Chinese (NMT Systems Only)** | | | | | | | | | |
| ④ | mT-large | 419M | 33.9 | 29.2 | 86.64 | 27.5 | 24.9 | 82.26 | 22.5 | 21.6 |
| ⑤ | ④ + QE | 419M | 34.0 | 29.1 | 87.04 | 27.6 | 25.0 | 82.77 | 22.7 | 22.0 |
| ⑥ | ⑤ + FT | 419M | 34.0 | 29.1 | 87.00 | 27.8 | 25.0 | 82.53 | 22.9 | 21.6 |

Table 2: Translation quality on NTREX, FLORES, and WMT test sets for the English→Japanese and Japanese→Chinese directions. 'FT' denotes directional Fine-Tuning, and 'QE' represents using QE-filtered training data. We use percentage for COMET scores.

| ID | Methods | #Param | FLORES-Devtest | | | NTREX | | | WMT-24 Official | |
|----|---------|--------|------|------|-------|------|------|-------|------|------|
| | | | BLEU | ChrF | COMET | BLEU | ChrF | COMET | BLEU | ChrF |
| | **NMT Systems** | | | | | | | | | |
| ① | mT-large | 419M | 42.2 | 35.0 | 84.94 | 33.3 | 28.7 | 79.20 | - | - |
| ② | ① + QE | 419M | 43.8 | 36.0 | 86.21 | 34.7 | 29.7 | 81.21 | - | - |
| ③ | ② + QE + FT | 419M | 43.9 | 36.2 | 86.12 | 35.0 | 29.7 | 80.95 | 33.5 | 31.6 |
| | **LLM Systems** | | | | | | | | | |
| ④ | FT-LLama2 | 7B | 34.6 | 31.2 | 86.60 | - | - | - | - | - |
| ⑤ | Tower-Instruct | 7B | 42.3 | 37.4 | 88.09 | 35.2 | 31.1 | 85.42 | 36.2 | 33.2 |
| ⑥ | Tower-Instruct | 13B | 43.2 | 38.0 | 88.12 | 36.2 | 32.0 | 85.36 | 38.5 | 35.3 |
| | **NMT + LLM Systems** | | | | | | | | | |
| ⑦ | Self-Refined PED (⑤) | 7B | 40.3 | 36.1 | 85.61 | 34.1 | 30.4 | 83.79 | 36.0 | 33.0 |
| ⑧ | PED (③ + ⑤) | 7.42B | 39.7 | 35.8 | 83.68 | 31.3 | 28.3 | 78.80 | 38.1 | 34.9 |
| ⑨ | QE-based PED (③ + ⑤) | 7.42B | 40.7 | 36.1 | 86.22 | 32.5 | 29.2 | 81.40 | 38.2 | 35.3 |

Table 3: Translation quality on NTREX, FLORES-200, and WMT-24 test sets for the English→Chinese direction. For WMT-24, we report BLEU and ChrF scores as returned by the OCELoT submission system.

## 5 Results and Analyses

In this section, we present the final results of our experiments and discuss the findings. Table 3 and 2 show the results of English→Chinese and the other two directions (en→ja and ja→zh) on Flores-devtest, Ntrex, and WMT24 official test sets.

### 5.1 Quality-Estimation Filtering improves NMT systems

Our key finding is that implementing Quality-Estimation (QE) Filtering effectively reduces low-quality data samples, leading to improved NMT system performance. Specifically, we observed BLEU score improvements of +4.4 and +0.2 for the English→Japanese and Japanese→Chinese directions, respectively, on the WMT24 official test sets. For the English→Chinese direction, we ob-

served BLEU gains of +1.6 on the Flores-devtest and +1.4 on the Ntrex test sets. Similar positive performance improvements were also noted across other metrics, such as ChrF and COMET. These results indicate that filtering training samples with low COMET scores enables our supervised NMT system to generate higher-quality translations.

### 5.2 Fine-tuned LLaMA2 and Data Quality

We conduct experiments on LLaMA2-7B in English to Chinese translation direction, where we collect 300K parallel samples from the training set, controlling the QE scores are all higher than 87. In Table 3, ④ shows the results. It is easy to see that the fine-tuned LLaMA2 results in the best COMET performance (86.60) on the Flores benchmark. However, the results on surface-level metrics, such as BLEU and ChrF, significantly lag

| Language | Data | BLEU | COMET |
|---|---|---|---|
| LLama2-7B | 10k (Cleaned) | 28.0 | 82.7 |
| LLama2-7B | 100k (Cleaned) | 35.7 | 85.6 |
| LLama2-7B | 53k (COMET > 87) | **36.1** | 86.1 |
| LLama2-7B | 53k (COMET > 89) | 33.5 | 84.3 |
| LLama2-7B | 100k (COMET > 87) | 35.5 | 85.7 |
| LLama2-7B | 300k (COMET > 87) | 34.6 | **86.6** |

Table 4: Evaluation results of fine-tuned LLama2-7B models for the English→Chinese direction on the Flores-devtest set. 'Cleaned' indicates random sampling from the 'Cleaned' training dataset, while 'COMET>x' refers to the sampling of data with COMET scores greater than x.

behind encoder-decoder-based NMT systems by 7.6 and 3.8 points, respectively.

We further control the fine-tuning data quality to show the impact. We select 10K and 100K samples from the cleaned dataset (See Table 1). To further improve the quality of parallel semantic alignment, we score all of the 39M cleaned training samples using COMET, and then we construct fine-tuning sets under the following settings:

- We selected all 53k samples with very high COMET scores, using a threshold of 89.

- We then lowered the score threshold to 87 and selected another 53k samples.

- We extend the number of samples with scores higher than 87 to 100k and 300k.

Table 4 shows the corresponding results after fine-tuning using datasets with different qualities. We observe that: 1) Simply extending the fine-tuning set from 10k to 100k largely improves the resulting performance. 2) However, no clear improvements can be observed when further raising the fine-tuning data QE quality. E.g., using 100k trivial samples (after data cleaning, QE score lower than 80) achieves comparable performance to that of using 100k samples with a QE score higher than 87. Additionally, fine-tuning with samples that have extremely high QE scores (COMET > 89) even resulted in a decline in translation quality compared to using 53k samples with relatively lower QE scores (COMET > 87). 3) Further extending the fine-tuning size from 100k to 300k yields no clear improvements.

Our experiments suggest that simply enhancing the quality of fine-tuning data for LLMs, at least when using COMET as the central measure of quality, is not a promising approach.

## 5.3 Post-Editing Enhances Translation Quality

As shown in Table 3, using the Tower-Instruct 7B LLM to post-edit the outputs of our strongest supervised NMT model (PED (③ + ⑤)) resulted in large improvements, with BLEU and ChrF gains of +4.6 and +3.3, respectively, over the NMT model alone on the WMT24 official test set. Notably, this post-editing approach also outperformed direct translation with Tower-Instruct 7B, achieving additional gains of +1.9 BLEU and +1.7 ChrF. In contrast, applying the Tower-Instruct model to post-edit its own generated translations (self-refined PED) resulted in negative improvements across all test sets. These findings suggest that integrating supervised NMT models with LLMs is a promising strategy for enhancing translation quality by leveraging the strengths of both systems.

Furthermore, Table 3 demonstrates that the QE-guided PED system (QE-based PED (③ + ⑤)) can further improve translation quality, as evidenced by the positive performance gains across the Flores-devtest, Ntrex, and WMT24 official test sets. In particular, the QE-guided PED system, utilizing Tower-Instruct 7B as the LLM backbone, achieved performance on par with Tower-Instruct 13B in the ChrF metric on the WMT24 official test set.

Despite the promising results on the WMT-24 Official test set, we found this Post-Editing approach delivered negative performance improvements on Flores and Ntrex sets (Table 3).

## 6 Conclusions

In this paper, we investigate three aspects of using LLMs for translation: 1) Comparison with Encoder-Decoder NMT Systems: directly fine-tuning LLaMA2 on a relatively small bitext dataset (100K) yields COMET results comparable to those of strong encoder-decoder NMT systems trained on over 50 million parallel sentence pairs. However, this approach significantly underperforms in surface-level metrics such as BLEU and ChrF. 2) Impact of Data Quality: properly filtering samples with low COMET scores largely improves encoder-decoder systems, however, no clear improvements can be observed for LLMs when further controlling the fine-tuning set with higher COMET scores. 3) Combining NMT Systems with LLMs: lastly, we show that combining NMT systems with LLMs via post-editing generally yields the best performance in our experiments.

## References

Duarte Alves, Nuno Guerreiro, João Alves, José Pombal, Ricardo Rei, José de Souza, Pierre Colombo, and André FT Martins. 2023. Steering large language models for machine translation with finetuning and in-context learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11127–11148.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Abhik Bhattacharjee, Tahmid Hasan, Wasi Ahmad, Yuan-Fang Li, Yong-Bin Kang, and Rifat Shahriyar. 2023. Crosssum: Beyond english-centric cross-lingual summarization for 1,500+ language pairs. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2541–2564.

Pinzhen Chen, Zhicheng Guo, Barry Haddow, and Kenneth Heafield. 2023. Iterative translation refinement with large language models. *arXiv preprint arXiv:2306.03856*.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Qianlong Du, Chengqing Zong, and Jiajun Zhang. 2023. Mods: Model-oriented data selection for instruction tuning. *arXiv preprint arXiv:2311.15653*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2020. Ccaligned: A massive collection of cross-lingual web-document pairs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5960–5969.

Tahmid Hasan, Abhik Bhattacharjee, Md Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M Sohel Rahman, and Rifat Shahriyar. 2021. Xl-sum: Large-scale multilingual abstractive summarization for 44 languages. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4693–4703.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hérve Jégou, and Tomas Mikolov. 2016. Fasttext. zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*.

Huda Khayrallah and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83.

Dayeon Ki and Marine Carpuat. 2024. Guiding large language models to post-edit machine translation with error annotations. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4253–4273.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2023. Findings of the 2023 conference on machine translation (wmt23): Llms are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions*, pages 177–180. Association for Computational Linguistics.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Auguste Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, et al. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.

Viet Lai, Chien Nguyen, Nghia Ngo, Thut Nguyn, Franck Dernoncourt, Ryan Rossi, and Thien Nguyen. 2023. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 318–327.

Teven Le Scao, Angela Fan, Christopher Akiki, El-lie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2023. Bloom: A 176b-parameter open-access multilingual language model.

Jiahuan Li, Hao Zhou, Shujian Huang, Shanbo Cheng, and Jiajun Chen. 2024. Eliciting the translation ability of large language models via multilingual finetuning with translation instructions. *Transactions of the Association for Computational Linguistics*, 12:576–592.

Davis Liang, Hila Gonen, Yuning Mao, Rui Hou, Na-man Goyal, Marjan Ghazvininejad, Luke Zettle-moyer, and Madian Khabsa. 2023. Xlm-v: Overcoming the vocabulary bottleneck in multilingual masked language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13142–13152.

Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Na-man Goyal, Shruti Bhosale, Jingfei Du, et al. 2022. Few-shot learning with multilingual generative language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9019–9052.

Graham Neubig, Yosuke Nakata, and Shinsuke Mori. 2011. Pointwise prediction for robust, adaptable japanese morphological analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 529–533.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53.

Xingyuan Pan, Luyang Huang, Liyan Kang, Zhicheng Liu, Yu Lu, and Shanbo Cheng. 2024. G-dig: Towards gradient-based diverse and high-quality instruction data selection for machine translation. *arXiv preprint arXiv:2405.12915*.

Jan-Thorsten Peter, David Vilar, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, and Markus Freitag. 2023. There's no data like better data: Using qe metrics for mt data filtering. In *Proceedings of the Eighth Conference on Machine Translation*, pages 561–577.

Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. Xcopa: A multilingual dataset for causal commonsense reasoning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2362–2376.

Maja Popović. 2017. chrf++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618.

Surangika Ranathunga, Nisansa De Silva, Velayuthan Menan, Aloka Fernando, and Charitha Rathnayake. 2024. Quality does matter: A detailed look at the quality and utility of web-mined parallel corpora. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 860–880.

Vikas Raunak, Amr Sharaf, Yiren Wang, Hany Awadalla, and Arul Menezes. 2023. Leveraging gpt-4 for automatic translation post-editing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12009–12024.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702.

Uri Shaham, Maha Elbayad, Vedanuj Goswami, Omer Levy, and Shruti Bhosale. 2023. Causes and cures for interference in multilingual translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15849–15863.

Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. 2022. Language models are multilingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057*.

David Stap and Ali Araabi. 2023. Chatgpt is not a good indigenous translator. In *Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP)*, pages 163–167.

David Stap, Eva Hasler, Bill Byrne, Christof Monz, and Ke Tran. 2024. The fine-tuning paradox: Boosting translation quality without sacrificing llm abilities. *arXiv preprint arXiv:2405.20089*.

Shaomu Tan, Di Wu, and Christof Monz. 2024. Neuron specialization: Leveraging intrinsic task modularity for multilingual machine translation. *arXiv preprint arXiv:2404.11201*.

Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Na-man Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2021. Multilingual translation from denoising pre-training. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3450–3466.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti

Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting palm for translation: Assessing strategies and performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15406–15427.

Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. Is chatgpt a good nlg evaluator? a preliminary study. In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11.

Di Wu, Shaomu Tan, David Stap, Ali Araabi, and Christof Monz. 2023. Uva-mt's participation in the wmt 2023 general translation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 175–180.

Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. Less: Selecting influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024a. A paradigm shift in machine translation: Boosting translation performance of large language models. In *The Twelfth International Conference on Learning Representations*.

Wenda Xu, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, Biao Zhang, Zhongtao Liu, William Yang Wang, Lei Li, and Markus Freitag. 2024b. Llmrefine: Pinpointing and refining large language models via fine-grained actionable feedback. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 1429–1445.

Wen Yang, Chong Li, Jiajun Zhang, and Chengqing Zong. 2023. Bigtranslate: Augmenting large language models with multilingual translation capability over 100 languages. *arXiv preprint arXiv:2305.18098*.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

Wenhao Zhu, Yunzhe Lv, Qingxiu Dong, Fei Yuan, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Extrapolating large language models to non-english by aligning languages. *arXiv preprint arXiv:2308.04948*.

# TOWER-V2:
# Unbabel-IST 2024 Submission for the General MT Shared Task

**Ricardo Rei**[*1] , **José Pombal**[*1,2,4] , **Nuno M. Guerreiro**[*1,2,4,5] , **João Alves**[*1] , **Pedro H. Martins**[*1]
**Patrick Fernandes**[2,3,4] , **Helena Wu**[1] , **Tania Vaz**[1] , **Duarte M. Alves**[2,4] , **Amin Farajian**[1]
**Sweta Agrawal**[2] , **Antonio Farinhas**[2,4] , **José G.C. de Souza**[1], **André F. T. Martins**[1,2,4]

[1]Unbabel
[2]Instituto de Telecomunicações  [3]Carnegie Mellon University
[4]Instituto Superior Técnico & Universidade de Lisboa (Lisbon ELLIS Unit)
[5]MICS, CentraleSupélec, Université Paris-Saclay

## Abstract

In this work, we present TOWER-V2, an improved iteration of the state-of-the-art open-weight TOWER models, and the backbone of our submission to the WMT24 General Translation shared task. TOWER-V2 introduces key improvements including expanded language coverage, enhanced data quality, and increased model capacity up to 70B parameters. Our final submission combines these advancements with quality-aware decoding strategies, selecting translations based on multiple translation quality signals. The resulting system demonstrates significant improvement over previous versions, outperforming closed commercial systems like GPT-4O, CLAUDE-SONNET-3.5, and DEEPL even at a smaller 7B scale.

## 1 Introduction

Large Language Models (LLMs) are making strides towards becoming the *de facto* solution for multilingual machine translation (MMT). Many works have shown that it is possible to adapt LLMs for translation and achieve state-of-the-art results (Zhang et al., 2023; Wei et al., 2023; Alves et al., 2023; Reinauer et al., 2023; Zhu et al., 2024).

One such example is our recent work on TOWER (Alves et al., 2024), which demonstrates that open NMT models like NLLB200 can be outperformed by adapting an LLM to translation. Specifically, we continue the pre-training of LLaMA-2 (Touvron et al., 2023) on both monolingual and parallel data, and fine-tune the resulting model on high-quality instructions covering several MT-related tasks. This approach requires much less parallel training data than traditional NMT and preserves the general capabilities of the LLM to respond to various prompts.

For the WMT24 General Translation task (Kocmi et al., 2024a), we enhance TOWER by significantly improving its training data, by

extending its language support from 10 to 15 languages — including low-resource ones like Icelandic —, and by scaling the underlying model to 70 billion parameters. Furthermore, because the WMT24 General Translation task focuses on paragraph-level translation instead of sentence-level, we also experiment with full-document translation and longer contexts, where TOWER originally struggled. These key improvements result in TOWER-V2 7B and 70B.

For our primary submission, we combine TOWER-V2 70B with Quality-Aware Decoding (QAD) strategies (Fernandes et al., 2022), such as Minimum Bayes Risk decoding (MBR) and Tuned Reranking (TRR). These techniques use reward models during inference to select the best candidate from a set of generated samples, enhancing the overall output quality.

We report our results, including the human evaluation and final submission, in Section 5. By outperforming strong commercial systems like GPT-4, CLAUDE-SONNET-3.5, and DEEPL across the board, TOWER-V2 — even at 7B parameters — challenges the belief that in MMT there must be a trade-off in performance between high- and low-resource language pairs (Fernandes et al., 2023).

Our contributions are:

- We show that expanding from 10 to 15 languages maintains the quality of translations for the initial 10 and significantly improves the newly added languages.

- We significantly improve the paragraph- and document-level translation capabilities of the previous TOWER.

- We demonstrate that scaling the model from 7 to 70B parameters yields improvements, indicating that increased capacity benefits not only general LLM abilities but also task-specific performance.

---

*Core Contributor. ✉ ai-research@unbabel.com

185

- We analyze the impact of QAD on larger models than those studied by Fernandes et al. (2022), showing that MBR decoding outperforms TRR according to both automatic metrics and human evaluation.

## 2 Overview of the Shared Task

The primary aim of the general machine translation shared task is to evaluate the ability of various models to translate across different domains, genres, and possibly modalities (e.g., speech). This year's shared task, compared to previous editions, emphasizes English→X (en→xx) and Non-English→Non-English (xx→yy) language pairs.[1]

The WMT24 test sets include source sentences from four domains: news articles, social media posts, speech (machine-generated transcripts), and literary texts. Additionally, all test sets from this year are focusing on the paragraph level rather than sentence-level.

Throughout this paper we will evaluate several of our models using both automatic and human evaluation; yet, for the shared task only primary submissions are evaluated, and final results are based solely on human evaluation using the ESA protocol (Kocmi et al., 2024c).

## 3 TOWER-V2: A New Translation LLM

We create TOWER-V2 by improving upon the original TOWER recipe: continued pre-training of a base model on a multilingual dataset with billions of tokens and subsequent supervised fine-tuning for translation-related tasks.

We focus on three key areas: 1) careful refinement of the training data; 2) expansion of language coverage to support all of the shared task's languages; 3) scaling up model capacity.

**Improving the training data.**   To enhance the general translation capabilities of TOWER, we mainly focus on improving the quality of its training data, be it for translation, post-translation, or general instructions.

For continued pre-training (CPT), we train on monolingual data from sources of superior quality, and apply more aggressive quality and length filters on the parallel data.

Regarding the supervised fine-tuning (SFT) phase, we use data created by humans — similarly to the previous version of TOWER— and introduce high-quality synthetic data. Human translations are sourced from well-known translation benchmarks. We go beyond simple sentence-level translation by transforming sentence-level to document-level data or into multi-parallel translation data (translating a single source sentence into multiple languages). When language variants are available, we include them in the training prompt (e.g. Chinese (simplified) vs Chinese (Taiwan)). All datasets were carefully filtered[2] and converted to instructions using a diverse set of templates.

**Improving post-translation data and general instructions.**   Data from tasks like APE, MQM evaluation, and translation ranking are carefully filtered using several quality signals. Similarly to XTOWER (Treviso et al., 2024), APE and MQM evaluation always expect the model to return a "translation correction," so we always ensure that the post-edition (PE) is deemed better than the original translation according to several metrics. For translation ranking, we choose only samples where there is significant alignment between human annotations and automatic metrics.

Like in the previous TOWER version, we aim to build a model that adheres to different prompts and can work as a general multilingual LLM. Thus, we include filtered and adapted multilingual general-purpose instruction data from publicly available high quality datasets such as AYA (Singh et al., 2024).

**Going from 10 to 15 Languages.**   We extend the language support of TOWER-V2 to Czech, Icelandic, Hindi, Ukrainian, and Japanese by adding training data of these languages to both CPT and SFT stages. For CPT, we add monolingual and parallel training data, increasing the total number of training tokens considerably. Aside from to-/from-English language pairs, we also include Czech-Ukrainian and Japanese-Chinese (and vice-versa) parallel data. In the SFT stage, we mostly add translation data for the new language pairs.

**More Paragraphs/Documents.**   In addition to the sentence-level parallel data we also add parallel documents to the CPT stage. For SFT, we sample high quality monolingual documents and per-

---

[1]The complete list of language pairs for this year's task includes: Czech→Ukrainian, Japanese→Chinese, and English→Chinese, Czech, German, Hindi, Icelandic, Japanese, Russian, Spanish (Latin America), Ukrainian

[2]We found low-quality translations even on datasets built by professionals.

| | WMT24 | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | en→de | en→es | en→cs | en→ru | en→uk | en→is | en→ja | en→zh | en→hi | cs→uk | ja→zh |
| **Baselines** | | | | | | | | | | | |
| NLLB-54B | 7.23 9 | 7.05 9 | 8.63 9 | 7.51 9 | 8.42 8 | 9.66 9 | 5.46 8 | 10.18 8 | 4.31 6 | 4.16 7 | 11.33 9 |
| GPT-4o | 1.41 6 | 1.57 7 | 1.48 6 | 1.39 6 | 1.42 6 | 2.31 7 | 1.04 5 | 1.65 5 | 1.19 4 | 0.94 4 | 3.42 6 |
| Claude-Sonnet-3.5 | 1.33 5 | 1.52 6 | 1.34 5 | 1.27 5 | 1.30 5 | 2.19 6 | 0.95 4 | 1.53 4 | 1.14 3 | 0.86 3 | 3.11 4 |
| DeepL | 1.81 8 | 2.10 9 | 1.71 7 | 2.21 8 | 1.44 6 | — | 3.95 7 | 2.22 7 | — | 1.40 5 | 7.36 9 |
| **Tower** | | | | | | | | | | | |
| Tower-v1 13B | 1.61 7 | 1.67 8 | — | 1.64 7 | — | — | — | 1.82 6 | — | — | — |
| Tower-v2 7B | 1.41 6 | 1.42 5 | 1.39 5 | 1.41 6 | 1.36 5 | 1.90 5 | 1.10 5 | 1.71 5 | 1.57 5 | 0.82 3 | 3.66 7 |
| Tower-v2 70B | 1.26 4 | 1.33 4 | 1.27 4 | 1.18 4 | 1.16 4 | 1.70 4 | 0.93 4 | 1.52 4 | 1.55 5 | 0.81 3 | 3.27 5 |
| **Tower + QAD** | | | | | | | | | | | |
| Tower-v2 70B+MBR | 0.93 2 | 0.96 2 | 0.83 2 | 0.80 2 | 0.72 2 | 1.20 2 | 0.71 2 | 1.20 2 | 0.97 2 | 0.61 2 | 2.64 2 |
| Tower-v2 70B+TRR | 1.07 3 | 1.05 3 | 0.96 3 | 0.91 3 | 0.87 3 | 1.27 3 | 0.82 3 | 1.27 3 | 1.07 3 | **0.59 1** | 2.88 3 |
| Tower-v2 70B 2-step | **0.91 1** | **0.94 1** | **0.77 1** | **0.76 1** | **0.70 1** | **1.14 1** | **0.68 1** | **1.17 1** | **0.94 1** | **0.57 1** | **2.59 1** |

Table 1: Translation quality (via METRICX-QE-XXL) on the WMT24 test set. TOWER-V2 with MBR/TRR ranks first across all language pairs. Even with Greedy decoding TOWER-V2-70B still ranks above other strong systems like CLAUDE-SONNET-3.5, GPT-4O and DEEPL except in en→hi and ja→zh where CLAUDE-SONNET-3.5 has similar scores.

formed full document translations using previous TOWER models while controlling for translation quality using COMETKIWI (Rei et al., 2022). At the end, we are left with more data for document-level than segment-level, further contributing to improved performance on paragraph- and document-level translation.

**Model suite.** TOWER-V2 now comes in two sizes: a 7B parameter model based on MISTRAL-7B (Jiang et al., 2023) and a larger 70B model based on LLAMA-3-70B (AI@Meta, 2024).

## 4  Quality-aware decoding with TOWER-V2

On LLM-based MT, translations are typically generated through lightweight decoding strategies such as greedy or nucleus sampling. Nevertheless, strategies informed by quality metrics such as Minimum Bayes Risk Decoding (MBR) and Tuned Reranking (TRR) consistently perform better compared to other methods (Fernandes et al., 2022; Freitag et al., 2022; Nowakowski et al., 2022; Farinhas et al., 2023). As such for our submission, we experiment with MBR and TRR. For both methods, we use a candidate pool of 100 samples and $\epsilon$-sampling (Freitag et al., 2023a) with $\epsilon = 0.02$, and COMET22 as the target objective. For TRR, we use

the WMT23 test set for tuning the weights[3]. The translation quality features used include: model log probabilities, COMET-QE-20, COMETKIWI22, COMETKIWI-XL, and xCOMET-QE-XL.

To leverage the strengths of both approaches, we also experiment with a second step of refinement. After obtaining translations from both MBR and TRR, we select the TRR translation only if all quality features (except the model log probabilities) agree that the TRR translation is better than the MBR translation; otherwise, we retain the MBR translation[4].

## 5  Experimental Setup

### 5.1  Evaluation Setup

During the development of TOWER-V2, we used WMT23 as our validation set. For our final analysis, we use WMT24 test set source sentences and report only QE metrics: COMETKIWI-XXL (Rei et al., 2023), METRICX-QE-XXL (Juraska et al., 2023), and xCOMET-QE-XXL (Guerreiro et al., 2023). Additionally, we add the official preliminary results to the Appendix which include METRICX (reference-based) (Kocmi et al., 2024b).

We use evaluation metrics to develop and op-

---

[3]We sample 5000 sentences from the WMT23 test set to train the weights more efficiently.

[4]According to both automatic and human evaluation (Table 2 and Table 3 respectively) results of MBR translations are generally better.

| Models | en→xx | | | xx→yy | | |
|---|---|---|---|---|---|---|
| | METRICX ↓ | xCOMET↑ | COMETKIWI ↑ | METRICX ↓ | xCOMET↑ | COMETKIWI ↑ |
| **Baselines** | | | | | | |
| NLLB-54B | 7.61 7 | 66.90 7 | 57.01 7 | 7.74 8 | 48.21 6 | 56.14 7 |
| GPT-4O | 1.50 6 | 83.74 6 | 77.04 5 | 2.18 5 | 70.44 2 | 76.19 4 |
| CLAUDE-SONNET-3.5 | 1.40 5 | 84.85 5 | 78.09 4 | 1.98 4 | 69.73 2 | 76.77 4 |
| DEEPL | — | — | — | 4.38 6 | 56.19 4 | 68.33 6 |
| **TOWER** | | | | | | |
| TOWER-V2 7B | 1.48 5 | 83.77 5 | 77.02 5 | 2.24 5 | 67.44 4 | 75.86 4 |
| TOWER-V2 70B | 1.32 4 | 84.87 4 | 78.29 4 | 2.04 4 | 69.20 3 | 76.70 4 |
| **TOWER + QAD** | | | | | | |
| TOWER-V2 70B+MBR | 0.92 2 | 88.78 2 | 81.39 3 | 1.62 2 | 69.88 2 | 78.28 2 |
| TOWER-V2 70B+TRR | 1.03 3 | 87.95 3 | 82.13 2 | 1.73 2 | **71.95 1** | 79.38 2 |
| TOWER-V2 70B 2-step | **0.89 1** | **89.25 1** | **82.54 1** | **1.58 1** | 70.85 2 | **79.69 1** |

Table 2: Translation quality aggregated by language pairs on the WMT24 test set (without testsuites). We omit DEEPL from the en→xx averages because it does not support two language pairs. All metrics are their XXL variant.

timize our models (e.g., using MBR and/or TRR during inference), with the exception of metrics of the METRICX family. Thus, to mitigate potential biases, we report METRICX-QE-XXL as our main evaluation metric and conduct human evaluation for English→German and English→Chinese. For the human evaluation, we use SQM quality levels with full document context. The annotators are in-house expert linguists familiar with evaluating MT outputs.

On Table 1, we report performance clusters based on statistically significant performance gaps at a 95% confidence threshold. On Table 2, we create per-language groups for systems with similar performance, following Freitag et al. (2023b), and obtain system-level rankings using a normalized Borda count (Colombo et al., 2022), which is defined as an average of the obtained clusters.

Regarding baselines, we report three commercial systems, GPT-4O, CLAUDE-SONNET-3.5, and DEEPL, along with an open-source NMT model, NLLB 54B. While little is known about the commercial systems, they show top performance on the WMT23. All models are evaluated in a 0-shot setting, unless stated otherwise.

## 5.2 Main Results

Table 1 shows our main results on English→X language pairs according to METRICX-QE-XXL (↓). Table 2 shows aggregated scores for English→X and X→Y according to different metrics. From Table 1, we observe that even the 7B model



Figure 1: Improvement in MT quality after adding new languages to TOWER-V2; measured in negative METRICX-XXL-QE so taller bars equate to better quality.

with greedy decoding outperforms, or is on par, with the best baseline, CLAUDE-SONNET-3.5, for English→X. Scaling to 70B brings consistent improvements across all language pairs, and both TRR and MBR decoding bring METRICX-QE-XXL further down. Our final submission (2-step) ranks first for all language pairs with statistical significance.

## 5.3 Impact of Adding 5 Languages

To evaluate the impact of adding 5 languages, we train two 7B models: one with the initial 10 languages of TOWER; another with the 10 languages

Figure 2: Win rates margin by length of the tokenized source of TOWER-V2-7B (squares) and TOWER-V2-70B (triangles) against an older iteration that was not trained on long-context translation training data. All language pairs of the WMT23 dataset that intersect with WMT24 are considered. We define a (sentence-level) win if the delta between two systems is superior to $1 \times 10^{-3}$ METRICX-XXL points

plus Hindi, Japanese, Ukrainian, Czech, and Icelandic. The data distribution for CPT remains unchanged, but we increase the number of training tokens of the second model to accommodate the additional languages. For SFT, we extend the dataset by incorporating human-translated data from several sources.

Figure 1 illustrates the absolute difference in 0-shot translation quality between the two models. As expected, the model with additional support performs considerably better on the new languages[5]. Perhaps more interestingly, its performance on the initially supported languages — which is already state-of-the-art (Table 1) — remains largely unchanged.

### 5.4 Beyond sentence by sentence translation

Figure 2 compares the new versions of TOWER-V2 (7B and 70B) with an older TOWER version that had yet to be trained on data specifically tailored to improve long-context translation. Not only do TOWER-V2 models vastly outperform the older version, but the quality gap widens as source length increases.

Further to this point, we created a paragraph-level version of the WMT23 dataset, by joining

---

[5]We note that the initial version of TOWER has ability to translate to other languages outside the supported ones, especially when given few-shot examples (Richburg and Carpuat, 2024) Still, their zero-shot performance is weak for languages like Hindi or Icelandic, which are less represented in the pre-training of the base models like LLaMA-2.

| Decoding | en→de | en→zh |
|---|---|---|
| **Batch 1** | | |
| Greedy | 85.43 | 84.11 |
| TRR | 87.16 | 85.55* |
| MBR | 88.50* | 85.47* |
| **Batch 2** | | |
| TRR | — | 68.55 |
| MBR | — | 72.76* |

Table 3: SQM quality evaluation for three different decoding methods using TOWER-V2 70B. Numbers marked with an asterisk (*) are statistically significant. For English→Chinese, since the results of the first batch were not significant, we conducted a second batch comparison between TRR and MBR.

segments of the same document into paragraphs with at most 4 sentences. Results in Table 4 show that our final models are considerably better at translating paragraphs than their older counterpart.

### 5.5 Putting all together into 70B parameters

The gains from scaling up the number of parameters are clear from Tables 1 and 2, where we show that TOWER-V2-70B consistently outperforms all baselines in all language pairs, except ja→zh. Coupling TOWER-V2-70B with QAD methods yields state-of-the-art results for all languages and metrics considered. Remarkably, Figure 2 shows that the 70B model considerably improves upon its 7B counterpart suggesting that the benefits of scaling up are particularly noticeable when translating longer sources.

### 5.6 Human Evaluation: Greedy vs TRR vs MBR

To validate our findings with automatic metrics, we conducted a small-scale human evaluation for English→German and English→Chinese (Table 3). In a first phase, linguists annotated 100 samples from TOWER-V2-70B with different decoding strategies on the WMT24 test. For both language pairs, annotators scored greedy decoding lower than the other two methods. While there was a noticeable quality difference between MBR and TRR for English→German, this distinction was not evident for English→Chinese, with both decoding strategies achieving similar results. Therefore, we conducted a second round of annotations for English→Chinese, comparing only TRR with MBR. This provided more concrete results that favored MBR outputs.

| Models | WMT23-Paragraphs | | | | | |
| | en→xx | | | xx→yy | | |
| | METRICX ↓ | COMET ↑ | CHRF ↑ | METRICX ↓ | COMET ↑ | CHRF ↑ |
|---|---|---|---|---|---|---|
| TOWER (older) | 5.14 | 79.11 | 50.93 | 6.99 | 75.45 | 53.29 |
| TOWER-V2-7B | 2.72 | 84.45 | 54.35 | 1.87 | 87.57 | 61.36 |
| TOWER-V2-70B | **2.40** | **84.87** | **55.06** | **1.72** | **87.75** | **62.29** |

Table 4: Performance of different TOWER versions on our paragraph-level version of WMT23 (measured by METRICX-XXL, COMET-22, and CHRF). TOWER (older) is a version prior to the interventions we ultimately made on the training data of TOWER-V2 to make it better at translating longer sources. These changes led to major improvements in paragraph-level translation for TOWER-V2-7B, which are further realized with TOWER-V2-70B.

## 5.7 Context-aware translation

| Models | en→xx | |
| | METRICX ↓ | XCOMET ↑ |
|---|---|---|
| TOWER-V2-70B 0-shot | 0.510 | 96.96 |
| TOWER-V2-70B 5-shot | 0.495 | 96.89 |
| | xx→en | |
| TOWER-V2-70B 0-shot | 1.051 | 94.84 |
| TOWER-V2-70B 5-shot | 0.766 | 95.54 |

Table 5: Translation quality of TOWER-V2-70B on the development set of the WMT24 Chat Shared Task. Using a prompt that incorporates conversational context (see Appendix A), the model provides high-quality translations, especially with examples (5-shot).

To evaluate TOWER-V2 in a different domain, we tested it on chat translation data. In this domain, the model translates a segment based on the context of previous conversation turns. Ignoring this context can result in subpar translations with pronoun mistakes and lexical inconsistencies (Läubli et al., 2018; Toral et al., 2018). Table 5 shows that TOWER-V2-70B excels at chat translation, even without specific training for this task. Using the prompt in Appendix A, which includes the conversation context, the model provides high-quality translations, especially when given domain-specific examples.

## 6 Conclusion

In this paper, we describe the joint submission from Unbabel and IST to the WMT24 General MT shared task. Our new model, TOWER-V2, significantly improves upon previous versions by expanding language coverage from 10 to 15 languages

and enhancing translation quality for longer paragraphs. Our largest model, with 70 billion parameters, combined with QAD strategies, achieved first place on the WMT24 test set according to both reference-free automatic evaluation, which we employed, and reference-based evaluation, as reported in the preliminary results from the WMT24 organizers (Kocmi et al., 2024b).

## Limitations

This paper highlights the key improvements in TOWER-V2 compared to previous versions and benchmarks it against other commercial state-of-the-art systems like GPT-4O, CLAUDE-SONNET-3.5, and DEEPL. However, our submission is "unconstrained and closed," meaning the information provided is not sufficient for full system replication. Furthermore, our comparisons primarily focus on translation quality and do not consider factors like inference speed, training budget, or model efficiency.

We also disclose the number of parameters in our models, from the 7B version to the final 70B version, to facilitate a clearer understanding of their scale. However, these comparisons with other systems do not account for differences in model parameters and other operational metrics.

## Acknowledgements

# References

AI@Meta. 2024. Llama 3 model card.

Duarte Alves, Nuno Guerreiro, João Alves, José Pombal, Ricardo Rei, José de Souza, Pierre Colombo, and Andre Martins. 2023. Steering large language models for machine translation with finetuning and in-context learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11127–11148, Singapore. Association for Computational Linguistics.

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks.

Pierre Colombo, Nathan Noiry, Ekhine Irurozki, and Stéphan Clémençon. 2022. What are the best systems? new perspectives on nlp benchmarking. In *Advances in Neural Information Processing Systems*.

António Farinhas, José de Souza, and Andre Martins. 2023. An empirical study of translation hypothesis ensembling with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11956–11970, Singapore. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Patrick Fernandes, Behrooz Ghorbani, Xavier Garcia, Markus Freitag, and Orhan Firat. 2023. Scaling laws for multilingual neural machine translation. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 10053–10071. PMLR.

Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023a. Epsilon sampling rocks: Investigating sampling strategies for minimum Bayes risk decoding for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9198–9209, Singapore. Association for Computational Linguistics.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023b. Results of wmt23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, Singapore. Association for Computational Linguistics.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024a. Findings of the 2024 conference on machine translation (WMT24). In *Proceedings of the Ninth Conference on Machine Translation*, Miami, Florida. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024b. Preliminary wmt24 ranking of general mt systems and llms.

Tom Kocmi, Vilém Zouhar, Eleftherios Avramidis, Roman Grundkiewicz, Marzena Karpinska, Maja Popović, Mrinmaya Sachan, and Mariya Shmatova. 2024c. Error span annotation: A balanced approach for human evaluation of machine translation.

Samuel Läubli, Rico Sennrich, and Martin Volk. 2018. Has machine translation achieved human parity? a case for document-level evaluation. *arXiv preprint arXiv:1808.07048*.

Artur Nowakowski, Gabriela Pałka, Kamil Guttmann, and Mikołaj Pokrywka. 2022. Adam Mickiewicz University at WMT 2022: NER-assisted and quality-aware neural machine translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 326–334, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, JosÃ© Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Raphael Reinauer, Patrick Simianer, Kaden Uhlig, Johannes E. M. Mosig, and Joern Wuebker. 2023. Neural machine translation models can learn to be few-shot learners.

Aquia Richburg and Marine Carpuat. 2024. How multilingual are large language models fine-tuned for translation?

Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F. Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura OMahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Souza Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergün, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Minh Chien, Sebastian Ruder, Surya Guthikonda, Emad A. Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. 2024. Aya dataset: An open-access collection for multilingual instruction tuning.

Antonio Toral, Sheila Castilho, Ke Hu, and Andy Way. 2018. Attaining the unattainable? reassessing claims of human parity in neural machine translation. *arXiv preprint arXiv:1808.10432*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Marcos Treviso, Nuno M. Guerreiro, Sweta Agrawal, Ricardo Rei, José Pombal, Tania Vaz, Helena Wu, Beatriz Silva, Daan van Stigt, and André F. T. Martins. 2024. xtower: A multilingual llm for explaining and correcting translation errors.

Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li, Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhiwei Cao, Binbin Xie, Tianxiang Hu, Shangjie Li, Binyuan Hui, Bowen Yu, Dayiheng Liu, Baosong Yang, Fei Huang, and Jun Xie. 2023. Polylm: An open source polyglot large language model.

Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhengrui Ma, Yan Zhou, Langlin Huang, Mengyu Bu, Shangtong Gui, Yunji Chen, Xilin Chen, and Yang Feng. 2023. Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2765–2781, Mexico City, Mexico. Association for Computational Linguistics.

## A  Appendix

### A.1  Metrics for QAD

The translation quality features used include: model log probabilities, COMET-QE-20[6], COMETKIWI22[7], COMETKIWI-XL[8], and XCOMET-QE-XL[9].

### A.2  Chat Translation Prompt

Given a source (SRC) to be translated from SRC_LANG to TGT_LANG, and previous turns

---

[6] Unbabel/wmt20-comet-qe-da
[7] Unbabel/wmt22-cometkiwi-da
[8] Unbabel/wmt23-cometkiwi-da-xl
[9] Unbabel/XCOMET-XL

in a conversation between two agents (TURN_i), the 0-shot prompt used was:

Context: <TURN_1>\n <TURN_2>...\n <TURN_k>.\n\nTranslate the <SRC_LANG>source text to <TGT_LANG>, given the context.\n<SRC_LANG>: <SRC>\n<TGT_LANG>:

When using five in-context examples, the prompt is repeated six times separated by two new lines; five times with a reference translation at the end, and one times exactly as written above.

### A.3 Further analysis on long-context translation

Compared to the first version of TOWER, the ability of TOWER-V2 to translate long sources has greatly improved. Whereas the translation quality of latter fell behind GPT-4 for longer sources, TOWER-V2-70B is superior across the board compared to the current best closed model for translation, CLAUDE-SONNET-3.5. In fact, the performance gap tends to widen as source length increases. TOWER-V2-7B is also competitive for the first 4 quantiles of length, but falls slightly behind on the last one.



Figure 3: Win rates margin by length of the tokenized source of TOWER-V2-7B (squares) and TOWER-V2-70B (triangles) against CLAUDE-SONNET-3.5. All language pairs of the WMT23 dataset that intersect with WMT24 are considered. We define a (sentence-level) win if the delta between two systems is superior to $1 \times 10^{-3}$ METRICX-XXL points

### A.4 Preliminary Results from Kocmi et al. (2024b)

See Tables 6 to 16 for the official automatic evaluation conducted by WMT 24 organizers. Our submission, Unbabel-Tower70B, ranks first on all language pairs and metrics.

## Czech-Ukrainian

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 0.9 | 0.719 | ✓ |
| Claude-3.5 § | 1.7 | 1.0 | 0.683 | ✓ |
| IOL-Research | 1.9 | 1.3 | 0.681 | ✓ |
| CommandR-plus § | 1.9 | 1.3 | 0.677 | ✓ |
| GPT-4 § | 2.0 | 1.4 | 0.677 | ✓ |
| Gemini-1.5-Pro | 2.0 | 1.2 | 0.668 | ✓ |
| ONLINE-W | 2.3 | 1.4 | 0.661 | ✓ |
| Mistral-Large § | 2.3 | 1.6 | 0.666 | |
| IKUN | 2.3 | 1.6 | 0.664 | ✓ |
| Aya23 | 2.5 | 1.9 | 0.665 | ✓ |
| TranssionMT | 2.6 | 1.5 | 0.648 | |
| ONLINE-B | 2.6 | 1.6 | 0.648 | |
| ONLINE-A | 2.6 | 1.5 | 0.647 | |
| Llama3-70B § | 2.6 | 2.0 | 0.661 | |
| ONLINE-G | 2.8 | 1.8 | 0.639 | |
| CUNI-Transformer | 3.0 | 2.0 | 0.639 | ✓ |
| IKUN-C | 3.0 | 2.4 | 0.648 | ✓ |
| Phi-3-Medium § | 9.1 | 6.5 | 0.425 | |
| BJFU-LPT † | 11.5 | 7.6 | 0.321 | |
| CycleL | 21.0 | 19.5 | 0.146 | |

Table 6: Preliminary WMT24 General MT automatic ranking for Czech-Ukrainian.

## English-Czech

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 1.8 | 0.732 | ✓ |
| Claude-3.5 § | 2.1 | 2.4 | 0.693 | ✓ |
| CUNI-MH | 2.1 | 2.3 | 0.690 | ✓ |
| CUNI-GA | 2.3 | 3.7 | 0.726 | ✓ |
| Gemini-1.5-Pro | 2.6 | 2.8 | 0.678 | ✓ |
| GPT-4 § | 2.6 | 2.9 | 0.682 | ✓ |
| IOL-Research | 2.8 | 3.0 | 0.676 | ✓ |
| ONLINE-W | 2.8 | 2.8 | 0.669 | ✓ |
| CommandR-plus § | 2.9 | 2.9 | 0.669 | ✓ |
| SCIR-MT | 3.2 | 3.3 | 0.664 | ✓ |
| TranssionMT | 3.5 | 3.5 | 0.655 | |
| ONLINE-A | 3.6 | 3.4 | 0.648 | |
| Mistral-Large § | 3.7 | 3.6 | 0.647 | |
| IKUN | 3.9 | 3.7 | 0.638 | ✓ |
| ONLINE-B | 4.0 | 3.9 | 0.640 | |
| Llama3-70B § | 4.1 | 4.0 | 0.640 | ✓ |
| Aya23 | 4.3 | 4.0 | 0.630 | ✓ |
| CUNI-DocTransformer | 4.4 | 4.0 | 0.621 | ✓ |
| IKUN-C | 4.7 | 4.3 | 0.618 | ✓ |
| CUNI-Transformer † | 4.7 | 4.3 | 0.614 | |
| ONLINE-G | 5.7 | 5.2 | 0.592 | |
| NVIDIA-NeMo † | 7.6 | 6.5 | 0.536 | |
| Phi-3-Medium § | 15.0 | 11.4 | 0.305 | |
| TSU-HITs | 19.5 | 16.6 | 0.235 | |
| CycleL2 | 24.2 | 19.5 | 0.077 | |
| CycleL | 27.0 | 22.5 | 0.031 | |

Table 7: Preliminary WMT24 General MT automatic ranking for English-Czech.

# English-German

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 1.1 | 0.723 | ✓ |
| Dubformer | 1.8 | 1.2 | 0.694 | ✓ |
| TranssionMT | 1.8 | 1.4 | 0.699 | ✓ |
| GPT-4 | 1.8 | 1.4 | 0.700 | ✓ |
| ONLINE-B | 1.8 | 1.4 | 0.698 | ✓ |
| Claude-3.5 | 1.9 | 1.4 | 0.695 | ✓ |
| CommandR-plus | 2.0 | 1.4 | 0.696 | ✓ |
| Mistral-Large | 2.0 | 1.5 | 0.694 | ✓ |
| Gemini-1.5-Pro | 2.2 | 1.5 | 0.688 | ✓ |
| ONLINE-W | 2.2 | 1.5 | 0.689 | |
| IOL-Research | 2.3 | 1.6 | 0.692 | ✓ |
| Llama3-70B § | 2.5 | 1.7 | 0.686 | ✓ |
| Aya23 | 2.7 | 1.8 | 0.680 | ✓ |
| IKUN | 3.0 | 1.8 | 0.668 | ✓ |
| ONLINE-A | 3.0 | 1.8 | 0.667 | |
| Phi-3-Medium § | 3.4 | 2.0 | 0.657 | |
| ONLINE-G | 3.5 | 2.1 | 0.662 | |
| IKUN-C | 3.8 | 2.0 | 0.641 | ✓ |
| CUNI-NL | 4.2 | 2.1 | 0.624 | |
| AIST-AIRC | 7.2 | 3.3 | 0.551 | |
| NVIDIA-NeMo † | 7.4 | 3.5 | 0.558 | |
| Occiglot | 8.2 | 3.8 | 0.539 | |
| MSLC | 11.9 | 4.4 | 0.390 | |
| TSU-HITs | 13.3 | 5.6 | 0.395 | |
| CycleL2 | 27.0 | 11.5 | 0.091 | |
| CycleL | 27.0 | 11.5 | 0.091 | |

Table 8: Preliminary WMT24 General MT automatic ranking for English-German.

## English-Spanish

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 1.9 | 0.745 | ✓ |
| GPT-4 | 1.9 | 2.5 | 0.712 | ✓ |
| Dubformer | 2.0 | 2.2 | 0.700 | ✓ |
| CommandR-plus | 2.1 | 2.6 | 0.706 | ✓ |
| Claude-3.5 | 2.1 | 2.6 | 0.705 | ✓ |
| Mistral-Large | 2.2 | 2.7 | 0.707 | ✓ |
| IOL-Research | 2.3 | 2.8 | 0.701 | ✓ |
| Gemini-1.5-Pro | 2.4 | 2.8 | 0.696 | ✓ |
| Llama3-70B § | 2.6 | 3.0 | 0.693 | ✓ |
| ONLINE-B | 2.7 | 3.1 | 0.690 | |
| ONLINE-W | 2.7 | 3.0 | 0.682 | |
| TranssionMT | 2.8 | 3.2 | 0.689 | |
| IKUN | 2.8 | 3.3 | 0.687 | ✓ |
| Phi-3-Medium § | 3.0 | 3.4 | 0.685 | |
| ONLINE-A | 3.0 | 3.3 | 0.676 | |
| Aya23 | 3.1 | 3.5 | 0.681 | |
| ONLINE-G | 3.2 | 3.6 | 0.674 | |
| IKUN-C | 3.4 | 3.5 | 0.666 | ✓ |
| NVIDIA-NeMo † | 4.5 | 4.4 | 0.631 | |
| Occiglot | 5.9 | 5.4 | 0.583 | |
| MSLC | 7.4 | 6.4 | 0.532 | ✓ |
| TSU-HITs | 16.3 | 14.2 | 0.289 | |
| CycleL | 24.0 | 20.9 | 0.072 | |

Table 9: Preliminary WMT24 General MT automatic ranking for English-Spanish.

## English-Hindi

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 3.1 | 0.657 | ✓ |
| Claude-3.5 § | 1.2 | 3.3 | 0.649 | ✓ |
| TranssionMT | 1.3 | 3.3 | 0.644 | ✓ |
| ONLINE-B | 1.4 | 3.3 | 0.641 | ✓ |
| Gemini-1.5-Pro § | 1.6 | 3.6 | 0.635 | ✓ |
| GPT-4 § | 2.1 | 4.5 | 0.628 | ✓ |
| IOL-Research | 2.1 | 4.3 | 0.622 | ✓ |
| Llama3-70B § | 2.1 | 4.6 | 0.630 | ✓ |
| CommandR-plus § | 2.3 | 4.4 | 0.612 | |
| Aya23 | 3.2 | 5.4 | 0.591 | ✓ |
| ONLINE-A | 3.5 | 6.2 | 0.590 | |
| ONLINE-G | 4.2 | 7.4 | 0.583 | |
| Mistral-Large § | 5.0 | 7.7 | 0.541 | |
| IKUN-C | 5.5 | 7.1 | 0.499 | ✓ |
| NVIDIA-NeMo † | 5.8 | 8.9 | 0.530 | |
| Phi-3-Medium § | 7.4 | 10.7 | 0.483 | |
| IKUN | 7.7 | 9.4 | 0.428 | |
| ONLINE-W | 15.3 | 20.9 | 0.296 | |
| CycleL | 20.0 | 23.4 | 0.083 | |

Table 10: Preliminary WMT24 General MT automatic ranking for English-Hindi.

# English-Icelandic

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 2.5 | 0.740 | ✓ |
| Claude-3.5 § | 2.3 | 3.6 | 0.697 | ✓ |
| Dubformer | 2.5 | 3.4 | 0.685 | ✓ |
| IKUN | 3.2 | 4.3 | 0.666 | ✓ |
| GPT-4 | 3.4 | 4.7 | 0.673 | ✓ |
| AMI | 3.7 | 4.9 | 0.663 | ✓ |
| IKUN-C | 3.7 | 4.9 | 0.657 | ✓ |
| TranssionMT | 4.2 | 5.5 | 0.653 | |
| ONLINE-B | 4.2 | 5.5 | 0.652 | |
| IOL-Research | 4.3 | 5.7 | 0.655 | ✓ |
| ONLINE-A | 5.5 | 6.4 | 0.603 | |
| Llama3-70B § | 6.7 | 8.0 | 0.586 | ✓ |
| ONLINE-G | 6.9 | 7.9 | 0.573 | |
| CommandR-plus § | 9.8 | 10.6 | 0.487 | |
| Mistral-Large § | 10.4 | 10.9 | 0.465 | |
| Aya23 § | 15.2 | 14.9 | 0.311 | |
| Phi-3-Medium § | 16.2 | 15.7 | 0.278 | |
| ONLINE-W | 18.1 | 19.5 | 0.296 | |
| TSU-HITs | 19.2 | 18.4 | 0.192 | |
| CycleL | 21.0 | 20.2 | 0.148 | |

Table 11: Preliminary WMT24 General MT automatic ranking for English-Icelandic.

## English-Japanese

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 2.0 | 0.762 | ✓ |
| ONLINE-B | 1.4 | 2.4 | 0.750 | ✓ |
| Claude-3.5 | 1.5 | 2.3 | 0.744 | ✓ |
| Gemini-1.5-Pro | 1.7 | 2.5 | 0.734 | ✓ |
| GPT-4 | 1.7 | 2.7 | 0.740 | ✓ |
| Team-J | 1.9 | 2.9 | 0.740 | ✓ |
| NTTSU | 1.9 | 2.6 | 0.731 | ✓ |
| CommandR-plus | 1.9 | 2.7 | 0.730 | ✓ |
| IOL-Research | 2.3 | 3.1 | 0.724 | ✓ |
| Aya23 | 2.3 | 3.1 | 0.719 | ✓ |
| Llama3-70B § | 2.6 | 3.5 | 0.714 | ✓ |
| DLUT-GTCOM | 2.6 | 3.0 | 0.697 | |
| Phi-3-Medium § | 2.8 | 3.6 | 0.709 | |
| ONLINE-W | 2.9 | 3.6 | 0.700 | |
| Mistral-Large § | 2.9 | 3.8 | 0.707 | |
| ONLINE-A | 3.0 | 3.6 | 0.699 | |
| IKUN | 3.1 | 3.7 | 0.696 | |
| IKUN-C | 3.9 | 4.3 | 0.669 | ✓ |
| ONLINE-G | 6.4 | 6.6 | 0.599 | |
| AIST-AIRC | 6.6 | 6.5 | 0.583 | |
| UvA-MT | 6.7 | 6.7 | 0.589 | |
| NVIDIA-NeMo † | 6.9 | 6.9 | 0.582 | |
| CycleL | 24.0 | 22.4 | 0.101 | |

Table 12: Preliminary WMT24 General MT automatic ranking for English-Japanese.

# English-Russian

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 2.4 | 0.742 | ✓ |
| Dubformer | 1.9 | 2.8 | 0.701 | ✓ |
| Yandex | 1.9 | 2.9 | 0.705 | ✓ |
| Claude-3.5 | 2.0 | 3.0 | 0.706 | ✓ |
| ONLINE-G | 2.2 | 3.3 | 0.706 | ✓ |
| GPT-4 | 2.3 | 3.4 | 0.703 | ✓ |
| Gemini-1.5-Pro | 2.3 | 3.2 | 0.697 | ✓ |
| CommandR-plus § | 2.4 | 3.4 | 0.693 | ✓ |
| ONLINE-W | 2.6 | 3.5 | 0.688 | |
| IOL-Research | 2.6 | 3.7 | 0.694 | ✓ |
| Mistral-Large § | 2.7 | 3.7 | 0.692 | |
| Llama3-70B § | 3.1 | 4.1 | 0.681 | ✓ |
| ONLINE-B | 3.1 | 3.9 | 0.673 | |
| TranssionMT | 3.1 | 3.9 | 0.673 | |
| IKUN | 3.2 | 4.1 | 0.675 | ✓ |
| Aya23 | 3.3 | 4.2 | 0.669 | ✓ |
| ONLINE-A | 3.4 | 4.1 | 0.663 | |
| Phi-3-Medium § | 3.9 | 4.7 | 0.654 | |
| IKUN-C | 3.9 | 4.7 | 0.649 | ✓ |
| CUNI-DS | 5.9 | 6.2 | 0.584 | |
| NVIDIA-NeMo † | 7.2 | 7.3 | 0.549 | |
| TSU-HITs | 10.8 | 9.8 | 0.421 | |
| CycleL | 24.3 | 22.2 | 0.062 | |
| CycleL2 | 25.0 | 22.4 | 0.027 | |

Table 13: Preliminary WMT24 General MT automatic ranking for English-Russian.

# English-Ukrainian

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 2.2 | 0.732 | ✓ |
| Dubformer | 1.8 | 2.7 | 0.691 | ✓ |
| Claude-3.5 | 2.0 | 3.0 | 0.693 | ✓ |
| ONLINE-W | 2.1 | 2.8 | 0.679 | ✓ |
| Gemini-1.5-Pro | 2.2 | 3.0 | 0.677 | ✓ |
| CommandR-plus § | 2.3 | 3.2 | 0.678 | ✓ |
| GPT-4 | 2.3 | 3.3 | 0.682 | ✓ |
| ONLINE-G | 2.3 | 3.1 | 0.670 | |
| IOL-Research | 2.4 | 3.4 | 0.675 | ✓ |
| Mistral-Large § | 2.4 | 3.4 | 0.675 | |
| IKUN | 2.8 | 3.7 | 0.661 | ✓ |
| ONLINE-B | 3.1 | 3.9 | 0.646 | |
| TranssionMT | 3.1 | 4.0 | 0.646 | |
| Llama3-70B § | 3.2 | 4.2 | 0.647 | |
| Aya23 | 3.3 | 4.2 | 0.642 | |
| ONLINE-A | 3.3 | 4.1 | 0.634 | |
| IKUN-C | 3.9 | 4.7 | 0.622 | ✓ |
| NVIDIA-NeMo † | 6.2 | 7.0 | 0.537 | |
| Phi-3-Medium § | 11.1 | 11.3 | 0.339 | |
| CycleL | 21.0 | 22.4 | 0.037 | |

Table 14: Preliminary WMT24 General MT automatic ranking for English-Ukrainian.

# English-Chinese

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 2.3 | 0.726 | ✓ |
| Claude-3.5 | 1.7 | 3.0 | 0.703 | ✓ |
| ONLINE-B | 1.7 | 2.9 | 0.697 | ✓ |
| IOL-Research | 1.8 | 3.1 | 0.700 | ✓ |
| Gemini-1.5-Pro | 1.8 | 3.1 | 0.698 | ✓ |
| GPT-4 | 2.0 | 3.3 | 0.693 | ✓ |
| CommandR-plus | 2.2 | 3.3 | 0.681 | ✓ |
| ONLINE-W | 2.2 | 3.2 | 0.677 | |
| HW-TSC | 2.3 | 3.4 | 0.675 | ✓ |
| Mistral-Large § | 2.8 | 4.0 | 0.665 | |
| Llama3-70B § | 2.8 | 3.9 | 0.662 | ✓ |
| Aya23 | 3.0 | 4.1 | 0.655 | ✓ |
| IKUN | 3.1 | 4.0 | 0.646 | ✓ |
| Phi-3-Medium § | 3.1 | 4.2 | 0.648 | |
| ONLINE-A | 3.3 | 4.1 | 0.636 | |
| IKUN-C | 3.5 | 4.2 | 0.624 | ✓ |
| UvA-MT | 4.3 | 5.2 | 0.607 | |
| ONLINE-G | 4.8 | 5.5 | 0.588 | |
| NVIDIA-NeMo † | 7.3 | 7.6 | 0.494 | |
| CycleL | 20.1 | 20.1 | 0.086 | |
| CycleL2 | 22.0 | 22.1 | 0.030 | |

Table 15: Preliminary WMT24 General MT automatic ranking for English-Chinese.

# Japanese-Chinese

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 3.2 | 0.622 | ✓ |
| Claude-3.5 | 1.7 | 3.5 | 0.603 | ✓ |
| Gemini-1.5-Pro | 1.9 | 3.5 | 0.595 | ✓ |
| DLUT-GTCOM | 2.0 | 3.3 | 0.586 | ✓ |
| GPT-4 | 2.1 | 3.8 | 0.597 | ✓ |
| IOL-Research | 2.2 | 3.9 | 0.593 | ✓ |
| CommandR-plus | 2.8 | 4.1 | 0.576 | ✓ |
| Team-J | 2.8 | 4.0 | 0.570 | ✓ |
| Llama3-70B § | 3.1 | 4.7 | 0.578 | ✓ |
| Mistral-Large § | 3.5 | 4.9 | 0.568 | |
| Aya23 | 3.7 | 5.0 | 0.563 | ✓ |
| NTTSU | 3.7 | 5.3 | 0.566 | ✓ |
| Phi-3-Medium § | 4.0 | 5.1 | 0.552 | |
| IKUN | 4.4 | 5.4 | 0.544 | ✓ |
| ONLINE-B | 5.2 | 5.5 | 0.518 | |
| UvA-MT | 5.2 | 6.3 | 0.534 | |
| ONLINE-W | 5.3 | 6.0 | 0.522 | |
| IKUN-C | 5.5 | 6.2 | 0.519 | ✓ |
| ONLINE-A | 6.8 | 6.8 | 0.484 | |
| MSLC | 8.9 | 8.8 | 0.450 | |
| ONLINE-G | 10.3 | 9.6 | 0.413 | |
| CycleL | 23.0 | 21.5 | 0.202 | |

Table 16: Preliminary WMT24 General MT automatic ranking for Japanese-Chinese.

# TSU HITS's Submissions to the WMT 2024 General Machine Translation Shared Task

**Vladimir Mynka♠, Nikolay Mikhaylovskiy♠◇**
♠Higher IT School, Tomsk State University, Tomsk, Russia
◇NTR Labs, Moscow, Russia
vladimirmynka34821@gmail.com, nickm@ntrlab.com

## Abstract

This paper describes the TSU HITS team's submission system for the WMT'24 general translation task. We focused on exploring the capabilities of discrete diffusion models for the English-to-{Russian, German, Czech, Spanish} translation tasks in the constrained track. Our submission system consists of a set of discrete diffusion models for each language pair. The main advance is using a separate length regression model to determine the length of the output sequence more precisely.

## 1 Introduction

This report gives an overview of TSU HITS submissions in the WMT 2024 general machine translation tasks. We focused on exploring the capabilities of discrete diffusion models for the English-to-{Russian, German, Czech, Spanish} translation tasks in the constrained track. Our main contributions are

1. the use of regression-based output length prediction model
2. the use of the input length as a key feature for the output length prediction

The report is organized as follows. In the Section 2, we provide a general description of the discrete diffusion approach to machine translation, as it is not yet very widespread. In the Section 3, we describe the experimental setting and training processes. Section 4 discusses the results.

## 2 Discrete Diffusion Approach to Machine Translation

### 2.1 Diffusion: Preliminaries

Diffusion approaches (Sohl-Dickstein et al., 2015 , Ho et al, 2020) to generating objects (for example images) include forward (data to noise) and reverse (noise to data) diffusion processes. In the forward process, a small amount of noise is gradually added to the data. In the classical direct diffusion process, the original object $x_0$ is repeatedly and additively perturbed by a small Gaussian random noise, and in a fixed number of steps $T$ goes into state $x_T$ with a normal distribution (and thus is converted to noise):

$$q(x_t|x_{t-1}) = \mathcal{N}\left(x_t;\ x_{t-1}\sqrt{1-\beta_t},\beta_t\right), \quad (1)$$

where $\forall\, t = \overline{1..T}\ \beta_t \in (0;1]$ are the hyperparameters that regulate the diffusion rate.

During the reverse diffusion process, the machine learning model step by step reconstructs the object's states from $x_T$ to $x_0$, and this denoising restores an object from the original distribution:

$$p_\theta(x_{t-1}|x_t) \sim \mathcal{N}\left(x_{t-1};\ \mu_\theta(x_t,t), \sigma_\theta(x_t,t)\right), (2)$$

where $\theta$ are the model's trainable weights.

Texts in typical representations do not have the property of continuity and are a sequence of tokens with discrete values that do not have an order relation and correspond to the categorical data type. Thus, we follow the path of adapting the diffusion processes to categorical data - such approaches are called discrete diffusion.

### 2.2 Discrete Diffusion for Text Generation

Diffusion models with discrete state spaces were first introduced by Sohl-Dickstein et al. (2015), who considered a diffusion process over binary

Figure 1: Overview of the system

random variables. Hoogeboom et al. (2021) extended the model class to categorical random variables with transition matrices characterized by uniform transition probabilities. We follow Austin et al. (2021) to define a discrete diffusion model for texts.

Namely, we consider each text token $x_t$ to be a discrete random variable with $K$ categories. For text data, $K = |V|$ is the size of the vocabulary. (He et al., 2023). The forward transition probabilities can be represented by matrices: $[Q_t]_{ij} = q(x_t = j | x_{t-1} = i)$. The process of adding noise can then be written as

$$q(x_t | x_{t-1}) = Cat(x_t; p = x_{t-1} Q_t) \quad (3)$$

where $Cat(\cdot)$ is a category distribution (Austin et al., 2021).

## 2.3 Masked Language Models and Discrete Diffusion

He et al. (2023) noted the relationship between the discrete diffusion process and the task of pretraining of masked language modeling (MLM) encoder models. Namely, they suggested incorporating an absorbing state, e.g., [MASK] for BERT, into the Markov process of diffusion. In particular, each token in the sequence either stays the same or transitions to [MASK] with some probability. Formally, each entry of the transition matrix at step $t$ is as follows,

$$[Q_t]_{ij} = \begin{cases} 1 & if\ i = j = [M] \\ \beta_t & if\ j = [M], i \neq [M] \quad (4) \\ 1 - \beta_t & if\ i = j \neq [M] \end{cases}$$

where [M] is short for [MASK].

Such a Markov process converges to a stationary distribution $q(x_T)$ that places all the probability mass on the sequence with all [MASK] tokens.

The most common transformer (Vasvani et al., 2017) models pre-trained for the MLM task are models from the BERT family (Devlin et al, 2019). He et al. (2023) suggested DiffusionBERT that uses a pretrained BERT model as an encoder due to the similarity of the tasks. The length of the output sequence of the DiffusionBERT model is fixed and is set to different values depending on the problem solved.

## 2.4 Discrete Diffusion for Translation

Reid et al. (2023) suggested a diffusion model using Levenstein operations for machine translation. They have tested the model on WMT14 EN-DE dataset. It is unclear from the paper how do the authors determine the target length of the output sequence.

Zheng et al. (2023) suggest a reparameterized discrete diffusion (RDM) approach to text generation, and report results for the machine translation task on the IWSLT14 DE-EN, WMT14 EN-DE and WMT16 EN-RO datasets. To determine the translation length, the authors of RDM trained a separate model similar to the one of Ghazvininejad et al. (2020). They pose the problem of determining the length of the output sequence as a classification problem, selecting $k$ best options out of $N$ possible, where $N$ is the maximum text length that the model used can process. Similarly to Gao et al. (2024), several options are selected and the best one is chosen based on the metrics of the overall text quality.

Ye et al. (2023) explore the possibilities of increasing applicability domain of discrete diffusion approaches, while considering an approach similar to DiffusionBERT, except that instead of the BERT encoder, the authors use the RoBERTa model (Liu et al., 2019). The quality of machine translation is assessed on the IWSLT14 DE-EN and WMT14 EN-DE data sets, using the same quality metrics and the same idea for determining the length as in the RDM approach.

## 3  System Overview

### 3.1  General Translation Process

The general translation process is presented on Figure 1. Our system consists of a discrete diffusion model and an output length prediction model.

On each diffusion step, a concatenation of the source text and output is used as the input to the generative model, but the absorbing tokens are distributed only within the output part. We do not use any special separation tokens, but just use the prompt "{Source Language}: {Source Text} \n {Target Language}: $\{x_t\}$".

Since we use XLM-RoBERTa's (Conneau et al, 2020) positional embedding model as an encoder and are forced to fit the input sequence of the model into 512 tokens, we apply punctuation splitting of the source texts, limiting the maximum size of the source text to 200 tokens, and then glue the results back. We also do not use the extended context to improve translation; this is left for the future work.

We take a fixed number of the diffusion steps $T$ equal to 50. Tokens that were unmasked in the previous steps are likely to be replaced with subsequent ones, just like in DiffusionBERT (He et al., 2023). The standard argmax approach is used as a sampling method. We do not use temperature and do not limit the number of tokens to choose from.

### 3.2  Generative Model

We largely follow Ye et al. (2023) and use XLM-RoBERTa (Conneau et al, 2020) family pre-trained model that includes a multilayer transformer encoder and a single-layer MLM head.

We fine-tune both the encoder and the head for discrete diffusion text generation that differs from MLM mainly by the percentage of the masked tokens. We use the cross-entropy weighted relative to the diffusion step $t$ loss proposed by Zheng et al. (2023):

$$L_t = -\lambda_{t-1} \sum_i^N y_i \log p_i \qquad (5)$$

| Generative Model | |
| --- | --- |
| Architecture | XLM-RoBERTa-Large |
| Optimizer | AdamW($\beta_1 = 0.9, \beta_2 = 0.98$) |
| Weights decay | 0.01 |
| Learning Rate Schedule | Cosine |
| Max learning rate | 5E-05 |
| Batch size | 16 |
| Accumulation step | 8 |
| Steps | 30000 |
| Warmup ratio | 0.01 |
| Loss | (Section 3.2) |
| Number format | FP16 |
| **Length Model** | |
| Hidden size | 1024 |
| Optimizer | AdamW($\beta_1 = 0.9, \beta_2 = 0.999$) |
| Learning Rate Schedule | OneCycleLR (Smith et al, 2017), two phases |
| Max learning rate | 7E-07 |
| Batch size | 8 / 16 |
| Steps | 30000 |
| Embedding calculation | Mean pooling |
| Activation | ELU |
| Loss function | MSE |
| Number format | FP16 |

Table 1: Hyperparameters of the models

where $y_i$ is the true probability (0 or 1) of token with index $i$ in model dictionary, $p_i$ is the predicted probability, $N$ is the size of the dictionary, $\lambda_{t-1}$ is the parameter that depends on the percentage of the masked tokens at the steps $t$ and $t-1$.

Following Chang et al. (2022), we use the cosine noise schedule:

$$\beta_t = \cos(\frac{\pi t}{2T}) \qquad (6)$$

### 3.3  Length Predictor

Our length predictor also consists of an encoder and a task-specific head. Although our length prediction model is based on the same XLM-RoBERTa, physically these two models are completely separate. We tried not to fine-tune the encoder for the length problem and to use the standard XLM-RoBERTa, but we got worse metrics on the test data.

We use a regression predictor of the output length, unlike other works that use classifiers with the number of categories equal to the length of the context, for example, 512 tokens. Our regression head is a two-layer perceptron with ELU-activation. Standard MSE loss is used when the length predictor is trained.

| | #tokens | #model parameters |
|---|---|---|
| EN-DE | 68,333 | 444,158,849 |
| EN-RU | 91,932 | 492,511,151 |
| EN-ES | 31,380 | 368,444,201 |
| EN-CS | 65,514 | 438,382,718 |

Table 2: Numbers of tokens and model parameters after pruning the tokenizer

The main improvement in length prediction is because of the use of the input length. There is a fairly strong relationship between the length of the text in the source language and the length of its translation, which, in general, is almost linear. We suggest taking this into account when the target is defined. Our model predicts the ratio of the input and output lengths, normalized by the average ratio for the training set. We employ standard mean pooling to convert the matrix of token embeddings obtained from the encoder into a common embedding of text, which will be used as features for the length head.

### 3.4 Training Data

The WikiMatrix dataset (Schwenk et al., 2021) was used as a train dataset for EN-DE, EN-RU, EN-CS language pairs; Neulab-TedTalks (Tiedemann, 2012) was used for EN-ES. The training sets were trimmed to 480 thousand examples when training the generation model and to 240 thousand when training a length prediction model.

### 3.5 Pruning the tokenizer

Due to the computational limitations we reduce the token set of our models for each pair of languages to the minimum required (all the other tokens are replaced with [UNK]). The effect of reduction on the number of model parameters is demonstrated in Table 2. According to our observations, it increases the quality of models when tested on validation datasets for the selected language pair, but may degrade the quality of general translation when tested on complex examples.

Pruning the tokenizer was made before trimming the training sets to keep as much tokens as possible.

### 4 Results

The official automatic scores of our system on the test data are presented in the Table 3. The gap

| | AutoRank | MetricX | CometKiwi |
|---|---|---|---|
| EN-DE | 13.3 | 5.6 | 0.395 |
| EN-RU | 10.8 | 9.8 | 0.421 |
| EN-ES | 16.3 | 14.2 | 0.289 |
| EN-CS | 19.5 | 16.6 | 0.235 |

Table 3: System official scores

between our results and the leading system is significant.

### 4.1 Model size

We used XLM-Roberta-Large with 561 million parameters as the main model for generating translation, while other systems participating in the competition this and last years had tens of billions of parameters. This makes our model largely uncompetitive. Unfortunately, today there are no pretrained open-weight encoder models comparable to leading open-weight decoder models in terms of parameters number and pretrain token count.

### 4.2 Quantity and quality of training data

Due to technical limitations, we used only a small part of the translation datasets provided, no more than 480 thousand examples for each language pair. Increasing the training set and better cleaning should significantly improve the quality, especially when using a larger pretrained model.

### Acknowledgments

### References

Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. 2021. Structured denoising diffusion models in discrete state-spaces. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual,* pages 17981–17993.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Zhujin Gao, Junliang Guo, Xu Tan, Yongxin Zhu, Fang Zhang, Jiang Bian, and Linli Xu. 2024. Empowering Diffusion Models on the Embedding Space for Text Generation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4664–4683, Mexico City, Mexico. Association for Computational Linguistics.

Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. 2019. Mask-Predict: Parallel Decoding of Conditional Masked Language Models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6112–6121, Hong Kong, China. Association for Computational Linguistics.

Zhengfu He, Tianxiang Sun, Qiong Tang, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. 2023. DiffusionBERT: Improving Generative Masked Language Models with Diffusion Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4521–4534, Toronto, Canada. Association for Computational Linguistics.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. 2021. Argmax flows and multinomial diffusion: Towards non-autoregressive language models. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems, NeurIPS 2021, December 6-14, 2021, virtual,* pages 12454-12465.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv preprint 1907.11692

Machel Reid, Vincent J. Hellendoorn, and Graham Neubig. 2023. Diffuser: Discrete diffusion via edit-based reconstruction. In *Proceedings of The Eleventh International Conference on Learning Representations,* https://openreview.net/forum?id=nG9RF9z1yy3

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021. WikiMatrix: Mining 135M Parallel Sentences in 1620 Language Pairs from Wikipedia. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume,* pages 1351–1361, Online. Association for Computational Linguistics.

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning, 2015,* pages 2256–2265.

Jörg Tiedemann. 2012. Parallel Data, Tools and Interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems, NeurIPS 2017,* pages 5998–6008

Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu and William T. Freeman. "MaskGIT: Masked Generative Image Transformer." 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022): 11305-11315.

Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Quanquan Gu. 2023. Diffusion language models can perform many tasks with scaling and instruction-finetuning. arXiv preprint arXiv:2308.12219.

Lin Zheng, Jianbo Yuan, Lei Yu, and Lingpeng Kong. 2023. A reparameterized discrete diffusion model for text generation. arXiv preprint arXiv:2302.05737.

Leslie N. Smith and Nicholay Topin. 2017. Super-Convergence: Very Fast Training of Residual Networks Using Large Learning Rates. arXiv preprint arXiv:1708.07120.

# Document-level Translation with LLM Reranking: Team-J at WMT 2024 General Translation Task

[*]**Keito Kudo** [1,2], [*]**Hiroyuki Deguchi** [3], [*]**Makoto Morishita** [4,1], [*]**Ryo Fujii** [4], [*]**Takumi Ito** [1,5],
[*]**Shintaro Ozaki** [3], [*]**Koki Natsumi** [3], **Kai Sato** [1], **Kazuki Yano** [1], **Ryosuke Takahashi** [1],
**Subaru Kimura** [1], **Tomomasa Hara** [1], **Yusuke Sakai** [3], **Jun Suzuki** [1,2]
[1]Tohoku University   [2]RIKEN   [3]NAIST   [4]Future Corporation   [5]Langsmith Inc.

## Abstract

We participated in the constrained track for English-Japanese and Japanese-Chinese translations at the WMT 2024 General Machine Translation Task. Our approach was to generate a large number of sentence-level translation candidates and select the most probable translation using minimum Bayes risk (MBR) decoding and document-level large language model (LLM) re-ranking. We first generated hundreds of translation candidates from multiple translation models and retained the top 30 candidates using MBR decoding. In addition, we continually pre-trained LLMs on the target language corpora to leverage document-level information. We utilized LLMs to select the most probable sentence sequentially in context from the beginning of the document.

## 1 Introduction

This paper details Team-J's system submission for the WMT 2024 Shared Task: General Machine Translation. We participated in the English-Japanese (En→Ja) and Japanese-Chinese (Ja→Zh) translation tasks under the constrained track.

As with last year's competition, the use of publicly available pre-trained models and metrics evaluated in the WMT Metrics shared tasks, such as COMET (Rei et al., 2020), was permitted. Following the Kudo et al.'s (2023) system, we employed multiple machine translation (MT) models to generate numerous candidate sentences for each source text. We then applied minimum Bayes risk (MBR) decoding (Fernandes et al., 2022) using the COMET metric to select the optimal translations.

Additionally, contrary to the previous years, the use of large language models (LLMs) was also permitted this year. Our primary objective was to use these LLMs to achieve consistent document-level machine translation. Specifically, we aimed

to develop models based on LLMs and also implemented a reranking system. Figure 1 provides an overview of our system. The following sections describe its components in detail.

## 2 Dataset Construction

In this section, we describe the training data, the process of synthetic data generation, and the data cleaning methodologies.

### 2.1 Provided Data

Since we participated in the constrained track, we solely used the data officially provided by the organizer.

**Bitext data.** We used all the provided bitext data. For English to Japanese translation, we used JParaCrawl v3.0 (Morishita et al., 2022a), News Commentary v18, Wiki Titles v3, WikiMatrix (Schwenk et al., 2021), Japanese-English Subtitle Corpus (JESC) (Pryzant et al., 2018), The Kyoto Free Translation Task (KFTT) Corpus (Neubig, 2011), and TED Talks (Cettolo et al., 2012). For Japanese to Chinese translation, we used JParaCrawl Chinese (Nagata et al., 2024), News Commentary v18, Linguatools Wiki Titles, WikiMatrix, OPUS, and Neulab TED Talks (Tiedemann, 2012).

**Monolingual data.** We also used the following provided monolingual data for Japanese and Chinese: News Crawl, News Commentary, Leipzig Corpora (Goldhahn et al., 2012), Common Crawl (Buck et al., 2014), and Extended Common Crawl (Conneau et al., 2020; Wenzek et al., 2020). For the continual pre-training of the language models, we only used the Common Crawl and Extended Common Crawl due to the limited availability of document-level data beyond these two datasets.

**Development data.** We used NTREX-128 (Federmann et al., 2022), Flores-200 (Team et al., 2022;

---

[*]: Equal contributions.

Figure 1: System overview

Goyal et al., 2022; Guzmán et al., 2019) and the past WMT test sets as development data. These datasets were also employed to fine-tune the models.

## 2.2 Synthetic Data

We constructed synthetic data to augment the training dataset. We used the synthetic data created by Kudo et al. (2023) for the En→Ja task, and newly created data for the Ja→Zh task as follows. For preprocessing, we tokenized the bitext (Section 2.1) into truecased[1] subwords using a unigram language model with Sentencepiece (Kudo and Richardson, 2018), with "byte_fallback", and "split_digits" options enabled following Touvron et al. (2023); Dubey et al. (2024); Kudo et al. (2023). After that, we created a back translation model (Sennrich et al., 2016), which we call an initial translation model using the training configurations in Table 7 (Appendix C) and trained it on the bitext. Then, we translated the Chinese monolingual data (Section 2.1) with a beam size of 10 and a length penalty of 1.0.

## 2.3 Data Cleaning

We conducted data cleaning on the corpus. Specifically, we applied several rules to clean and filter out noisy sequences using HojiChar (Shinzato, 2023). HojiChar is a text preprocessing tool that mainly supports monolingual corpus in Japanese

and English, with typical filters preinstalled. We first extended HojiChar to make it work with parallel corpus and implemented a variety of rules with careful investigation of the provided data. Table 1 shows the list of data cleaning methods we applied on the bitext and monolingual data. Table 2 shows the amount of data after filtering.

The following provides a detailed explanation of the cleaning rules that were mainly implemented using tools other than HojiChar.

**Character count-based filtering.** We qualitatively examined the Common Crawl and Extended Common Crawl datasets. Our analysis revealed that shorter sequences tend to be noisy. Therefore, we discarded sequences that were less than or equal to 200 characters for Japanese and 100 characters for Chinese, respectively (see (26) in Table 1). This threshold also helps us retain document-level data that is suitable for the continual pre-training of LLMs. To efficiently filter out shorter sequences, we used the `awk` command.

**Toxic content cleaning.** Qualitative analysis of the Common Crawl data revealed a significant amount of low-quality toxic contents, such as adult material, are included in the corpus. To address this, we applied a toxic content filter to exclude such samples from our training data ((9) in Table 1). For the Japanese data, we used filters originally implemented in HojiChar.[2] For the Chinese corpus, we defined a list of toxic words based on

---

[1] https://github.com/moses-smt/mosesdecoder/blob/master/scripts/recaser/truecase.perl

[2] `DiscardAdultContentJa` in HojiChar.

| Filter & Cleaner | Ja | Zh | En-Ja | Ja-Zh |
|---|---|---|---|---|
| (1) Discard content having identical source and target | | | ✓ | ✓ |
| (2) Discard content with invalid unicode characters | ✓ | ✓ | ✓ | ✓ |
| (3) Remove non-printable unicode characters | ✓ | ✓ | ✓ | ✓ |
| (4) Apply NFKC normalization | ✓ | ✓ | ✓ | ✓ |
| (5) Normalize space-like characters to half-width spaces | ✓ | ✓ | ✓ | ✓ |
| (6) Restore escaped HTML symbols | ✓ | ✓ | ✓ | ✓ |
| (7) Discard content like progress bars | ✓ | ✓ | ✓ | ✓ |
| (8) Discard content having many square brackets | ✓ | ✓ | ✓ | ✓ |
| (9) Discard content containing keywords for porn contents | ✓ | ✓ | | |
| (10) Discard content containing keywords for online bulletin boards | ✓ | ✓ | | |
| (11) Discard content containing part of sequences like word lists | ✓ | ✓ | ✓ | ✓ |
| (12) Discard content containing having many punctuations | ✓ | ✓ | ✓ | ✓ |
| (13) Discard content containing having many numbers | ✓ | ✓ | ✓ | ✓ |
| (14) Reduce repeated space and punctuation characters | ✓ | ✓ | ✓ | ✓ |
| (15) Discard content having many same consecutive characters | ✓ | ✓ | ✓ | ✓ |
| (16) Discard content having many same consecutive N-grams | ✓ | ✓ | ✓ | ✓ |
| (17) Discard content having less punctuations | ✓ | ✓ | | |
| (18) Discard content having no punctuations in a sliding window of specified length | ✓ | ✓ | | |
| (19) Discard content having low compression ratio with zlib compression | ✓ | ✓ | | |
| (20) Discard content not in expected languages | ✓ | ✓ | ✓ | ✓ |
| (21) Remove ellipsis symbols | ✓ | ✓ | ✓ | ✓ |
| (22) Remove open bracket end symbols at the end of the sentence | ✓ | ✓ | ✓ | ✓ |
| (23) Remove parentheses with no content inside | ✓ | ✓ | ✓ | ✓ |
| (24) Remove Unicode control characters | ✓ | ✓ | ✓ | ✓ |
| (25) Remove content starts with "&" | ✓ | ✓ | ✓ | ✓ |
| (26) Discard too short content | ✓ | ✓ | | |
| (27) Convert traditional Chinese to simplified Chinese | | | | ✓ |
| (28) Exact deduplication | ✓ | ✓ | ✓ | ✓ |
| (29) Fuzzy deduplication | ✓ | ✓ | | |
| (30) Discard too long content | | | ✓ | ✓ |
| (31) Discard content having too large source/target token ratio | | | ✓ | ✓ |
| (32) Discard content having too large token/char ratio | | | ✓ | ✓ |
| (33) Discard semantically irrelevant translations | | | ✓ | ✓ |

Table 1: List of data cleaning rules.

those used for the ChineseWebText (Chen et al., 2023) dataset.

**Compression rate-based cleaning.** We used a cleaning method based on the compression rate to remove non-textual data ((19) in Table 1).[3] Samples with a high compression rate typically contain excessive repetitions, while those with a low compression rate often consist of random strings. Specifically, we calculated the compression rate for each sample and removed those that did not fall within a specified range.

**Language detection.** To ensure the collection of data in the target language, we used language detection (20) in Table 1. Simple heuristic language detection methods are implemented in Hojichar, such as a method that checks for the presence of *hiragana* or *katakana*. Alongside these simple methods, we also used FastText-based language detection (Joulin et al., 2017b,a).

**Conversion of traditional Chinese to simplified Chinese.** We converted Chinese data written in traditional characters to simplified characters to augment the bitext data ((27) in Table 1). We used OpenCC[4] for these conversions.

**Deduplication.** Duplicate data in training sets can negatively impact the performance of language models (Lee et al., 2022). To mitigate this, we performed exact deduplication using the `sort` command ((28) in Table 1) and fuzzy deduplication using MinHash (Broder, 1997) ((29) in Table 1). We used the `text-dedup` tool (Mou et al., 2023) for implementation.

**Bitext similarity cleaning.** We performed cleaning based on bitext similarity using LaBSE (Feng et al., 2022) to filter out semantically irrelevant pairs ((33) in Table 1). We set the lenient threshold of 0.5 for bitext and more strict threshold of 0.7 to synthetic data.

|                  | # samples | # tokens |
|------------------|-----------|----------|
| **LLMs**         |           |          |
| Monolingual Ja   | 88.4M     | 35.8B    |
| Monolingual Zh   | 137.4M    | 29.9B    |
| Parallel En-Ja   | 29.8M     | 4.0B     |
| Parallel Ja-Zh   | 3.8M      | 506.3M   |
| **Encoder-Decoder** |        |          |
| Synthetic En-Ja  | 587M      | 12.9B    |
| Synthetic Ja-Zh  | 291M      | 10.3B    |
| Parallel En-Ja   | 28.2M     | 730.0M   |
| Parallel Ja-Zh   | 6.3M      | 163.6M   |

Table 2: The amount of training data used for LLMs and Encoder-Decoder MT models. The token count for LLMs is based on the tokenizer of Mistral-7B, and the count for Encoder-Decoder MT models is based on the subwords on the target side.

## 3 Primary Translation Models

We developed translation models using two architectures: Encoder-Decoder and Decoder-only (LLMs).

### 3.1 Encoder-Decoder MT Models

For En→Ja, we used the existing translation models created by Morishita et al. (2022b); Kudo et al. (2023). For Ja→Zh, we newly constructed translation models through pre-training and fine-tuning.

**Pre-training.** We trained the pre-training model using the pre-training configuration in Table 7 (Appendix C). For the training data, we used the bitext (Section 2.1) and the synthetic data (Section 2.2) after applying data cleaning (Section 2.3). We performed upsampling to achieve a $1 : 4.7$ ratio between the bitext and the synthetic data. Moreover, we applied the tagged back-translation technique (Caswell et al., 2019), adding a special token `<BT>` at the beginning of the source sentences in the synthetic data and storing this tag in the vocabulary dictionary.

**Fine-tuning.** After pre-training, we conducted fine-tuning using the development data (Section 2.1) with the fine-tuning configuration in Table 7 (Appendix C).

### 3.2 LLM-based MT Models

We used the Llama2-13B (Touvron et al., 2023) and Mistral-7B (Jiang et al., 2023), which are permitted for use in the constrained track. These LLMs were used only for the En→Ja direction and not for the Ja→Zh direction. For Mistral-7B, we also prepared a variant with an expanded vocabulary to improve

its Japanese generation capability. For more details on vocabulary expansion, please refer to Section B.

**Continual pre-training.** Although the datasets used for training Llama2 and Mistral are not publicly disclosed, it is generally believed that they are predominantly in English. Consequently, continual pre-training has been conducted to enhance performance on Japanese tasks (Fujii et al., 2024a; Okazaki et al., 2024). This approach has been reported to improve English-Japanese translation performance. To further boost Japanese language capability, we also performed continual pre-training using the cleaned monolingual corpus detailed in Section 2.3. The training configurations are shown in Table 8, 9, and 10.

**Supervised fine-tuning** After continual pre-training, we conducted supervised fine-tuning for the translation task. In this phase, we used the cleaned bitext corpus and development data described in Section 2. Initially, we fine-tuned the model using the bitext corpus, followed by additional fine-tuning with the development data which is relatively clean. To prepare for the Stepwise MBR-Enhanced LLM decoding detailed in Section 4.2, we used all combinations of the first $n$ sentences from each document as training samples for the development data fine-tuning. Figure 2 shows the prompt template, and Table 8, 9, and 10 shows hyperparameters used in the training process.

**Preference learning.** To align the translation results with human preferences, we conducted preference learning for Mistral-7B. [5] We used Contrastive Preference Optimization (CPO) (Xu et al., 2024) as the preference learning algorithm. In preliminary experiments, we also tried Direct Preference Optimization (DPO) (Rafailov et al., 2023) as an alternative to CPO. However, despite the decrease in loss during training, we observed that the DPO often resulted in output collapse (complete loss of input-output correspondence) during decoding. Therefore, we selected CPO as our preference learning.

Let $L_{\mathrm{NLL}}(\pi_\theta)$ and $L_{\mathrm{pref}}(\pi_\theta)$ be the negative log-likelihood of $\pi_\theta$ and preference of output given by

---

[5] Due to computational resource limitations, we applied LoRA fine-tuning (Hu et al., 2022).

次の英語を日本人のネイティブのように日本語に翻訳してください。 原文：{src} 訳文：{tgt}

Figure 2: The general prompt for supervised fine-tuning. {src} denotes the source sentence. {tgt} denotes the target sentence.

$\pi_\theta$, respectively, that is:

$$L_{\text{NLL}}(\pi_\theta) = -\mathbb{E}_{(s,r)\sim\mathcal{D}} \left[\log \pi_\theta(r \mid s)\right]$$
$$L_{\text{pref}}(\pi_\theta) = -\mathbb{E}_{(s,r,y_r)\sim\mathcal{D}} \left[\log \sigma(\beta d)\right] \quad , \quad (1)$$
$$d = \log \pi_\theta(r \mid s) - \log \pi_\theta(\hat{y} \mid s)$$

where $\sigma$ is the Sigmoid function. Then, CPO minimizes the following objective function during training:

$$\min_\theta \left[L_{\text{pref}}(\pi_\theta) + \alpha L_{\text{NLL}}(\pi_\theta)\right] \quad . \quad (2)$$

Here, $\mathcal{D} = \left\{\left(s^{(i)}, r^{(i)}, \hat{y}^{(i)}\right)\right\}_{i=1}^{N}$ represents the dataset. $\pi_\theta$ denotes a parameterized policy, and $\alpha$ and $\beta$ are hyperparameters. We used the development data for training in preference learning. In this context, $s$ corresponds to the source text from the development data, $r$ to the reference text from the development data, and $\hat{y}$ to the output of the model before preference learning.

To prevent the model output from collapsing, we introduced a minor modification to the CPO objective function. Specifically, we implemented a warm-up phase to reduce the impact of the preference learning loss at the beginning of training. This approach is formulated as follows:

$$\min_\theta \left[\min\left(1, \frac{i}{i_{\text{w}}}\right) L_{\text{pref}}(\pi_\theta) + \alpha L_{\text{NLL}}(\pi_\theta)\right] . \quad (3)$$

Here, $i$ represents the number of training steps, and $i_{\text{w}}$ denotes the number of warm-up steps for the preference learning loss.

## 4 Decoding

This year's test set consists of segments with multiple sentences in context. Since most bitext corpora are at the sentence level, translating larger segments in one shot is not preferable. Thus, we initially divided each segment in the test set into individual sentences using spaCy (Honnibal et al., 2020).[6] In case the resulting split was overly short, we combined texts from its adjacent splits.

| | hypotheses | pseudo-references | |
|---|---|---|---|
| | | top-p sampling | epsilon sampling |
| En→Ja | 1272.15 | 3288.5 | 3421.99 |
| Ja→Zh | 261.84 | 884.11 | 3108 |

Table 3: The average number of hypotheses and pseudo references for each source sentence generated by the Encoder-Decoder MT models. Note that due to errors during decoding, the number of hypotheses and pseudo-references generated for a single source sentence varies.

### 4.1 MBR Decoding

We apply minimum Bayes risk (MBR) decoding (Eikema and Aziz, 2020) to select high-quality translations from the set of hypotheses generated by the multiple translation models using MBRS (Deguchi et al., 2024). Let $\mathcal{Y}$ be the output space of translation models. We use the Monte Carlo method to estimate the expected utility (Eikema and Aziz, 2022), as follows:

$$y^{\text{MBR}} = \underset{h\in H}{\arg\max} \; \mathbb{E}_{\hat{r}\in\hat{R}} \left[u(h, \hat{r})\right],$$
$$= \underset{h\in H}{\arg\max} \; \frac{1}{|\hat{R}|} \sum_{\hat{r}\in\hat{R}} u(h, \hat{r}), \quad (4)$$

where $y^{\text{MBR}}$ is the selected translation by MBR decoding, $H \subseteq \mathcal{Y}$ is the hypotheses set, and $\hat{R}$ is the multiset (a.k.a bag) of translation samples[7], called "pseudo-references". $u: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ is the utility function that returns scores of the translation quality of the hypothesis under the given pseudo-references, which is formally defined as $h \succeq h' \iff u(h, \hat{r}) \geq u(h', \hat{r})$ where $\succeq$ denotes the preference relation. We employ COMET-22[8] (Rei et al., 2020, 2022) for the utility function $u$. Therefore, the MBR decoding using COMET-22 is formulated as follows:

$$y^{\text{MBR}} = \underset{h\in\mathcal{H}}{\arg\max} \; \frac{1}{|\hat{R}|} \sum_{\hat{r}\in\hat{R}} \texttt{COMET-22}(s, h, \hat{r}). \quad (5)$$

Note that COMET-22 also takes the source sentence $s$ as input.

---

[6]We used "en_core_web_lg" model for English and "ja_core_news_lg" model for Japanese.

[7]The support set is a subset of the output space, i.e., $\text{Supp}(\hat{R}) \subseteq \mathcal{Y}$

[8]https://huggingface.co/Unbabel/wmt22-comet-da

In our system, we select the 30-best translations using MBR decoding instead of selecting the 1-best translation as shown in Equation 4 to determine the final decision using another algorithm than MBR decoding. In other words, MBR decoding is used to prune translation hypotheses. We generate hypotheses for each source sentence using an ensemble of Encoder-Decoder MT models with beam search decoding. In addition, we prepare two types of pseudo-references by decoding with top-p sampling ($p = 0.9$) and epsilon sampling (Freitag et al., 2023) ($\epsilon = 0.02$). The number of hypotheses and pseudo-references used in MBR decoding is presented in Table 3.

### 4.2 Stepwise MBR-Enhanced LLM Decoding

---

**Algorithm 1:** Stepwise MBR-Enhanced LLM decoding

**Input:** $D_{\mathrm{src}} = \{s_0, s_1, \ldots, s_n\}$
**Output:** $D_{\mathrm{hyp}} = \{h_0, h_1, \ldots, h_n\}$
1 $D_{tgt} \leftarrow \{\}$;
2 $S_{hist} \leftarrow \{\}$;
3 **for** $i \leftarrow 0$ **to** $n$ **do**
   // Generate candidates for $s_i$
4   $H \leftarrow \mathrm{LLMs_{MT}}(s_i, S_{hist}, D_{tgt})$;
5   $h_i \leftarrow \mathrm{MBR}(H, H)$;
6   $D_{tgt} \leftarrow D_{\mathrm{hyp}} \cup \{h_i\}$;
7   $S_{hist} \leftarrow S_{hist} \cup \{s_i\}$;
8 **return** $D_{\mathrm{hyp}}$

---

During our preliminary experiments with fine-tuned LLMs, we observed frequent issues where some sentences were skipped during decoding. This led to discrepancies in the number of sentences between the source and the translated output. Additionally, we observed samples where the same token was generated repeatedly. To address these issues, we propose a decoding method called `Stepwise MBR-Enhanced LLM Decoding` (Algorithm 4.2). This method translates documents sentence by sentence, considering the overall document context (see Figure 3). This approach resolves the issue of mismatched sentence counts between the source and hypothesis. Furthermore, we applied MBR decoding to achieve high-quality sentence-level translation without repeated tokens or other errors (line 5 of Algorithm 4.2). We used the outputs of four LLMs for this method. Specifically, we used four LLMs

with different settings: Mistral-7B with and without vocab expansion and with and without preference learning.

## 5  LLM Reranking

As mentioned in Section 3 and Section 4, primary translation models decode at the sentence level. To improve the overall document-level consistency of the translation results, we performed reranking using LLMs. We used the top 30 highest-scoring hypotheses from MBR decoding as the candidate pool and reranked them based on context-aware scoring. Specifically, we used the LLMs fine-tuned for the translation task described in Section 3.2 to calculate the likelihood of each hypothesis with context information. We repeated this process to select hypotheses with the highest likelihood scores, resulting in the final translation output. The details are described in Algorithm 2. In our system, we use supervised fine-tuned Mistral-7B as the reranker, and we set the beam size to $b = 2$.

---

**Algorithm 2:** LLM Reranking Algorithm

**Input:** $D_{\mathrm{src}} = \{s_0, s_1, \ldots, s_n\}$
**Input:** $D_{\mathrm{hyps}} = \{H_0, H_1, \ldots, H_m\}$
**Input:** $b$: Beam size
**Output:** $D_{\mathrm{hyp}} = \{h_0, h_1, \ldots, h_n\}$
1 $\mathcal{C}_{\mathrm{beam}} \leftarrow \{(\emptyset, -\infty)\}$;
2 $P \leftarrow \emptyset$;
3 **for** $H \in D_{\mathrm{hyps}}$ **do**
4   **for** $(\mathbf{c}, \_) \in \mathcal{C}_{beam}$ **do**
5     **for** $h \in H$ **do**
6       $p_h \leftarrow \mathrm{LLM_{MT}}(D_{\mathrm{src}}, \mathbf{c} \cup \{h\})$;
7       $P \leftarrow P \cup \{(\mathbf{c}, h, p_h)\}$;
8   $\mathcal{T}_b \leftarrow \mathrm{Top}_b(P, \text{ with respect to } p_h)$;
9   $\mathcal{C}_{\mathrm{beam}} \leftarrow \{(\mathbf{c} \cup \{h\}, p_h) \mid (\mathbf{c}, h, p_h) \in \mathcal{T}_b\}$;
10 $(\mathbf{c}^*, p_c^*) \leftarrow \arg\max_{(\mathbf{c}, p_c) \in \mathcal{C}_{\mathrm{beam}}} p_c$;
11 $D_{\mathrm{hyp}} \leftarrow \mathbf{c}^*$;
12 **return** $D_{\mathrm{hyp}}$

---

## 6  Post processing

Finally, we applied the following postprocessing rules to the selected translations. The rules are designed based on alignment errors commonly seen in the model translations of the development sets.

- Apply NFKC normalization

次の英語を日本人のネイティブのように日本語に翻訳してください。
原文：{src0} {src1} {src2} 訳文：{hyp0} {hyp1}

Figure 3: The prompt for stepwise MBR-enhanced LLM decoding from English to Japanese. This is an example for translating {src2}. {src0} and {src1} correspond to $S_{hist}$ in Algorithm 1, and {src2} corresponds to $s_i$ in Algorithm 1. Line breaks are added for readability; there are no them in the actual prompt.

- Append an emoji to the end of the hypotheses if it's present at the end of the source sentence
- Replace Japanese brackets (「」) to its Chinese counterparts ("") (Ja→Zh only)
- Replace Japanese commas (、) to its Chinese counterparts (,) (Ja→Zh only)
- Remove whitespaces before and after parentheses
- Remove whitespaces before and after commas, periods, exclamations, and question marks
- Fix letter case of alphabets in the hypotheses to match its counterparts in the source sentence
- Fix punctuations in the hypotheses to match their counterparts in the source sentence

## 7  Post Evaluation

We evaluated the performance of our system using automatic evaluation metrics. Specifically, using this year's test set as the evaluation data, we conducted the evaluation using COMET-22[9] (Rei et al., 2022), MetricX-XL[10] (Juraska et al., 2023), and CometKiwi-XL[11] (Rei et al., 2023) as the evaluation metrics. Note that, since several segments in this year's WMT test set contain multiple sentences, the scores could not be computed at the sentence level.

The results of the post-evaluation from En→Ja are presented in Table 4, while those for the Ja→Zh direction are shown in Table 5. In these tables, "VE" refers to the vocabulary-expanded model, and "CPO" refers to the model where Contrastive Preference Optimization was performed. Additionally, "EncDec" represents outputs from Encoder-Decoder MT models, "MBR (top-p)" refers to the case where MBR decoding was performed using pseudo references generated by top-p sampling, and "MBR (epsilon)" refers to the case where epsilon sampling was used.

**Performance of the LLM-based MT models.** Table 4 shows that the translation performance of Llama2-13B is lower than that of Mistral-7B. One potential reason for this is the limited amount of data used for continual pre-training of Llama2-13B due to constraints in computational resources.

**Efficiency of vocabulary expansion.** Comparing the models with and without vocabulary expansion ((b) vs. (d)), there is no significant difference in performance. However, as shown in Table 13, the model with vocabulary expansion requires fewer training tokens than the model without it in our settings. The generation speed is also faster for the model with vocabulary expansion compared to the one without it. Thus, we believe vocabulary expansion could be a good option for improved inference efficiency.

**CPO is effective but challenging.** Comparing the performance before and after preference learning, the model with vocabulary expansion shows improvement across all evaluation metrics ((d) vs. (e)). On the other hand, the model without vocabulary expansion exhibits a significant decrease in performance for COMET-22 and CometKiwi-XL ((b) vs. (c)), leading to inconsistent results.

Qualitative analysis of outputs from the model without vocabulary expansion (i.e., (c)) revealed instances where decoding of byte-fallbacked text failed, resulting in text being replaced with replacement characters. This may be due to insufficient adjustment of the hyperparameters during CPO training.

**Difference in pseudo references for MBR decoding.** Comparing settings (A) vs. (B) and (C), we observe that the performance improves when using MBR decoding compared to the 1-best output from the ensemble of models[12]. The difference in performance with regard to the pseudo-reference generation algorithms ((i) vs. (j) and (B) vs. (C)) was not significant.

[12]In the En→Ja, we use results from multiple models with different vocabularies for MBR decoding; hence we cannot compare the performance with the 1-best output from the ensemble of all transformers.

|  | COMET-22↑ | MetricX-XL↓ | CometKiwi-XL↑ |
|---|---|---|---|
| (a) Llama2-13B | 0.820 | 3.050 | 0.677 |
| (b) Mistral-7B | 0.841 | 2.806 | 0.711 |
| (c) Mistral-CPO-7B | 0.651 | 2.254 | 0.557 |
| (d) Mistral-VE-7B | 0.836 | 2.881 | 0.695 |
| (e) Mistral-VE-CPO-7B | 0.866 | 2.254 | 0.732 |
| (f) NT5 ([Morishita et al., 2022b](#)) | 0.847 | 2.697 | 0.718 |
| (g) Stepwise MBR-Enhanced LLM Decoding | 0.882 | **2.052** | 0.729 |
| (i) EncDec → MBR (top-p) | **0.885** | 2.263 | 0.737 |
| (j) EncDec → MBR (epsilon) | 0.884 | 2.264 | **0.743** |
| (k) EncDec → MBR (top-p) → LLM Reranking | 0.881 | 2.269 | 0.740 |

Table 4: Results of post evaluation in En→Ja.

|  | COMET-22↑ | MetricX-XL↓ | CometKiwi-XL↑ |
|---|---|---|---|
| (A) EncDec ensemble | 0.818 | 3.550 | 0.548 |
| (B) EncDec → MBR (top-p) | **0.841** | **3.168** | **0.570** |
| (C) EncDec → MBR (epsilon) | **0.841** | 3.230 | 0.566 |

Table 5: Results of post evaluation in Ja→Zh.

**Performance of stepwise MBR-enhanced LLM decoding.** `Stepwise MBR-Enhanced LLM Decoding` achieves the highest score on MetricX-XL. Additionally, compared to using a single LLM, the scores of COMET-22 and MetricX-XL improved. This improvement is likely because generating hypotheses at each step with MBR decoding helps eliminate obvious errors, such as repeated tokens.

**Effectiveness of LLM reranking.** LLM Reranking did not result in any significant improvements according to automatic evaluation metrics. However, we noted improved consistency within segments qualitatively. We intend to evaluate performance through human evaluation as part of future work.

## 8 Submission System

For the final submission system, we adopted system (k) for the En→Ja direction and system (B) for the Ja→Zh direction. However, particularly in the En→Ja direction, different systems ranked highest across various automatic evaluation metrics, leaving us uncertain about which system to select even after post-evaluation. Thus, further refinement of automatic evaluation metrics is essential to develop a superior system.

## 9 Negative Results and Discarded Trials

**Poor performance of LLMs for Japanese-to-Chinese translation.** We conducted continual pre-training and supervised fine-tuning of LLMs for Ja→Zh translation. However, the translation performance did not meet our expectations, leading us to exclude it from the submission system (see Table 5 for post evaluation results). This shortfall likely resulted from our computational resource constraints, which limited continual pre-training to Chinese datasets only. For further details, please refer to Section A.

**Use of LLM outputs as candidates for MBR decoding.** We also explored the inclusion of LLM outputs in the candidate pool for MBR Decoding. However, we observed a decrease in translation quality when these outputs were included, leading us to exclude this approach from the final system. This decline in quality can be attributed to two main factors: i). a substantial difference in the distribution between the outputs generated by LLMs and the pseudo references produced by Encoder-Decoder MT models, and ii). inadequate tuning of hyperparameters during decoding with LLMs.

## 10 Conclusion

This paper described our systems for the constrained track of the WMT 2024 Shared Task: Gen-

eral Machine Translation. We developed translation systems for En→Ja and Ja→Zh. To achieve consistent document-level machine translation, we concentrated on investigating the application of LLMs, which have become available for use this year, employing methods such as LLM Reranking and Stepwise MBR-Enhanced LLM Decoding.

Our submitted system consists of the following steps: i) First, we generate translations using multiple Encoder-Decoder MT models. ii) Next, we narrow down the generated candidates by selecting the optimal translation through MBR decoding. iii) Finally, we apply LLM reranking to incorporate contextual information in order to determine the final output (only for En→Ja). The results from the post-evaluation did not provide quantitative confirmation of the final submission system's effectiveness. However, we did observe a qualitative improvement in consistency within the documents. We hope for future research on better automatic evaluation metrics that can assess these document-level translation performances.

## Acknowledgments

## Contributions

**Keito Kudo** conducted the cleaning of the monolingual data, trained and decoded LLM-based MT models, developed the Ja→Zh Encoder-Decoder MT models, and performed post-evaluations.
**Hiroyuki Deguchi** conducted MBR decoding.
**Makoto Morishita** cleaned the monolingual and bitext data, pre-trained and fine-tuned the Ja→Zh translation model.
**Ryo Fujii** designed and implemented rules for filtering and post-processing, and performed qualitative evaluation of the resulting translations.
**Takumi Ito** designed and customized Hojichar for data cleaning, and designed and implemented Section 4.2.

**Shintaro Ozaki**, **Koki Natsumi** conducted pre-training and fine-tuning of the Ja→Zh Encoder-Decoder MT models, along with back-translation.
**Kai Sato**, **Kazuki Yano** implemented filters for data cleaning.
**Ryosuke Takahashi** implemented filters for data cleaning, conducted preference learning, and prepared scripts for decoding.
**Subaru Kimura** conducted the cleaning of the monolingual data, implemented checkpoint averaging, fine-tuned the LLM-based MT models, and implemented post-processing.
**Tomomasa Hara** implemented filters for the cleaning of the monolingual data and performed hyperparameter tuning for LLM-based MT models.
**Yusuke Sakai** managed the training process for the Ja→Zh Encoder-Decoder MT models.
**Jun Suzuki** provided the primary computational budget and overall project advice and carried out the decoding of NT5.

## References

A.Z. Broder. 1997. On the resemblance and containment of documents. In *Proceedings. Compression and Complexity of SEQUENCES 1997 (Cat. No.97TB100171)*, pages 21–29.

Christian Buck, Kenneth Heafield, and Bas van Ooyen. 2014. N-gram Counts and Language Models from the Common Crawl. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3579–3584, Reykjavik, Iceland. European Language Resources Association (ELRA).

Isaac Caswell, Ciprian Chelba, and David Grangier. 2019. Tagged back-translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 53–63, Florence, Italy. Association for Computational Linguistics.

Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. WIT3: Web Inventory of Transcribed and Translated Talks. In *Proceedings of the 16th Annual Conference of the European Association for Machine Translation*, pages 261–268, Trento, Italy. European Association for Machine Translation.

Jianghao Chen, Pu Jian, Tengxiao Xi, Dongyi Yi, Qianlong Du, Chenglin Ding, Guibo Zhu, Chengqing Zong, Jinqiao Wang, and Jiajun Zhang. 2023. Chinesewebtext: Large-scale high-quality chinese web text extracted with effective evaluation model. *Preprint*, arXiv:2311.01149.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised

Cross-lingual Representation Learning at Scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, and Taro Watanabe. 2024. mbrs: A library for minimum bayes risk decoding. *Preprint*, arXiv:2408.04167.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Bryan Eikema and Wilker Aziz. 2020. Is MAP decoding all you need? the inadequacy of the mode in neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Bryan Eikema and Wilker Aziz. 2022. Sampling-based approximations to minimum Bayes risk decoding for neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10978–10993, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – News Test References for MT Evaluation of 128 Languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online. Association for Computational Linguistics.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT Sentence Embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023. Epsilon sampling rocks: Investigating sampling strategies for minimum Bayes risk decoding for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9198–9209, Singapore. Association for Computational Linguistics.

Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. 2024a. Continual Pre-Training for Cross-Lingual LLM Adaptation: Enhancing Japanese Language Capabilities. *Preprint*, arXiv:2404.17790.

Kazuki Fujii, Taishi Nakamura, and Rio Yokota. 2024b. llm-recipes.

Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. 2012. Building large monolingual dictionaries at the Leipzig corpora collection: From 100 to 200 languages. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 759–765, Istanbul, Turkey. European Language Resources Association (ELRA).

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 Evaluation Benchmark for Low-Resource and Multilingual Machine Translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc'Aurelio Ranzato. 2019. The FLORES Evaluation Datasets for Low-Resource Machine Translation: Nepali–English and Sinhala–English. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6098–6111, Hong Kong, China. Association for Computational Linguistics.

Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. *Preprint*, arXiv:2310.06825.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Herve Jegou, and Tomas Mikolov. 2017a. FastText.zip: Compressing text classification models.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017b. Bag of Tricks for Efficient Text Classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain. Association for Computational Linguistics.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google Submission to the WMT 2023 Metrics Shared Task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Seungduk Kim, Seungtaek Choi, and Myeongho Jeong. 2024. Efficient and effective vocabulary expansion towards multilingual large language models. *CoRR*, abs/2402.14714.

Keito Kudo, Takumi Ito, Makoto Morishita, and Jun Suzuki. 2023. SKIM at WMT 2023 general translation task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 128–136, Singapore. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 66–71. Association for Computational Linguistics.

Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8424–8445, Dublin, Ireland. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2017. SGDR: stochastic gradient descent with warm restarts. In *5th International Conference on Learning Representations, ICLR, Toulon, France, April 24-26, 2017, Conference Track Proceedings*.

Makoto Morishita, Katsuki Chousa, Jun Suzuki, and Masaaki Nagata. 2022a. JParaCrawl v3.0: A Large-scale English-Japanese Parallel Corpus. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6704–6710, Marseille, France. European Language Resources Association.

Makoto Morishita, Keito Kudo, Yui Oka, Katsuki Chousa, Shun Kiyono, Sho Takase, and Jun Suzuki. 2022b. NT5 at WMT 2022 General Translation Task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 318–325, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Chenghao Mou, Chris Ha, Kenneth Enevoldsen, and Peiyuan Liu. 2023. Chenghaomou/text-dedup: Reference snapshot.

Masaaki Nagata, Makoto Morishita, Katsuki Chousa, and Norihito Yasuda. 2024. A Japanese-Chinese Parallel Corpus Using Crowdsourcing for Web Mining. *Preprint*, arXiv:2405.09017.

Graham Neubig. 2011. The Kyoto Free Translation Task. http://www.phontron.com/kftt.

Naoaki Okazaki, Kakeru Hattori, Hirai Shota, Hiroki Iida, Masanari Ohi, Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Rio Yokota, and Sakae Mizuki. 2024. Building a Large Japanese Web Corpus for Large Language Models. *Preprint*, arXiv:2404.17733.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A Fast, Extensible Toolkit for Sequence Modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Reid Pryzant, Youngjoo Chung, Dan Jurafsky, and Denny Britz. 2018. JESC: Japanese-English Subtitle Corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems*, volume 36, pages 53728–53741. Curran Associates, Inc.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 Submission for the Metrics Shared Task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, JosÃ© Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A Neural Framework for MT Evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021. Wiki-Matrix: Mining 135M Parallel Sentences in 1620 Language Pairs from Wikipedia. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361, Online. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving Neural Machine Translation Models with Monolingual Data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 86–96. Association for Computational Linguistics.

Kenta Shinzato. 2023. HojiChar: The text processing pipeline.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No Language Left Behind: Scaling Human-Centered Machine Translation. *Preprint*, arXiv:2207.04672.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *Preprint*, arXiv:2307.09288.

Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. 2020. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl.

Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2020. CCNet: Extracting High Quality Monolingual Datasets from Web Crawl Data. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4003–4012, Marseille, France. European Language Resources Association.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive Preference Optimization: Pushing the Boundaries of LLM Performance in Machine Translation. In *Forty-first International Conference on Machine Learning*.

请像中国本地人一样将以下日语翻译成中文。
原文: {src} 译文: {tgt}

Figure 4: The general prompt for supervised fine-tuning. {src} denotes the source sentence. {tgt} denotes the target sentence. Line breaks are added for readability; there are no them in the actual prompt.

| | COMET-22↑ | MetricX↓ | CometKiwi↑ |
|---|---|---|---|
| Llama2-13B | 0.754 | 4.763 | 0.503 |
| Mistral-7B | 0.795 | 4.410 | 0.547 |
| EncDec ensemble | 0.818 | 3.550 | 0.548 |

Table 6: Post evaluation results of the LLM trained for Ja→Zh translation. Compared to the ensemble of Encoder-Decoder MT models, the performance of the LLM for Ja→Zh translation was not sufficient.

## A  Japanese-Chinese LLM

**Training configurations.**    We trained LLMs for Ja→Zh translation, although these were not included in the final system. Due to time and computational resource constraints, we only conducted continual pre-training and supervised fine-tuning on Chinese monolingual corpora. During supervised fine-tuning, we used the template shown in Figure 4. Table 12, 13 lists the hyperparameters used for training in the Ja→Zh direction.

**Post evaluation.**    We conducted evaluations for the LLMs trained for the Ja→Zh translation. Table 6 presents the results. The performance of the LLMs in the Ja→Zh translation was insufficient compared to the ensemble of Encoder-Decoder MT models. This is likely because we were limited to continual pre-training using only Chinese corpora due to computational resource constraints.

## B  Vocabulary Expansion for LLM

As described in Section 3.2, we aimed to improve the Japanese language generation capability of Mistral-7B by expanding the model's vocabulary. Here, we provide details on the vocabulary expansion.

**Construction of additional vocabulary.**    We first constructed a Japanese vocabulary using the unigram algorithm of the `Sentencepiece` tool (Kudo and Richardson, 2018). This vocabulary was trained on a subset of 30,000,000 samples from the Japanese Monolingual Corpus. We set the vocabulary size to 27,000. During vocabulary

training, we enabled the options "byte_fallback" and "split_digits".

**Vocabulary initialization.**    We initialized the embeddings for the additional vocabulary using the weighted average of the original Mistral embeddings. The weights were determined based on the similarity scores between the new and original Mistral vocabularies, computed by LaBSE (Feng et al., 2022). The process is described by the following equation:

$$
\begin{aligned}
\mathbf{v}_{\text{new}} &= \sum_{i=1}^{N} \left( \frac{\exp(w_i)}{\sum_{j=1}^{N} \exp(w_j)} \right) v_i \\
&= \sum_{i=1}^{N} \text{softmax}(w_i) v_i
\end{aligned}
\tag{6}
$$

Here, $\mathbf{v}_{\text{new}}$ represents the embedding for the additional vocabulary, $w_i$ is the similarity score between the additional vocabulary and vocabulary entry $i$ as calculated by LaBSE, $\mathbf{v}_i$ is the vector of the existing vocabulary entry $i$, and $n$ is the size of the original vocabulary. This method was also used to initialize the language modeling head.

Given our focus on the English-to-Japanese translation task, vocabularies other than English and Japanese are considered less critical. Therefore, we replaced any vocabulary not identified as Japanese, English, or special tokens with the new additional vocabulary. The determination of the language for each token followed these rules:

**Japanese:** Tokens consisting of *hiragana*, *katakana*, common-use *kanji*, symbols, JIS level 1 *kanji*, and ASCII characters
**English:** Tokens consisting solely of ASCII characters
**Special tokens:** Tokens split by byte fallback, as well as bos, eos tokens, etc.

Consequently, we expanded the vocabulary to 51,200.

**Vocabulary warmup training.**    To address inconsistencies introduced by adding new vocabulary, prior research has proposed gradually training the model while fixing specific parameters after adding the vocabulary (Kim et al., 2024). We adopted a similar method to resolve these inconsistencies. Initially, we fixed the parameters of all transformer layers except for the embedding layer and the language modeling head and conducted the training. The hyperparameters used during this initial training phase are detailed in Table 11.

## C Training Hyperparameters

The hyperparameters during the training of each model are shown in Table 7- 13.

| Initial Translation Model | |
|---|---|
| Subword Size | 32,000 |
| Architecture | Transformer (big) with 6 layers, Encoder and Decoder FFN size of 8,192 |
| Optimizer | Adam $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 1 \times 10^{-8}$, weight_decay $= 0.0$ |
| Learning Rate Schedule | Inverse square root decay, Cosine |
| Warmup Steps | 4,000 |
| Max Learning Rate | 0.001 |
| Dropout | 0.1 |
| Gradient Clip | 1.0 |
| Batch Size | 1,048,576 tokens |
| Max Number of Updates | 50,000 steps |
| Averaging | Save a checkpoint every 500 steps and average the last ten |
| Implementation | `fairseq` (Ott et al., 2019) |

| Pre-training Configuration | |
|---|---|
| Subword Size | 16,000 |
| Architecture 1 | Transformer (big) with 9 layers, Encoder FFN size of 16,384, and Decoder FFN size of 4,096 |
| Architecture 2 | Transformer (big) with 9 layers, Encoder and Decoder FFN size of 8,192 |
| Optimizer | Adam $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 1 \times 10^{-8}$, weight_decay $= 0.0$ |
| Learning Rate Schedule | Inverse square root decay, Cosine |
| Warmup Steps | 4,000 |
| Max Learning Rate | 0.001 |
| Dropout | 0.1 |
| Gradient Clip | 0.1 |
| Batch Size | 1,048,576 tokens |
| Max Number of Updates | 50,000 steps |
| Averaging | Save a checkpoint every 500 steps and average the last ten |
| Implementation | `fairseq` (Ott et al., 2019) |

| Fine-tuning Configuration | |
|---|---|
| Learning Rate Schedule | Fixed |
| Warmup Steps | N/A |
| Max Learning Rate | $1 \times 10^{-5}$ |
| Dropout | 0.2 |
| Gradient Clip | 1.0 |
| Batch Size | 14,400 tokens |
| Max Number of Updates | 1,000 steps |
| Averaging | Save a checkpoint every ten steps and average the last ten |

Table 7: List of hyper-parameters. We used the initial translation model to generate synthetic data, the pre-training configuration to build the models described in Section 3.1, and the fine-tuning configuration to develop the models for submission. We created two models for pre-training and fine-tuning, labeled as "Architecture 1" or "Architecture 2," and used them for ensembling. The hyperparameters listed in the fine-tuning configuration represent only the differences from the pre-training configuration.

**Llama2-13B Pretraining**

| | |
|---|---|
| Vocab Size | 32,000 |
| Train Steps | 10,000 |
| Batch Size | 1,572,864 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 250 |
| Max Learning Rate | $2 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ |
| | $\epsilon = 1 \times 10^{-6},$ |
| | weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 100 steps and average the last five |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

**Llama2-13B Supervised Finetuning**

| | |
|---|---|
| Vocab Size | 32,000 |
| Train Steps | 3,500 |
| Batch Size | 1,310,720 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 175 |
| Max Learning Rate | $3 \times 10^{-6}$ |
| Min Learning Rate | $3 \times 10^{-7}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ |
| | $\epsilon = 1 \times 10^{-6},$ |
| | weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 100 steps and average the last three |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

Table 8: A list of hyperparameters used when training Llama2-13B on the En→Ja task.

**Mistral-7B Pretraining**

| | |
|---|---|
| Vocab Size | 32,000 |
| Train Steps | 20,000 |
| Batch Size | 1,310,720 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 500 |
| Max Learning Rate | $2 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ |
| | $\epsilon = 1 \times 10^{-6},$ |
| | weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 200 steps and average the last five |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

**Mistral-7B Supervised Finetuning**

| | |
|---|---|
| Vocab Size | 32,000 |
| Train Steps | 3,100 |
| Batch Size | 1,310,720 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 155 |
| Max Learning Rate | $1 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ |
| | $\epsilon = 1 \times 10^{-6},$ |
| | weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 200 steps and average the last three |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

**Mistral-7B Preference Learning**

| | |
|---|---|
| Vocab Size | 32,000 |
| Train Steps | 250 |
| Batch Size | 144 samples |
| Learning Rate Schedule | Constant |
| Learning Rate | $1 \times 10^{-5}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.999,$ |
| | $\epsilon = 1 \times 10^{-8},$ |
| | weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| CPO $\beta$ | 0.1 |
| CPO $\alpha$ | 1.5 |
| $i_w$ (See Section 3.2) | 740 |
| Lora $r$ | 16 |
| Lora $\alpha$ | 32 |
| Lora Dropout | 0.1 |
| Lora Target Layetr | All linear layer |
| Implementation | Transformers (Wolf et al., 2020), TRL (von Werra et al., 2020) |

Table 9: A list of hyperparameters used when training Mistral-7B on the En→Ja task.

**Mistral-7B (vocab expanded) Pretraining**

| | |
|---|---|
| Vocab Size | 51,200 |
| Train Steps | 12,283 |
| Batch Size | 1,376,256 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 300 |
| Max Learning Rate | $2 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ $\epsilon = 1 \times 10^{-6},$ weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 200 steps and average the last five |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

**Mistral-7B (vocab expanded) Supervised Finetuning**

| | |
|---|---|
| Vocab Size | 51,200 |
| Train Steps | 2,000 |
| Batch Size | 1,310,720 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 100 |
| Max Learning Rate | $1 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ $\epsilon = 1 \times 10^{-6},$ weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 200 steps and average the last two |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

**Mistral-7B (vocab expanded) Preference Learning**

| | |
|---|---|
| Vocab Size | 51,200 |
| Train Steps | 250 |
| Batch Size | 144 samples |
| Learning Rate Schedule | Fixed |
| Learning Rate | $1 \times 10^{-5}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.999,$ $\epsilon = 1 \times 10^{-8},$ weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| CPO $\beta$ | 0.1 |
| CPO $\alpha$ | 1.5 |
| $i_w$ (See Section 3.2) | 740 |
| Lora $r$ | 16 |
| Lora $\alpha$ | 32 |
| Lora Dropout | 0.1 |
| Lora Target Layer | All linear layers |
| Implementation | Transformers (Wolf et al., 2020), TRL (von Werra et al., 2020) |

Table 10: A list of hyperparameters used when training Mistral-7B with vocabulary expansion on the En→Ja task.

**Mistral-7B (vocab extended) Vocabulary Warmup**

| | |
|---|---|
| Vocab Size | 51,200 |
| Train Steps | 1800 |
| Batch Size | 1,376,256 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 50 |
| Max Learning Rate | $2 \times 10^{-4}$ |
| Min Learning Rate | $6.6 \times 10^{-7}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ $\epsilon = 1 \times 10^{-6},$ weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

Table 11: A list of hyperparameters used when training Mistral-7B with vocabulary expansion for vocabulary warmup on the En→Ja task.

**Llama2-13B Pretraining**

| | |
|---|---|
| Vocab Size | 32,000 |
| Train Steps | 10,000 |
| Batch Size | 1,572,864 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 250 |
| Max Learning Rate | $2 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ $\epsilon = 1 \times 10^{-6},$ weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 100 steps and average the last five |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

**Llama2-13B Supervised Finetuning**

| | |
|---|---|
| Vocab Size | 32,000 |
| Train Steps | 500 |
| Batch Size | 1,310,720 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 25 |
| Max Learning Rate | $3 \times 10^{-6}$ |
| Min Learning Rate | $3 \times 10^{-7}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9, \beta_2 = 0.95,$ $\epsilon = 1 \times 10^{-6},$ weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 25 steps and average the last three |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

Table 12: A list of hyperparameters used when training Llama2-13B on the Ja→Zh task.

| Mistral-7B Pretraining | |
| --- | --- |
| Vocab Size | 32,000 |
| Train Steps | 20,000 |
| Batch Size | 1,310,720 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 500 |
| Max Learning Rate | $2 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9$, $\beta_2 = 0.95$, $\epsilon = 1 \times 10^{-6}$, weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 200 steps and average the last five |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

| Mistral-7B Supervised Finetuning | |
| --- | --- |
| Vocab Size | 32,000 |
| Train Steps | 420 |
| Batch Size | 1,310,720 tokens |
| Learning Rate Schedule | Cosine (Loshchilov and Hutter, 2017) |
| Warmup Steps | 25 |
| Max Learning Rate | $1 \times 10^{-5}$ |
| Min Learning Rate | $1 \times 10^{-6}$ |
| Optimizer | Adam |
| | $\beta_1 = 0.9$, $\beta_2 = 0.95$, $\epsilon = 1 \times 10^{-6}$, weight_decay $= 0.1$ |
| Gradient Clip | 1.0 |
| Averaging | Save a checkpoint every 10 steps and average the last five |
| Implementation | Transformers (Wolf et al., 2020), llm-recipies (Fujii et al., 2024b) |

Table 13: A list of hyperparameters used when training Mistral-7B on the Ja→Zh task.

# GTCOM and DLUT's Neural Machine Translation Systems for WMT24

**Hao Zong**[1]      **Chao Bei**[2]      **Conghu Yuan**[2]
**Wentao Chen**[2]      **Huan Liu**[2]      **Degen Huang**[1*]

[1]Dalian University of Technology
[2]Global Tone Communication Technology Co., Ltd.

zonghao@mail.dlut.edu.cn
{beichao, yuanconghu, chenwentao and liuhuan}@gtcom.com.cn
huangdg@dlut.edu.cn

## Abstract

This paper presents the submission from Global Tone Communication Co., Ltd. and Dalian University of Technology for the WMT24 shared general Machine Translation (MT) task at the Conference on Empirical Methods in Natural Language Processing (EMNLP). Our participation encompasses two language pairs: English to Japanese and Japanese to Chinese. The systems are developed without particular constraints or requirements, facilitating extensive research in machine translation. We emphasize back-translation, utilize multilingual translation models, and apply fine-tuning strategies to improve performance. Additionally, we integrate both human-generated and machine-generated data to fine-tune our models, leading to enhanced translation accuracy. The automatic evaluation results indicate that our system ranks first in terms of BLEU score for the Japanese to Chinese translation.

## 1 Introduction

In this study, we employ fairseq (Ott et al., 2019) as our development framework and adopt the transformer (Vaswani et al., 2017) as the main architecture. The primary ranking index for the submitted systems is BLEU (Papineni et al., 2002), which also serves as the evaluation metric for our translation system via sacreBLEU[1], consistent with our methodology from the previous year.

For data preprocessing, we conduct punctuation normalization, tokenization, and Byte Pair Encoding (BPE) (Sennrich et al., 2015) across all languages involved. Furthermore, we applied a true-case model for English, tailored to the specific linguistic features of each language. Regarding tokenization, we utilize Jieba[2] for Chinese, Mecab[3] for Japanese, and the Moses tokenizer.perl (Koehn

et al., 2007) for English. Additionally, we incorporate knowledge-based rules along with a language model to cleanse parallel data, monolingual data, and synthetic data.

For the multilingual translation model, we consolidate all languages into a single model and enhance it with an English to Chinese parallel corpus to enrich the language information.

The remainder of this paper is structured as follows: Section 2 discusses the translation task and provides dataset statistics. Section 3 describes our baseline systems and introduces the proposed multilingual translation model. The data selection methodology is elaborated in Section 4. Section 5 presents experiments conducted on all translation directions, addressing data filtering, model architectures, back-translation, joint training strategies, adaptations of the multilingual model, fine-tuning, data selection, and ensemble decoding. Section 6 analyzes the results, offering insights into the efficacy of various techniques. Finally, Section 7 concludes the paper.

## 2 Task Description

This task focuses on bilingual text translation, with the provided data elaborated in Table 1, which includes both parallel and monolingual data. For the English-Japanese directions, the primary sources of parallel data include WikiMatrix (Schwenk et al., 2019), CCAligned (Rozis and Skadiņš, 2017), JESC (Pryzant et al., 2017), JParaCrawl v3.0 (Morishita et al., 2022), LinguaTools-WikiTitles (Tiedemann, 2012), News Commentary v16, and XLEnt (Tiedemann, 2012). For the Japanese-Chinese direction, the main parallel data is sourced from CCAligned, JParaCrawl, LinguaTools-WikiTitles, News Commentary v16, WikiMatrix, and XLEnt. Monolingual data comprises News Crawl (Kocmi et al., 2022) in English, Japanese, and Chinese; News Commentary in English, Japanese, and Chinese; and Europarl v10 in English. We uti-

---

| Language | Number of Sentences |
|---|---|
| en-ja parallel data | 85.2M |
| ja-zh parallel data | 14.4M |
| en monolingual data | 168M |
| ja monolingual data | 22.8M |
| zh monolingual data | 23.9M |
| en-ja development set | 1000 |
| ja-zh development set | 1012 |

Table 1: Task Description

lized the provided development set from new-stest2020 for English-Japanese and the FLoRes101 (NLLB Team, 2022) dataset for Japanese-Chinese.

## 3 Bilingual Baseline Model and Multilingual Translation Model

To establish a robust baseline for comparison with the multilingual model, we utilize the transformer_wmt_en_de as our bilingual baseline model, consisting of 24 encoder layers and 24 decoder layers. The multilingual translation model is designed to closely resemble the GTCOM2023 (Zong, 2023) model, referred to as the X to X model. To achieve superior translation quality, we include the English-Chinese parallel corpus as the primary auxiliary language pair to enhance linguistic information. We train a single multilingual model that encompasses all translation directions while applying joint Byte Pair Encoding (BPE) separately for all languages.

## 4 Data Selection

Similar to the last year, we use source test sets to train a text classification model based on RoBERTa (Liu et al., 2019). Specifically, we treat the in-domain test set as positive examples and select an equivalent amount of sentence pairs from the out-of-domain test set as negative examples. We fine-tune RoBERTa on this labeled dataset to develop a binary classifier capable of effectively distinguishing between in-domain and out-of-domain data. This classifier aids in selecting domain-specific training data from the general training corpus, with the chosen in-domain training data subsequently used to fine-tune the multilingual neural machine translation model.

Additionally, we also use prompt learning to explore an alternative data selection method. We develop a prompt template and leverage the gen-erative capabilities of Meta-Llama-3-8B-Instruct [4] to create a domain classifier using loRA (Hu et al., 2021). The prompt template mirrors that used in GTCOM2023 from the last year, shows in Table 2. Specifically, we extract 800 sentences from the development set which belong to the news, social, e-commerce, or conversation domains. We manually select 200 sentences from the training set that do not match these domains or are of inferior quality, categorizing them as "other." We then utilize these 1,000 labeled examples to fine-tune the Meta-Llama-3-8B-Instruct model in loRA. The resulting prompt-based classifier effectively differentiates between domains in the training data. Sentences predicted as "News," "Social," "E-commerce," and "Conversation" are classified as in-domain data, while those labeled as "Other" are considered out-of-domain data.

## 5 Experiment

This section outlines the step-by-step experiments we conducted, with the entire workflow depicted in Figure 1.

- **Data Filtering:** The data filtering techniques largely replicate those utilized last year, incorporating human rules, language models, and repetition cleaning.

- **Baseline:** Our baseline is constructed using the transformer big architecture, which comprises 24 encoder layers and 24 decoder layers.

- **Back-translation:** We employ the best translation model to translate target sentences back to the source side, cleaning synthetic data using a language model. This process includes translating each language pair featured in the multilingual translation model. We combine the cleaned back-translation data with parallel sentences and train the multilingual translation model accordingly.

- **Joint Training:** We repeat the back-translation step using the optimal model until no further improvements are observed.

- **Multilingual Translation Model:** A single model is trained for all translation directions, with each direction utilizing joint BPE and a

---

[4]https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct

| | |
|---|---|
| Instructions | Please determine the domain to which the given sentence belongs based on the following criteria.<br>1. Sentence Correctness: If the sentence is incomplete, incoherent, or grammatically incorrect, label it as "Other" domain. If the sentence is complete, fluent, and grammatically correct, proceed to the next step.<br>2. Domain Identification: Analyze the content of the sentence to identify the possible domain it belongs to. Consider the following domains: News, Social, E-commerce, Conversation, and Other. If the sentence shows clear indications of being from a specific domain, label it accordingly, otherwise label it as "Other" domain.<br>Please label the sentence with the appropriate domain:<br>- If the sentence is from the News domain, label it as "News".<br>- If the sentence is from the Social domain, label it as "Social".<br>- If the sentence is from the E-commerce domain, label it as "E-commerce".<br>- If the sentence is from the Conversation domain, label it as "Conversation".<br>- If the sentence does not fit any specific domain or is incorrect, label it as "Other". |
| Sentence | Sunday Best: Enter 1880s New York in HBO's "The Gilded Age" |
| Domain | News |

Table 2: Prompt Template.



Figure 1: The work flow of GTCOM machine translation competition systems

shared vocabulary. The multilingual translation model consists of 24 encoder layers and 24 decoder layers, employing the transformer big architecture.

- **Fine-tuning:** The multilingual translation model is fine-tuned for each direction and bidirection separately. For instance, we fine-tuned en2ja and ja2en on the multilingual translation model and fine-tuned en2ja on the multilingual translation model for English to Japanese separately.

- **Data Selection:** The model described in the Data Selection section is employed to choose a domain-specific training dataset, which is then fine-tuned on the multilingual translation model.

- **Ensemble Decoding:** We utilize the GMSE Algorithm (Deng et al., 2018) to select models, aiming for optimal performance.

## 6 Results and Analysis

Table 3 displays the BLEU scores evaluated on the development set for English to Japanese and Japanese to Chinese. As indicated in the table, back-translation remains the most effective data augmentation technique for enhancing translation quality from a data perspective. The multilingual translation model also demonstrates significant improvements across all translation directions. As shown in Table 4, our prompt learning strategy is

| Model | en2ja | ja2zh |
|---|---|---|
| Baseline | 26.36 | 15.07 |
| + Back-translation | 27.26 | 20.75 |
| Multilingual Translation Model | 26.50 | 15.20 |
| + Back-translation | 27.40 | 21.24 |
| + Bilingual Fine-tuning | 27.51 | 21.34 |
| + Single Fine-tuning | 27.22 | 20.98 |
| Ensemble Decoding | 27.95 | 22.21 |

Table 3: BLEU scores for English to Japanese and Japanese to Chinese. Values are calculated based on word counts.

| Direction | BLEU | BLEU with DS |
|---|---|---|
| en-ja | 39.2 | 39.7 |
| ja-zh | 32.9 | 32.3 |

Table 4: The final online automatic evaluation BLEU with/without prompt learning in data selection.

still able to improve the BLEU score on the direction of English to Japanese, but there was some decline in the Japanese-to-Chinese direction.

## 7 Conclusion

This paper introduces the neural machine translation systems developed by GTCOM and DLUT for the WMT24 shared general MT task. We apply three primary techniques to enhance translation quality: back-translation, a multilingual translation model, and fine-tuning accompanied by data selection. Through these methods, we achieve notable improvements in automatic evaluation metrics, as illustrated in Table 5.

## Acknowledgments

| Direction | BLEU | CometKiwi |
|---|---|---|
| en-ja | 39.7 | 0.697 |
| ja-zh | 32.9 | 0.586 |

Table 5: Final online automatic evaluation results.

## References

Yongchao Deng, Shanbo Cheng, Jun Lu, Kai Song, Jingang Wang, Shenglan Wu, Liang Yao, Guchun Zhang, Haibo Zhang, Pei Zhang, et al. 2018. Alibaba's neural machine translation systems for wmt18. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 368–376.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions*, pages 177–180.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Makoto Morishita, Katsuki Chousa, Jun Suzuki, and Masaaki Nagata. 2022. Jparacrawl v3. 0: A large-scale english-japanese parallel corpus. *arXiv preprint arXiv:2202.12607*.

James Cross Onur Çelebi Maha Elbayad Kenneth Heafield Kevin Heffernan Elahe Kalbassi Janice Lam Daniel Licht Jean Maillard Anna Sun Skyler Wang Guillaume Wenzek Al Youngblood Bapi Akula Loic Barrault Gabriel Mejia Gonzalez Prangthip Hansanti John Hoffman Semarley Jarrett Kaushik Ram Sadagopan Dirk Rowe Shannon Spruit Chau Tran Pierre Andrews Necip Fazil Ayan Shruti Bhosale Sergey Edunov Angela Fan Cynthia Gao Vedanuj

---

Goswami Francisco Guzmán Philipp Koehn Alexandre Mourachko Christophe Ropers Safiyyah Saleem Holger Schwenk Jeff Wang NLLB Team, Marta R. Costa-jussà. 2022. No language left behind: Scaling human-centered machine translation.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Reid Pryzant, Yongjoo Chung, Dan Jurafsky, and Denny Britz. 2017. Jesc: Japanese-english subtitle corpus. *arXiv preprint arXiv:1710.10639*.

Roberts Rozis and Raivis Skadiņš. 2017. Tilde model-multilingual open data for eu languages. In *Proceedings of the 21st Nordic Conference on Computational Linguistics*, pages 263–265.

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2019. Wikimatrix: Mining 135m parallel sentences in 1620 language pairs from wikipedia. *arXiv preprint arXiv:1907.05791*.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In *Lrec*, volume 2012, pages 2214–2218. Citeseer.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.

Hao Zong. 2023. Gtcom and dlut's neural machine translation systems for wmt23. In *Proceedings of the Eighth Conference on Machine Translation*, pages 192–197.

# CUNI at WMT24 General Translation Task:
# LLMs, (Q)LoRA, CPO and Model Merging

**Miroslav Hrabal, Josef Jon, Martin Popel, Nam H. Luu, Danil Semin, Ondřej Bojar**

Charles University, Faculty of Mathematics and Physics

{hrabal,jon,popel,bojar}@ufal.mff.cuni.cz,
namhoang.luu700@student.cuni.cz, dsemin2305@gmail.com

## Abstract

This paper presents the contributions of Charles University teams to the WMT24 General Translation task (English to Czech, German and Russian, and Czech to Ukrainian) and the WMT24 Translation into Low-Resource Languages of Spain task. Our most elaborate submission, CUNI-MH for en2cs, is the result of fine-tuning Mistral 7B v0.1 for translation using a three-stage process: Supervised fine-tuning using QLoRA, Contrastive Preference Optimization, and merging of model checkpoints. We also describe the CUNI-GA, CUNI-Transformer and CUNI-DocTransformer submissions, which are based on our systems from the previous year.

Our en2ru system CUNI-DS uses a similar first stage as CUNI-MH (QLoRA for en2cs) and follows with transfer learning for en2ru.

For en2de (CUNI-NL), we experimented with an LLM-based speech translation system, to translate without the speech input.

For the Translation into Low-Resource Languages of Spain task, we performed QLoRA fine-tuning of a large LLM on a small amount of synthetic (backtranslated) data.

## 1 Introduction

This paper describes the CUNI submissions to the WMT24 General Translation task (from English to Czech, German and Russian, and from Czech to Ukrainian) and the Translation into Low-Resource Languages of Spain task.

Our underlying goal for this year was to test the applicability of primarily small open-source LLMs to the languages of interest, and we also provide our English-to-Czech systems from the previous years for comparison.

The setups for the various target languages differ considerably in the methods used. Table 1 provides an overview of the individual system highlights. In Section 2, we detail the basic building steps and methods across our systems (not all setups use all

of them). Section 3 describes the training and development data used across the target languages. In Section 4, we evaluate the systems and compare their results with various available baselines and benchmarks. Section 5 summarizes our future plans, and we conclude in Section 6.

## 2 Methods

For the CUNI-MH submission, we fine-tuned Mistral 7B v0.1 (Jiang et al., 2023) using three stages:

1. Supervised fine-tuning on CzEng 2.0 training dataset (Kocmi et al., 2020)[1], see Section 2.3.

2. Contrastive Preference Optimization (Xu et al., 2024b), see Section 2.4.

3. Averaging model checkpoints (Utans, 1996; Wortsman et al., 2022; Gueta et al., 2023), see Section 2.5.

CUNI-Transformer and CUNI-DocTransformer are the same systems as submitted last year (Jon et al., 2023), relying on standard NMT training with Block backtranslation (Section 2.1) and optionally document-level training (Section 2.2).

For CUNI-GA, in English-to-Czech, we used outputs from CUNI-Transformer and a genetic algorithm to combine and modify them, again in the same way as previous year (Section 2.8; Jon et al., 2023; Jon and Bojar, 2023). For coincidentally identically called CUNI-GA submission in Translation into Low-Resource Languages of Spain task, we fine-tune larger LLMs (Command-R and Aya-23), without applying the genetic algorithm.

For the CUNI-NL system, we fine-tuned Llama 2 7B (Touvron et al., 2023) for the speech translation task, while also adapting it for text-only translation at the same time; see Section 2.6.

Finally CUNI-DS starts as step 1 of CUNI-MH but continues with transfer learning to target Russian instead of Czech, see Section 2.7.

---

[1] http://ufal.mff.cuni.cz/czeng/

| Task | CUNI-* Model | Initial LLM | SFT Data | SFT Highlights (§2.3) | Final Stages |
|------|-------------|-------------|----------|----------------------|--------------|
| cs2uk | Transformer | - | Opus, CzEng | BlockBT §2.1 | - |
| en2cs | DocTransformer | - | CzEng 2.0 | BlockBT §2.1, doc-level §2.2 | - |
| en2cs | GA | - | - | - | GA §2.8 |
| en2cs | MH | Mistral 7B v0.1 | CzEng 2.0 | QLoRA, Packing, AdamW | CPO §2.4; Checkpoint Merging §2.5 |
| spa | GA | Command-R, Aya | PILAR BT | QLoRA | - |
| en2de | NL | HuBERT, Llama-2-7b | MuST-C | Text-only use of a speech translation system §2.6 | |
| en2ru | DS | Mistral 7B v0.1 | CzEng, Yandex, News Commentary | Transfer from en2cs §2.7 | - |

Table 1: Overview of CUNI systems in WMT24 General Translation task and Translation into Low-Resource Languages of Spain task (spa). Systems in the upper part of the table are our last year's baselines. §· refer to the methods in Section 2.

## 2.1 BlockBT

For training CUNI-Transformer and CUNI-DocTransformer, we used iterated Block backtranslation (BlockBT) (Popel, 2018; Popel et al., 2020; Gebauer et al., 2021; Jon et al., 2022) in a standard Transformer (Vaswani et al., 2017) NMT training from scratch. The BlockBT method organizes the training data, so that the model can optimize the balance between authentic English-to-Czech parallel texts (exhibiting more translationese artifacts) and synthetic data created by back-translating Czech-only texts) by averaging eight checkpoints reflecting more of the former or the latter domain. The use of eight checkpoints for averaging is derived from the original paper (Popel, 2018) and a study on hyperparametrs for training Transformers (Popel and Bojar, 2018).

## 2.2 Document-level training

The approach for training CUNI-DocTransformer is described in Popel et al. (2019). Starting with the initial sentence-level model (CUNI-Transformer), we continued training on sequences of consecutive sentences coming from a coherent text with at most 3000 characters, where both sides (en and cs) have the same number of sentences. The sentences are separated by a special token in each of the languages.

## 2.3 Supervised fine-tuning (SFT)

For the CUNI-MH submission, we used 4-bit QLoRA (Dettmers et al., 2023) with a large LoRA rank of $r = 512$. We used a batch size of 32, a learning rate of $2e - 5$, 20 warm-up steps, 8-bit AdamW (Loshchilov and Hutter, 2019) optimizer

and weight decay of 0.01. We also used a scheduler with linear learning rate decay. Starting from the freely available Mistral 7B v0.1 model, we trained in a language modeling fashion on individual sentences, calculating the loss on each token. To reduce the number of padding tokens, we also used packing: examples are concatenated with the EOS token as a separator to achieve a total sequence length of 1000. In Appendix A, we present our translation prompt template and example of its processed form with packing as used during training.

We trained for a single epoch on the authentic part of CzEng 2.0. In Figure 1, we show how the performance of the model develops during the first stage, starting from 100 steps. A notable observation is that the COMET22 and COMETKIWI22 scores seem to plateau relatively early, despite the evaluation loss steadily decreasing, while BLEU seems to be steadily increasing. This appears to be consistent with the results presented by Xu et al. (2024a), although we suspect it could also be the result of insufficient regularization.

For training, we used the HuggingFace Transformers and TRL libraries by Wolf et al. (2020) and von Werra et al. (2020). We also used the Unsloth library,[2] which provides speed and VRAM optimizations to Transformers and TRL libraries.

Another of our submissions that made use of a pre-trained LLM and SFT was CUNI-GA in the Translation into Low-Resource Languages of Spain task. We used 4-bit QLoRA with the rank of $r = 16$ and the learning rate of $4e - 4$ for fine-tuning the pretrained *Command-R* model, and $1e - 3$ for fine-tuning the *Aya* model, with an effective batch size

Figure 1: CUNI-MH Stage 1 – metrics during training.

of 32 and an AdamW optimizer with the weight decay value of 0.001.

## 2.4 Contrastive Preference Optimization (CPO)

CPO is a fine-tuning method introduced by Xu et al. (2024b) as an approximation of Direct Preference Optimization (Rafailov et al., 2024).

The goal of CPO is to fine-tune the model to directly optimize for preferences between translation candidates, rather than just optimizing the likelihood of the reference translations.

From a high-level point of view, the main difference between using SFT and CPO for translation is that for a given source text, we need two translations: *preferred* and *dis-preferred*. This means that the training dataset consists of triplets, rather than pairs as is typical for supervised training of NMT. For a more detailed description of the dataset we used and how it was created, see Section 3.2.

To apply CPO during the second stage of CUNI-MH training, we started two separate training runs from models we created during the first stage. One

of the runs starts from model ③ and the other from model ④ in Table 2.

We selected these models because they had the best COMET22 and COMETKIWI22 scores among the models we had available at the time, when evaluated on the sentence-level WMT22 validation set.

Because we wanted to use a smaller LoRA rank size comparable to those used in the original paper (Xu et al., 2024b), we merged LoRA adapters with the quantized model into a 16-bit model and added new, smaller adapters.

We trained for two epochs with the following parameters: LoRA rank $r = 32$, LoRA $\alpha = 64$, CPO $\beta = 0.1$. We trained two separate runs, starting from the checkpoints mentioned earlier. Similarly to the SFT stage, we used 8-bit AdamW, this time without learning rate decay. Our GPU memory capacity was limiting us to the batch size of 4, so to compensate, we used 64 gradient accumulation steps to simulate a larger effective batch size of 256.

| Stage | ID | Model | Checkpoint | COMET22 | COMETKIWI22 | BLEU |
|-------|----|-------|-----------|---------|-------------|------|
|  | ⓪ | Mistral 7B v0.1 5-shot |  | 67.16 | 59.79 | 17.35 |
| 1 | ① |  | 16000 | 85.59 | 79.04 | 33.46 |
| 1 | ② | SFT from ⓪ | 24000 | 86.10 | 79.40 | 34.35 |
| 1 | ③ |  | 103000 | 85.80 | 78.85 | 35.32 |
| 1 | ④ | SLERP merge of ① and ② |  | 86.16 | 79.44 | 35.15 |
| 2 | ⑤ | CPO from ④ | 150 | 89.76 | 82.71 | 32.56 |
| 2 | ⑥ | CPO from ③ | 100 | 89.93 | 83.04 | 34.43 |
| 2 | ⑦ | CPO from ⓪ | 400 | 83.21 | 76.54 | 18.33 |
| **3** | **⑧** | **Linear merge of ⑤ and ⑥** |  | **90.21** | **83.16** | **36.52** |

Table 2: CUNI-MH's training stages, models and their sentence-level scores on WMT23 (test set). The final CUNI-MH submission ⑧ is in bold.

Checkpoints were saved every 50 steps[3] and evaluated on the validation test set using COMETKIWI22. The performance peaked around checkpoint 150 for the first run, leading us to conclude that further training beyond 2 epochs was unnecessary. However, we acknowledge that the training parameters may not be optimal and could potentially be tweaked further for better results.

## 2.5 Checkpoint merging

To further improve the performance of the CUNI-MH model, we experimented with two methods for merging model weights: linear interpolation (Utans, 1996) and spherical linear interpolation (SLERP, Shoemake, 1985) in different training stages.

In particular, after the SFT stage, we merged two promising checkpoints from the same training run using SLERP, which led to a small improvement in all metrics, as can be seen by looking at model ④ in Table 2.

After the CPO stage, we once again experimented with model merging, this time we merged the best performing checkpoints from two different CPO training runs. This led to a further modest improvement in all COMET22, COMETKIWI22 and BLEU metrics, as shown by model ⑧ in Table 2.

For model merging using both SLERP and linear interpolation, we used the mergekit library by Goddard et al. (2024).

## 2.6 SFT from Speech Translation System (SFTSpeech)

The CUNI-NL system was adapted from a speech translation system, which features a frozen Hu-

BERT component (Hsu et al., 2021) and the Llama 2 7B (Touvron et al., 2023) LLM.

The original speech translation system applied the CTC collapsing strategy to extract the speech hidden features; these features would subsequently be given as the prompt to a LLM to generate the ASR transcription and its corresponding translation simultaneously.

For the purposes of the General Translation Task, we avoid any audio features during inference and directly prompt the LLM with the source language text. We expect the LLM to translate using that only information. The motivation for this experiment was to check if a LLM-based speech translation system remains versatile enough to support text-only translation.

The original speech translation system was a fine-tuned LLM using 4-bit QLoRA (Dettmers et al., 2023) adapters, with the rank of $r = 8$ and alpha of $\alpha = 8$. Other training hyperparameters included the batch size of $1$, the learning rate of $1e-4$ with 10 warmup steps, and an AdamW optimizer (Loshchilov and Hutter, 2019) with a cosine scheduler (Loshchilov and Hutter, 2017).

## 2.7 SFT for Transfer Learning

We used transfer learning across languages in the CUNI-DS system for English-to-Russian, transferring from English-to-Czech system.

### 2.7.1 Phase 1: en2cs Training

In the first phrase, we proceeded very similarly as described in Section 2.3. We started with the 4-bit quantized Mistral 7B v0.1 model (Jiang et al., 2023) and trained it using QLoRA (Dettmers et al., 2023) with a rank of 64 and an alpha of 128. The training followed Alpaca-like (Taori et al., 2023)

---

[3]Resulting in total of 7 checkpoints for each of the two runs.

instructions, with 20 warmup steps, a learning rate of $2e-5$, weight decay of $1e-2$, and a cumulative batch size of 32.

The model was trained on CzEng 2.0 for 24 hours, with segments packed into chunks of 2048 tokens. The final checkpoint was selected for the next phase.

### 2.7.2 Phase 2: en2ru Fine-Tuning

The model was then fine-tuned for en2ru translation using the Yandex Corpus for sentence-level data and the News Commentary v18.1 dataset for paragraph-level data. The datasets were shuffled and concatenated, and fine-tuning was conducted under the same conditions as the first stage, lasting 24 hours.

### 2.8 Genetic algorithm

For the CUNI-GA submission in English-to-Czech, we used a genetic algorithm to combine and modify n-best lists (Jon and Bojar, 2023) produced by CUNI-Transformer (at the sentence level), in the same manner as in Jon et al. (2023). We combined 5 metrics for the fitness function by a weighted average: BLEU (Papineni et al., 2002), chrF (Popović, 2015), wmt22-comet-da (Rei et al., 2022a), wmt22-cometkiwi-da (Rei et al., 2022b) and wmt23-cometkiwi-da-xl (Rei et al., 2023). The reference-based metrics use MBR decoding (Freitag et al., 2022) in place of the unknown reference.

## 3 Data

This section details the dataset used across the various training steps and language pairs.

### 3.1 SFT dataset

#### 3.1.1 English-Czech

For the first stage of the CUNI-MH training, we used the authentic part CzEng 2.0. We did not use any preprocessing, except for applying the prompt template and packing described in Appendix A.

#### 3.1.2 English-German

The CUNI-NL system was trained using the MuST-C dataset (Cattoni et al., 2021), a large multilingual corpus built from English TED Talks, containing the audio data, the English transcription of such audio, with its translation in multiple languages. Specifically, we used the en2de subset, consisting of approximately 400 hours of speech data.

During training, we randomly took 25% of the dataset, in which the input was the source transcript

itself, instead of the audio features, so that the system could know how to translate from text-only data.

We trained the system for two epochs, both checkpoints of which were then used for evaluating against the WMT23 test set.

#### 3.1.3 English-Russian

The initial phase of CUNI-DS system training (en2cs) utilized the first million segments from the CzEng 2.0 (Kocmi et al., 2020) dataset. In the second phase (en2ru), a combination of the Yandex Corpus[4] and the News Commentary v18.1[5] dataset was used, with the latter segmented into chunks of 10 sentences each.

#### 3.1.4 Translation into Low-Resource Languages of Spain

For the Translation into Low-Resource Languages of Spain task, we backtranslated the literary part (literary.txt) of the PILAR dataset (Galiano-Jiménez et al., 2024) into Spanish using Apertium (Forcada and Tyers, 2016), resulting in 230k, 25k and 24k sentence pairs for Aranese, Aragonese and Asturian, respectively. For Aranese, we also backtranslated the Aranese side of the parallel part of the corpus, while keeping the paragraphs whole up to the length of 30 sentences, resulting in 726k sentences in 4329 documents. To make use of the paragraph-level context, we employed a context-aware prompt shown in Appendix B.

### 3.2 CPO dataset

To create a dataset for CPO (Section 2.4), we need triplets: source segment, preferred output and dis-preferred output. We construct these triplets at the *paragraph level* (i.e. several sentences concatenated into a single segment) but sentence-level processing, inspired by the approach of (Xu et al., 2024b), is used in the preparation as described below.

Given a source segment, we select both preferred and dis-preferred translation from three candidates: our stage 1 output, our last year's constrained system and human reference. Our approach ensures that we still satisfy the requirements for a constrained submission.

Our CPO source segments (and their corresponding manual reference translations) are ran-

---

| Source text | Preferred translation | Dis-preferred translation |
|---|---|---|
| E6 goes further north along the west coast and through Norway to the Norwegian town Kirkenes at Barents Sea. | E6 pokračuje dále na sever podél západního pobřeží a přes Norsko do norského města Kirkenes u Barentsova moře. | E6 pokračuje dále na sever podél západního pobřeží a přes Norsko do norského města Kirkenes <u>v</u> Barentsově moři. |
| He became seriously ill in October 1914 and retired. | V říjnu 1914 vážně onemocněl a odešel do důchodu. | V říjnu 1914 ∅ onemocněl a odešel do důchodu. |
| This was published in June 1925, in a special issue of Poetry magazine. | Tato báseň byla publikována v červnu 1925 ve speciálním vydání časopisu Poetry. | Ta vyšla v červnu 1925 ve zvláštním čísle časopisu Poezie. |
| This convention has been ratified and acceded to by Ghana. | Tuto úmluvu ratifikovala a přistoupila k ní Ghana. | Tato úmluva byla ratifikována a <u>přistoupena</u> k ní Ghana. |

Table 3: Short examples from the CPO dataset. Errors (underlined) are, resp.: Kirkenes located *in* Barents Sea; missed the adverb *seriously*; and grammatically inacceptable form of passivization mentioning the subject Ghana. The third example's dis-preferred translation does not mention the detail that we are referring to a poem ("báseň"), although this fact is not explicit in the source either; other lexical variations are minor.

domly sampled documents from CzEng 2.0, a total of 47257 documents containing 200k sentences. We then used the best checkpoint from stage 1 (see model ④ in Table 2) together with our constrained model from the previous year, CUNI-DocTransformer, to generate translations for the samples.[6]

Because we want to consider the manual translation as one of the candidates for the (dis-)preferred translation, we cannot use it as the reference to select the better candidate. Therefore, we use the reference-free wmt20-comet-qe-da[7] model to rank the translations, selecting the one with the highest score as the preferred one and the one with the lowest score as the dis-preferred one.

Note that wmt20-comet-qe-da scores individual sentences, not complete paragraphs, so we do this for each sentence in the sampled dataset, while giving all preceding sentences in the corresponding document (as translated by the given system) as a context (DocCOMET, Vernikos et al., 2022).

Since this DocCOMET approach is currently not supported by the COMET project[8] for newer model architectures, such as those used by COMETKIWI22 and XCOMET, we have not tried to build the data set using these newer models.

To arrive back at paragraph-level segments for CPO, we concatenate all the sentences in each original document. The result is a dataset consisting of 47k paragraph-level triplets for CPO. Each triplet consists of the paragraph in source language and

two translations: preferred[9] and dis-preferred.[10] Due to the sentence-level selection, both preferred and dis-preferred translations may actually mix sentences from each of the three seed translations: human, our CUNI-DocTransformer and CUNI-MH Stage 1. We leave the analysis of document-level errors that arise in this process for future.

In Figure 2, we show which sentences were selected as preferred and dis-preferred. Note that this comparison is done on sentence-level, because the resulting paragraph-level examples can be composed of sentences from different sources. Interestingly, reference sentences were scored lowest by wmt20-comet-qe-da most frequently. We also show a few short examples from our dataset in Table 3. During training, the source sentences are formatted with the prompt template shown in Appendix A, similarly to how they are handled in the SFT stage Section 2.3.

We are aware that there are several potential issues with our method of preparing the dataset. First, there is a reason to be concerned about potential overfitting to a given metric (wmt20-comet-qe-da in our case) used to select the sentences. Second, our stage 1 CUNI-MH model did the translation in sentence-level fashion, potentially disregarding the relevant context. Third, we select sentences for preferred vs. dis-preferred class considering their preceding source-side context and their preceding target-side context as translated by the candidate system, not as selected so far within the document. This leaves document-level properties both in the positive and negative cases unhandled. Ideally, the preferred paragraph would avoid also any contextual errors, and for the dis-preferred paragraph, we

---

[6]For clarity, we note that we create only one CPO dataset, using translations by ④, and we apply the CPO method using this dataset three times, starting from three different models, see Table 2.

[7]https://huggingface.co/Unbabel/wmt20-comet-qe-da

[8]https://github.com/Unbabel/COMET

[9]Sometimes also called chosen or positive example.

[10]Sometimes also called rejected or negative example.

| Model | COMET22 | COMETKIWI22 | BLEU |
|---|---|---|---|
| CUNI-Transformer | 87.19 | 80.45 | 41.44 |
| CUNI-DocTransformer | 88.29 | 81.32 | 42.47 |
| CUNI-GA | 90.78 | 84.43 | 43.27 |
| GPT4-5shot | 89.36 | 82.82 | 37.76 |
| CUNI-MH | 90.21 | 83.16 | 36.52 |

Table 4: CUNI-MH's sentence-level scores on the en2cs WMT23 test set. Other systems' scores are taken from WMT23's automatic evaluation results.

| Model | COMET22 | COMETKIWI22 | BLEU |
|---|---|---|---|
| CUNI-Transformer | 81.13 | 68.24 | 42.27 |
| CUNI-DocTransformer | 83.52 | 70.69 | 43.29 |
| CUNI-GA | 86.15 | 73.56 | 43.83 |
| GPT4-5shot | 85.45 | 72.57 | 38.45 |
| CUNI-MH $k = 1$ | 87.35 | 73.30 | 37.47 |
| **CUNI-MH $k = 8$** | **87.73** | **74.82** | **35.42** |

Table 5: CUNI-MH's document-level scores on the en2cs WMT23-para test set. $k$ denotes how many sentences at most are translated together in one chunk. The CUNI-MH final submission is in bold.



Figure 2: CPO dataset - sources of preferred and dis-preferred translations.

could construct worse translations in two ways: (1) using worse individual segments, as we do, and (2) combining better or worse individual segments in a way that purposefully damages paragraph context. Fourth, because we sampled uniformly from the CzEng 2.0 documents, our final dataset actually has a large number of documents, namely 24744 out of 41835, that only consist of a single sentence. We opted for a trivial sampling because we were concerned that naive solutions aiming at having more longer documents could potentially have a negative impact on the diversity of the dataset, however this is something we would like to address in the future.

All in all, we believe that there is potential to make subsequent iterations of the dataset higher quality by alleviating some of these concerns.

### 3.3 Validation and test datasets

During training of CUNI-MH, we used the WMT22 test set as the validation data set and the WMT23 test set as the test data set. In particular, we used WMT22 when selecting the best checkpoints and hyperparameters and only used WMT23 to estimate the final performance compared to baselines.

To prepare for paragraph-level evaluation, we also concatenated all the sentences in each document to a long paragraph, creating what we call WMT22-para and WMT23-para data sets. For CUNI-GA in English-to-Czech, we did not use validation sets, we did not compare the possible configurations on validation set, we chose the parameters based on our experience. For CUNI-GA in Translation into Low-Resource Languages of Spain, we use FLORES+ validation set (NLLB Team et al., 2022).

## 4 Evaluation

### 4.1 English-Czech

We show the sentence-level metrics on the WMT23 test set for the CUNI-MH system in Table 4 and the document-level metrics on the WMT23 test set in Table 5. We used greedy decoding for this system.

Since our preliminary experiments on WMT22-

| Submission | $W_{BLEU}$ | $W_{CHRF}$ | $W_{CMT22}$ | $W_{QE22}$ | $W_{QE23-XL}$ | CHRF | BLEU | QE22 | QE23-XL | MetricX |
|---|---|---|---|---|---|---|---|---|---|---|
| CUNI-Transformer | - | - | - | - | - | 57.3 | 29.3 | X | 0.614 | 4.3 |
| CUNI-GA | 0.1 | 0.1 | 0.4 | 0.4 | 0 | 56.4 | 29.5 | 0.819 | 0.658 | - |
| CUNI-GA | 0 | 0 | 0.5 | 0.5 | 0 | 55.5 | 26.5 | 0.827 | 0.650 | - |
| **CUNI-GA** | **0** | **0** | **0.5** | **0** | **0.5** | **54.8** | **25.6** | **0.797** | **0.726** | **3.7** |

Table 6: Paragraph-level scores on WMT24 test set for the CUNI-GA submission, primary submission in bold. CUNI-Transfomer was used to produce the n-best lists which are combined and modified for the CUNI-GA submission.

| Model | COMET22 | COMETKIWI22 | BLEU |
|---|---|---|---|
| Baseline | 24.04 | 28.55 | 0.20 |
| CUNI-NL (epoch=1) | 81.07 | 77.23 | 29.61 |
| CUNI-NL (epoch=2) | 80.90 | 77.51 | 30.75 |

Table 7: CUNI-NL's sentence-level scores on the en2de WMT23 test set.

para showed that our model did not handle longer paragraphs or documents well, we used sentence-splitter from Moses[11] to split segments into sentences. We then concatenate these sentences into chunks of up $k$, which we translate together as a whole. We then concatenate all the chunks to the original segments.

By testing our model on the WMT22-para validation dataset, we chose to use $k = 8$ for our final submission to optimize for the highest COMET22 and COMETKIWI22 scores. This can also be seen in Table 5, where the model with $k = 8$ has better COMET22 and COMETKIWI22 scores than the one with $k = 1$, at the cost of worse BLEU score.

The submitted CUNI-MH system also seems to perform well according to the preliminary automatic rankings, where it surpasses most of our systems from previous years and closely matching the performance of another of our systems, CUNI-GA. These results are shown in Table 8.

However, since both systems use COMET or COMETKIWI metrics during either training or inference, raising potential concerns about overfitting, we are also awaiting the results of human evaluation (Kocmi et al., 2024).

We also tried to use CPO with our new dataset to train the base Mistral model directly, skipping the supervised fine-tuning stage. The results are shown in Table 2, see ⑦, which is the best performing checkpoint of the training run, according to its COMETKIWI22 score on the validation dataset. It can be seen that the performance of this model is significantly worse in all metrics, so the SFT stage

seems necessary in our setting.

We have also submitted CUNI-Transformer and CUNI-DocTransformer systems from previous year to provide reasonable constrained baselines for our newer models.

The CUNI-GA in this task submission combines hypotheses from CUNI-Transformer n-best lists created with beam sizes 4, 10 and 25 for each sentence. The resulting 39 translation candidates were processed by the genetic algorithm. The fitness (objective) function was a weighted combination of 5 metrics: BLEU, chrF, wmt22-comet-da (CMT22 in Table 6), wmt22-cometkiwi-da (QE22) and wmt23-cometkiwi-da-xl (QE23-XL). The weights and the obtained scores (chrF, BLEU, QE22, QE23-XL and MetricX (Juraska et al., 2023)) on the WMT24 test set are shown in Table 6. We did not use a development set due to high computational requirements of this approach, the weights are chosen based on our previous experience. An expected conclusion is that our approach allows us to easily optimize for the fitness metrics, which can be seen by comparing the QE23-XL scores of baseline translations (first row) and the score of the translations directly optimized for this metric (last row).

### 4.2 Czech-Ukrainian

We will add results for the Czech-Ukrainian submission in the camera-ready version.

### 4.3 English-German

For the CUNI-NL submission, we performed inference using the beam search algorithm, with the beam size of 2 for both checkpoints. We evaluated the performance of the two checkpoints of this system (as trained for speech translation), after epoch

---

[11]Wrapped by https://pypi.org/project/mosestokenizer/

# English-Czech

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation? |
|---|---|---|---|---|
| Unbabel-Tower70B | 1.0 | 1.8 | 0.732 | ✓ |
| Claude-3.5 § | 2.1 | 2.4 | 0.693 | ✓ |
| CUNI-MH | 2.1 | 2.3 | 0.690 | ✓ |
| CUNI-GA | 2.3 | 3.7 | 0.726 | ✓ |
| Gemini-1.5-Pro | 2.6 | 2.8 | 0.678 | ✓ |
| GPT-4 § | 2.6 | 2.9 | 0.682 | ✓ |
| IOL-Research | 2.8 | 3.0 | 0.676 | ✓ |
| ONLINE-W | 2.8 | 2.8 | 0.669 | ✓ |
| CommandR-plus § | 2.9 | 2.9 | 0.669 | ✓ |
| SCIR-MT | 3.2 | 3.3 | 0.664 | ✓ |
| TranssionMT | 3.5 | 3.5 | 0.655 | |
| ONLINE-A | 3.6 | 3.4 | 0.648 | |
| Mistral-Large § | 3.7 | 3.6 | 0.647 | |
| IKUN | 3.9 | 3.7 | 0.638 | ✓ |
| ONLINE-B | 4.0 | 3.9 | 0.640 | |
| Llama3-70B § | 4.1 | 4.0 | 0.640 | ✓ |
| Aya23 | 4.3 | 4.0 | 0.630 | ✓ |
| CUNI-DocTransformer | 4.4 | 4.0 | 0.621 | ✓ |
| IKUN-C | 4.7 | 4.3 | 0.618 | ✓ |
| CUNI-Transformer † | 4.7 | 4.3 | 0.614 | |
| ONLINE-G | 5.7 | 5.2 | 0.592 | |
| NVIDIA-NeMo † | 7.6 | 6.5 | 0.536 | |
| Phi-3-Medium § | 15.0 | 11.4 | 0.305 | |
| TSU-HITs | 19.5 | 16.6 | 0.235 | |
| CycleL2 | 24.2 | 19.5 | 0.077 | |
| CycleL | 27.0 | 22.5 | 0.031 | |

Table 8: Preliminary WMT24 General MT automatic ranking for English-Czech. **Closed systems** are highlighted with a dark gray background, **open systems** with a light gray background, and **constrained systems** are shown on a white background.

1 and after epoch 2 of en2de MuST-C corpus, with the latter performing better, so we chose it for the final evaluation against the test set this year. The results of the evaluation on the WMT23 test set are shown in Table 7.

### 4.4 English-Russian

For the CUNI-DS submission, we ran the evaluation on the paragraph level, i.e. the model needed to output the translation of the whole input at once. We used greedy decoding due to frequent emission of repeated tokens (sometimes called "spasm" by NMT practitioners) we observed with beam search. The outcomes of the CUNI-DS system's two-stage training are presented in Tables 9 and 10.

### 4.5 Translation into Low-Resource Languages of Spain

We compare Apertium and two open-source LLMs – Aya-23-8B and Command-R (35B version, quantized to 4 bits) – in translation from Spanish into the other languages of the task. We show the scores in Table 11. We fine-tuned both LLMs as a single joint model for all the languages on the back-translated literary data described in Section 3. We present BLEU, chrF and COMET-22 scores of the best-performing checkpoints after fine-tuning in Table 12. We submitted the translations produced using the Aya-23 model fine-tuned for 5000 steps. While the results are at best comparable to Apertium scores, we note that we only did a very lightweight fine-tuning on synthetic (backtranslated) data, which shows the potential of LLMs for translation into previously unsupported low-resource languages related to a language present in the training data. For instance, we obtained improvement from 46.7 to 70.2 ChrF (12.4 to 39.0 BLEU) in Aragonese by fine-tuning on 24k backtranslated sentence pairs from a different (literary) domain.

## 5 Future work

We have several ideas to improve the performance of the future iterations of our CUNI-MH model:

- Longer sequences: During our SFT stage, we trained on short sequences, mostly single sentences. In the future, we would like to experiment with training on larger sequences, so that the model is able to handle longer inputs in end-to-end fashion.

- Better CPO dataset: Our current dataset for CPO (Section 3.2) was created without including any filtering steps. The Stage 1 model we used to create one kind of translation candidates also translated in sentence-level fashion only. We think there is potential to create a higher quality dataset by using our final model, ensuring all translations are done with paragraph or document level context and possibly investigating means of filtering out lower quality examples.

- Better QLoRA initialization: During our SFT stage, we used the default initialization from the original LoRA paper (Hu et al., 2021). There are other initialization methods specifically for the combination of LoRA adapters and quantization, such as LoftQ (Li et al., 2023) which seems to consistently perform better for QLoRA. In the future, we would like to evaluate using this initialization method.

- Monolingual pretraining stage: Xu et al. (2024a) have shown promising results by including a stage where they continue pretraining Llama 2 7B and Llama 2 13B models on monolingual data covering their target languages. We think including such a stage before our SFT stage is worth considering in our future models.

- Optimization of model merging: Our experiments with checkpoint merging (Section 2.5) were extremely sparse. In the future, we would also like to evaluate SLERP and linear interpolation in comparable settings and a broader range of possible combined models (checkpoints from a single run vs. checkpoints across different run branches).

## 6 Conclusion

In this paper, we presented the CUNI submissions for the WMT24 General Translation task and the Translation into Low-Resource Languages of Spain task. Our primary focus was on using small open-source language models for various language pairs and providing comparisons with our systems from previous years.

The CUNI-MH system for English-to-Czech translation, based on Mistral 7B, showed promising results, possibly because of its CPO stage which led to a significant improvement of COMET and

| Dataset | COMET22 | COMETKIWI22 | BLEU |
|---|---|---|---|
| WMT22 | 84.24 | 78.21 | 24.30 |
| WMT23 | 75.33 | 74.81 | 21.63 |
| WMT23-para | 75.33 | 74.81 | 25.89 |

Table 9: CUNI-DS's segment-level scores for the first stage (en2cs training and en2cs evaluation) across different test datasets.

| Dataset | COMET22 | COMETKIWI22 | BLEU |
|---|---|---|---|
| WMT22 | 85.81 | 80.97 | 24.45 |
| WMT23 | 85.89 | 81.02 | 22.30 |
| WMT23-para | 72.27 | 78.21 | 21.63 |

Table 10: CUNI-DS's segment-level scores for the second stage (en2ru fine-tuning and en2ru evaluation) across different test datasets.

| Model | COMET | BLEU | chrF |
|---|---|---|---|
| **Apertium** | | | |
| Aragonese* | 0.788 | 65.3 | 82.0 |
| Aranese | 0.623 | 37.8 | 59.9 |
| Asturian | 0.652 | 16.9 | 50.6 |
| **Command-R 4-bit** | | | |
| Aragonese | 0.702 | 15.9 | 49.5 |
| Aranese | 0.576 | 4.5 | 33.3 |
| Asturian | 0.680 | 14.5 | 46.7 |
| **Aya-23** | | | |
| Aragonese | 0.685 | 12.4 | 46.7 |
| Aranese | 0.535 | 4.1 | 31.8 |
| Asturian | 0.645 | 9.0 | 40.3 |

Table 11: Scores of the baseline models on FLORES+ dev set in translation from Spanish into the given language. We note that the Aragonese part of the test set was created by post-editing Apertium translation, which is marked by the asterisk.

COMETKIWI scores, surpassing our previous systems. The model weights are available on Huggingface[12].

Our other submissions explored various techniques, such as transfer learning (CUNI-DS on en2ru), adaptation from speech translation (CUNI-NL on en2de) and creation of synthetic data using backtranslation to evaluate the feasibility of using LLMs for low-resource languages in the Translation into Low-Resource Languages of Spain task.

| Model | COMET | BLEU | chrF |
|---|---|---|---|
| **Command-R 4-bit (240)** | | | |
| Aragonese | 0.779 | 37.9 | 69.7 |
| Aranese | 0.634 | 33.1 | 57.4 |
| Asturian | 0.699 | 15.3 | 49.0 |
| **Aya-23 (5000)** | | | |
| Aragonese | 0.780 | 39.0 | 70.2 |
| Aranese | 0.632 | 35.0 | 58.1 |
| Asturian | 0.686 | 15.2 | 48.8 |

Table 12: Scores of the fine-tuned models on FLORES+ dev set in translation from Spanish into the given language. Number of fine-tuning steps in the parentheses.

# 7 Acknowledgments

---

[12]https://huggingface.co/wmt24-cuni/CUNI-MH

# References

Roldano Cattoni, Mattia Antonino Di Gangi, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. Must-c: A multilingual corpus for end-to-end speech translation. *Computer Speech and Language*, 66:101155.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.

Mikel L. Forcada and Francis M. Tyers. 2016. Apertium: a free/open source platform for machine translation and basic language technology. In *Proceedings of the 19th Annual Conference of the European Association for Machine Translation: Projects/Products*, Riga, Latvia. Baltic Journal of Modern Computing.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Pilar.

Petr Gebauer, Ondřej Bojar, Vojtěch Švandelík, and Martin Popel. 2021. CUNI systems in WMT21: Revisiting backtranslation techniques for English-Czech NMT. In *Proceedings of the Sixth Conference on Machine Translation*, pages 123–129, Online. Association for Computational Linguistics.

Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vlad Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. Arcee's mergekit: A toolkit for merging large language models. *arXiv preprint arXiv:2403.13257*.

Almog Gueta, Elad Venezian, Colin Raffel, Noam Slonim, Yoav Katz, and Leshem Choshen. 2023. Knowledge is a region in weight space for fine-tuned language models.

Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. ArXiv:2310.06825 [cs].

Josef Jon and Ondřej Bojar. 2023. Breeding machine translations: Evolutionary approach to survive and thrive in the world of automated evaluation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2191–2212, Toronto, Canada. Association for Computational Linguistics.

Josef Jon, Martin Popel, and Ondřej Bojar. 2022. CUNI-bergamot submission at WMT22 general translation task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 280–289, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Josef Jon, Martin Popel, and Ondřej Bojar. 2023. CUNI at WMT23 general translation task: MT and a genetic algorithm. In *Proceedings of the Eighth Conference on Machine Translation*, pages 119–127, Singapore. Association for Computational Linguistics.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Findings of the WMT24 general machine translation shared task: the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Martin Popel, and Ondřej Bojar. 2020. Announcing CzEng 2.0 Parallel Corpus with over 2 Gigawords. *arXiv:2007.03006*.

Yixiao Li, Yifan Yu, Chen Liang, Pengcheng He, Nikos Karampatziakis, Weizhu Chen, and Tuo Zhao. 2023. Loftq: Lora-fine-tuning-aware quantization for large language models.

Ilya Loshchilov and Frank Hutter. 2017. Sgdr: Stochastic gradient descent with warm restarts.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. ArXiv:1711.05101 [cs, math].

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti,

John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Martin Popel. 2018. CUNI transformer neural MT system for WMT18. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 482–487, Belgium, Brussels. Association for Computational Linguistics.

Martin Popel and Ondřej Bojar. 2018. Training Tips for the Transformer Model. *The Prague Bulletin of Mathematical Linguistics*, 110:43–70.

Martin Popel, Dominik Macháček, Michal Auersperger, Ondřej Bojar, and Pavel Pecina. 2019. English-Czech systems in WMT19: Document-level transformer. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 342–348, Florence, Italy. Association for Computational Linguistics.

Martin Popel, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature Communications*, 11(4381):1–15.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G.

C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ken Shoemake. 1985. Animating rotation with quaternion curves. In *Proceedings of the 12th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '85, page 245–254, New York, NY, USA. Association for Computing Machinery.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Joachim Utans. 1996. Weight averaging for neural networks and local resampling schemes. In *AAAI-96 Workshop on Integrating Multiple Learned Models*, page 133–138. AAAI Press.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.

Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 118–128, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. 2020. Trl: Transformer reinforcement learning. https://github.com/huggingface/trl.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024a. A paradigm shift in machine translation: Boosting translation performance of large language models.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024b. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation.

## A   CUNI-MH Model Prompt Template and Packing

We used the following prompt template for the model, inspired by the one used in Alpaca (Taori et al., 2023):

```
### Instruction:
Translate Input from English to Czech
### Glossary:

### Previous text:

### Input:
{source_text}
### Response:
{target_text}
```

The *Glossary* and *Previous text* sections were not used for the current task, so we left them empty. Since we trained only a single translation direction this time, the instruction remains constant.

Below is a shortened example of the packed[13] and tokenized training data, where <s> stands for the beginning of sequence token, </s> stands for the end of sequence token and \n stands for newline, the tokens are separated by spaces:

```
<s> ### Inst ruction : \n Trans late
Input from English to Czech
\n ### Gl oss ary : \n \n ### Pre
vious text : \n \n ### Input : \n It
had been bad enough , calling
Brother when she was with
him . \n ### Response : \n By lo
d ost z lé př iv ol at Br atra
, k dy ž byla s n ím . </s>
<s> ### Inst ruction : \n Trans late
Input from English to Czech
\n ### Gl oss ary : \n \n ### Pre
vious text : \n \n ### Input : \n To
do it now ? \n ### Response :
\n A le te ď ? </s> <s> ### Inst
ruction : \n Trans late Input from
English to Czech \n ### Gl oss
ary : \n \n ### Pre vious text :
\n \n ### Input : \n Here ? \n ###
Response : \n T ady ? </s>
```

---

[13]The packing itself is implemented by TRL's ConstantLengthDataset, see https://github.com/huggingface/trl/blob/e3fe28ee1a8bfab9739f849759c93d56776376e2/trl/trainer/utils.py#L431

## B CUNI-GA Model Prompt Template

We used the following prompt for context-aware translation in the Translation into Low-Resource Languages of Spain task, in order to make use of document-level context, while still keeping alignment on the sentence level, necessary for the evaluation:

```
We need to translate a single line from
conversation in Spanish into
{target_language}.  This is the
conversation: {src_context}

The start of the conversation is already
translated into English: {prev_context}
Translate the following line from
{src_lang} to {tgt_lang}.

Be very literal, and only translate the
content of the line, do not add any
explanations: {src_line}
```

# From General LLM to Translation: How we dramatically improve translation quality using human evaluation data for LLM finetuning

**Denis Elshin**  **Nikolay Karpachev**  **Alexander Antonov**  **Anton Chekashev**

**Alexander Chernyshev**  **Kirill Denisov**  **Ekaterina Enikeeva**  **Vera Frantsuzova**

**Ilya Golovanov**  **Boris Gruzdev**  **Georgy Ivanov**  **Ekaterina Latypova**

**Vladimir Layner**  **Vladislav Negodin**  **Dmitry Popov**  **Nickolay Skachkov**

## Abstract

This paper describes Yandex submission to the WMT2024 General Translation Task. More specifically, we present a novel pipeline designed to build a strong paragraph-level translation engine with an emphasis on video subtitles domain. In particular, we apply a multi-stage adaptaion pipeline on top of LLM pretraining to align the model for translation task and subsequently to the video subtitles format. Our submission ranks 3rd on the preliminary general translation leaderboard.

## 1 Introduction

In this paper, we present unconstrained system submitted by the Yandex LLC NLP team to the WMT 2024 General MT Translation track, focusing on English-to-Russian translation. Our approach involves training a YandexGPT[1] LLM-based model for translation tasks using a multi-stage process to ensure high-quality and contextually accurate translations.

We are not capable of revealing all the details of the model due to NDA reasons, however, we can say that it is a Yandex GPT-like model, specifically trained for the translation task.

Our multi-stage approach, which combines extensive pre-training, targeted fine-tuning, advanced prompt-tuning, and structure-preserving

---

[1] https://yandex.cloud/en/services/yandexgpt

techniques, ensures that our model delivers high-quality, fluent, and structurally consistent translations and performs well both in competitive benchmarks and real-world applications.

## 2 System Overview

### 2.1 Pretraining

The foundation of our approach is a robust pretraining phase involving a Large Language Model (LLM) trained on a vast corpus of clean texts in multiple languages, with a predominant focus on Russian and English. The quality of this pretrained model is evaluated using a comprehensive suite of benchmarks, including both automated metrics and human evaluation.

This initial phase ensures that the model captures a wide range of linguistic features and nuances across different languages, thereby establishing a strong base for subsequent fine-tuning.

### 2.2 Incorporating Parallel Data

Following the pretraining phase, we enhance the model by incorporating parallel data, where English and Russian texts are concatenated using a delimiter. This step is crucial for aligning the model's understanding of both languages in a translation context. We use a proprietary CommonCrawl-like parallel corpus of pages crawled from the Web. The data is meticulously curated to ensure high quality using Bicleaner-like Ramírez-Sánchez et al. (2020)

pipeline:

- Texts are selected using automated parallelism filters.

- Duplicates are removed to maintain a clean dataset.

This concatenation strategy enables the model to establish connections between two languages and to learn direct mappings from English to Russian and vice versa.

## 2.3 Sentence-level vs. Paragraph-level Translation

Our initial translation model primarily focuses on sentence-level translation. However, through extensive experimentation, we have observed that paragraph-level translation benefits significantly more from clean, coherent paragraph-level data. Unlike isolated sentences, paragraphs provide a broader context, which is essential for maintaining the flow and coherence in translations.

To leverage this, we gather texts that are inherently structured in paragraphs. These texts are preprocessed to ensure they meet our quality standards:

- Automated filters are employed to assess text parallelism and quality.

- Rigorous deduplication processes are applied to eliminate any repeated content, ensuring that the data fed into the model is both diverse and representative.

## 2.4 Structured content translation

Although the document-level translation system we have obtained using the pipeline above has high translation quality on generic textual data, it is incapable of consistently translating data in structured format, e.g. data in HTML format. Particularly, when presented texts with tags or other strict markup, model is prone to dropping or altering the markup and thus generating an invalid HTML page.

To handle this problem, we have designed a data augmentation strategy aimed at guiding the model towards HTML domain and such an augmentation have been incorporated into our document-level alignment stage.

## 2.5 Fine-Tuning LLM for Subtitle Translation

Building on a pre-trained LLM proficient in translating tagged web pages, we developed a method to train the model for subtitle translation. The key idea of this approach involves enclosing each speaker and dialogue in brackets, ensuring accurate parsing into individual dialogues.

This adaptation enhances the LLM's ability to meet the specific challenges of subtitle translation, ensuring contextually accurate outputs with proper segmentation by speaker and timing.

In the subsequent sections we further describe the main stages of our pipeline.

## 3 Supervised Fune-Tuning (SFT)

Firstly, we align the pretrained language model to the machine translation task. We conduct supervised fine-tuning (SFT) on an in-house dataset of parallel books fragments of up to 1000 tokens length.

We use multilayer prompt-tuning as in Liu et al. (2021) with each p-tuning block size of 100.

Overall LLM input consists of an English source text surrounded by two p-tuning blocks:



Figure 1: PTune blocks layout.

## 4 Human Feedback Alignment

Following the Supervised Fine-Tuning stage, we further improve core translation capabilities of the model using our internal Human Preferences dataset.

### 4.1 Data

We collect the training data using Side-By-Side human evaluation of paragraph-level translations, where an expert has to choose which of the two translations is better. The annotated data is presented in triplets (source, winner, loser), where 'winner' and 'loser' correspond to the compared translations. The source segments are sampled from various domains including books of different genres, web pages etc.
Our training dataset consists of the following parts:

**Sentence-level data**
Sentence part of the corpus consists of side-by-side comparisons between different model generations, in total 100.000 sentence triplets.

**Document-level data**

Document part of the corpus contains two primary sources of human feedback annotations.

Firstly, similarly to the sentence-level alignment data, we collect several thousands of document-level side-by-side comparisons between different versions of our model.

Secondly, we collect an additional contrastive triplet corpus aimed specifically at improving translation fluency.

Total document-level corpora size is several tens of thousands triplets.

## 4.2 Modeling

We fine-tune the model obtained at SFT stage using contrastive learning objective.
The model is trained using Contrastive Preference Optimization (CPO) loss function as in Xu et al. (2024).

$$\mathcal{L}(\pi_\theta; U) = \min_\theta \underbrace{\mathcal{L}(\pi_\theta, U)}_{\mathcal{L}_{\text{prefer}}} \underbrace{-E_{(x, y_w) \sim \mathcal{D}}[\log \pi_\theta(y_w | x)]}_{\mathcal{L}_{\text{NLL}}}.$$

(1)

where

$$\mathcal{L}_{\text{prefer}}(\pi_\theta, U) = -E_{(x, y_w, y_l) \sim \mathcal{D}} \Big[ \log \sigma \Big( \beta \log \pi_\theta(y_w | x) - \beta \log \pi_\theta(y_l | x) \Big) \Big].$$

(2)

We train with batch size of 64, 1 epoch and triangular learning rate schedule (warmup length of 0.1 epochs, peak learning rate 1e-6).
It is worth mentioning that, due to the dataset imbalance between sentences and documents, training on a uniform mixture yields results almost equal to only sentence-wise training. To handle this discrepancy between sources, we employ a variation of curriculum learning (Bengio et al. (2009)).
In particular, we implement an easy-to-hard schedule, where we start with training only on sentence-level data and shift towards longer documents to the end of the training. This enables more effective leveraging of low-resource document-level corpora.

## 5 Structured content translation

In this section, we explore the methodology developed to improve the translation of pages with structured data (e.g. web page or video subtitles data) by Large Language Models (LLMs). Traditional LLMs, when tasked with translating structured content, often exhibit significant hallucination level. This manifests as omission of tags, partial tag loss, or incorrect translation of tags. Our goal is to achieve a more robust and accurate translation of such content by ensuring the correct transfer of tags.

## 5.1 Current Challenge: Tag Hallucination

During free-form translation, LLMs struggle to maintain the integrity of HTML tags. This issue is critical as tags are essential for preserving the structure and formatting of HTML documents. A common problem observed is the complete omission of tags or their partial loss, which leads to a significant decrease in the quality of the translated document. An initial assessment showed a low percentage of correctly transferred tags. Tags are preserved only in 36% for CPO model that proves the need of a more reliable approach to tag preservation.

**Test data:** To test the accuracy of tag preservation we used a corpus of HTML-fragments. We collected innerHTML of block HTML tags from 10 Wikipedia pages.

**Proposed Solution:** Bracket Substitution and Model Adaptation
To address the issue of hallucination and improve tag preservation, we propose the following approach:

### 5.1.1 Tag Substitution with Brackets

Paired HTML tags are replaced with paired brackets (e.g., <div> becomes {, and </div> becomes }) to simplify the text structure for the model. Unpaired tags are also converted to a bracket format: every unpaired tag becomes a pair {}. This increases the proportion of sentences with retained tags to 76%.

```
a. I saw a cat.
b. <span><a>I</a> saw a <span>cat</span>.</span>
c. _< |span|>< |a|> |I| </ |a|> |_saw|_a|_< |span| > |cat| </ |span|>.</ |span|>
d. {[I] saw a [cat].}
```

Figure 2: a. Plain sentence. b. Sentence with html tags. c. Sentence with tags displayed as subword tokens processed by LLM. d. Sentence with tags replaced with braces.

### 5.1.2 Adaptation Using Parallel Corpus

We utilize a parallel corpus of HTML texts sourced from open repositories. This corpus serves as a foundation for generating synthetic data necessary for model fine-tuning.

### 5.1.3 Training Dual-Network System

**Train data:** We used the same parallel corpus as for SFT training but with tags aligned from original HTML documents. Sentence pairs with non-matching HTML tags were filtered out.
**First Network:** This network is trained to insert brackets and line breaks correctly into the text in the original language. This step helps to maintain the structural consistency of the text.
**Second Network:** Given a source text with tags and its translation without tags, this network learns to accurately re-insert the tags into the translated text. This network ensures that the translated content preserves the necessary HTML tags.

### 5.2 Synthetic Data Generation

By leveraging the dual-network system, we generate a substantial amount of synthetic data. This data includes the original text with brackets and line breaks, and the corresponding translated text with correctly inserted tags. Specifically, for the Contrastive Preference Optimization (CPO), we use:

1. The output of the first network as the source sentence in English.

2. The output of the second network on a good translation as the positive example.

3. The output of the second network on a poor translation as the negative example.

The good/poor translation pairs were obtained using human annotation as described above.

### 5.3 Results

Our experimental results demonstrate that the proposed methodology effectively increases the percentage of sentences with correctly transferred tags to 99%.
This substantial improvement underscores the effectiveness of our approach in reducing tag hallu-cination and ensuring a more stable and accurate translation of HTML content.
By substituting HTML tags with brackets, adapting the model using a parallel HTML corpus, and incorporating a dual-network system for synthetic data generation, we have developed a robust method to enhance HTML translation. This approach not only mitigates the problem of tag hallucination but also ensures the structural integrity of translated HTML documents. The success of this methodology paves the way for more reliable and efficient translation of structured data formats, significantly benefiting applications in web content translation and beyond.

## 6 Fine-Tuning LLM for Subtitle Translation

Building upon a pre-trained model that has demonstrated proficiency in translating tagged web pages, we have adapted the following approach to train a subtitles translation system. Its core idea is straightforward: we enclose each speaker and their corresponding dialogue in brackets, as shown in figure below.



```
[
{man 1:} {
    Hey, guys, Kevin here from snowboard pro camp, in this video I'm
    going to give you a list of the first ten tricks to learn on your
    snowboard.
}
{woman 1:} {
    These tricks are in order and each trick will teach you a skill that
    you'll use in the next trick on the list.
}
]
```

Figure 3: Subtitles input format.

This ensures that the translation preserves these brackets, allowing the entire text to be parsed into individual speaker dialogues.
The production version of the algorithm is somewhat more sophisticated, as it must align the translations of longer dialogues with their corresponding timestamps. However, for the purposes of this discussion, a more detailed description is unnecessary and is therefore omitted from this paper.
We fine-tune the model using publicly available subtitle corpora, which we preprocess to fit the above mentioned format. This additional training step has led to noticeable improvements in our human evaluation scores, particularly within the domain of movies and YouTube video subtitle translation. The reason for employing this model is that part of the competition data is presented in audio format, making effective subtitle translation a critical component of our approach.

By adapting the LLM in this manner, we enhance its ability to handle the unique challenges posed by subtitle translation, ensuring that the final outputs are both contextually accurate and properly segmented according to speaker and timing, which is crucial for maintaining the integrity of the original content in the translated version.

## 7  Evaluation Metrics and Results

**Ablation**

In order to estimate the effect of each stage of the pipeline, we compare our models using BLEURT-20 (Sellam et al. (2020)) and COMET (Rei et al. (2020)) automatic metrics, as well as BLEU. We rely primarily on neural metrics results as suggested in Freitag et al. (2022). Table 1 shows the scores on WMT-22 English to Russian testset.

| Model Ablation (wmt-22 fwd) | | | |
|---|---|---|---|
| Model Stage | BLEURT | COMET | BLEU |
| PTune | 0.76 | 0.836 | 31.3 |
| cpo-sents | **0.789** | **0.860** | **31.52** |
| cpo-curriculum-base | 0.787 | 0.855 | 24.8 |
| cpo-curriculum-tags | 0.784 | 0.855 | 27.1 |

Table 1: Metrics by stage (sentence-level).

| Model Ablation (wmt-22 fwd news) | | | |
|---|---|---|---|
| Model Stage | BLEURT | COMET | BLEU |
| PTune | 0.728 | 0.835 | **27.78** |
| cpo-sents | 0.733 | 0.847 | 25.55 |
| cpo-curriculum-base | **0.743** | **0.850** | 19.61 |

Table 2: Metrics by stage (document-level).

Firstly, the model trained only on parallel data (PTune) is already capable of generating decent quality translations. However, it exhibits bias towards literal translations and poor fluency.

During the alignment stage (cpo-curriculum-base) the model is exposed to a variety of high-quality translations (including contrastive triplets aimed specifically at improving fluency) and, hence, the model after initial CPO training is much more fluent, but prone to tags omission and format inconsistency.

Augmented CPO training solves the problem with format and tags without sacrificing the model's target language fluency capabilities.

Overall, the metrics ablation highlights the following:

1) BLEU correlates poorly with model quality, especially on document-level benchmarks due to high preference for literal translations.

2) On sentence-level evaluation contrastive learning model trained only on sentence data yields superior results both on neural and n-gram based metrics.

3) Tag-focused augmentation does not lead to quality degeneration on primary benchmarks whilst increasing model stability (see tag accuracy evaluations).

4) Contrastive learning phase with curriculum learning training improves the quality on document-level inputs, but only on neural metrics. We hypothesize that curriculum learning model increases fluency of the translations and introduces more complicated paraphrases that BLEU fails to score adequately.

**WMT'24 Results**

The quality of our system is assessed by the organizers using the following metrics:

MetricX-23-XL (Juraska et al. (2023)) – a reference-based metric built on top of the mT5 model. CometKiwi-DA-XL (Rei et al. (2023)) – a quality estimation metric built on the XLM-R XL model. Both metrics are among the top-performing metrics in the field (Freitag et al. (2023)). According to these metrics, our system currently ranks third on the leaderboard, with a MetricX score of 2.9 and a CometKiwi-DA-XL score of 0.705. The final leaderboard will be determined based on human evaluation results.

**Ethics Statement**

Our system was trained on the publicly available data. This unrestricted access to data allowed us to leverage a vast and diverse set of examples, enabling the model to learn from a wide array of linguistic patterns, contexts, and domains.

The absence of data limitations contributed to the development of a robust and versatile model, capable of generalizing well across various tasks and applications. By incorporating extensive datasets

from different sources, our system gained the ability to handle complex and varied scenarios, enhancing its overall performance and adaptability. This approach ensured that the model could effectively capture and respond to the nuances of different data types, ultimately leading to more accurate and reliable outputs in real-world applications.

## Acknowledgements

## References

Y. Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. volume 60, page 6.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André Martins. 2022. Results of wmt22 metrics shared task: Stop using bleu - neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation*, pages 46–68, Abu Dhabi.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *ArXiv*, abs/2110.07602.

Gema Ramírez-Sánchez, Jaume Zaragoza-Bernabeu, Marta Bañón, and Sergio Ortiz-Rojas. 2020. Bifixer and bicleaner: two open-source tools to clean your parallel data. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 291–298, Lisboa, Portugal. European Association for Machine Translation.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. pages 2685–2702.

Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. In *Proceedings of ACL*.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *arXiv preprint arXiv:2401.08417*.

# Cogs in a Machine, Doing What They're Meant to Do
# – The AMI Submission to the WMT24 General Translation Task

**Atli Jasonarson, Hinrik Hafsteinsson, Bjarki Ármannsson, Steinþór Steingrímsson**

The Árni Magnússon Institute for Icelandic Studies

Reykjavík, Iceland

`atli.jasonarson,hinrik.hafsteinsson,bjarki.armannsson,`
`steinthor.steingrimsson@arnastofnun.is`

## Abstract

This paper presents the submission of the Árni Magnusson Institute's team to the WMT24 General translation task. We work on the English→Icelandic translation direction. Our system comprises four translation models and a grammar correction model. For training our models we carefully curate our datasets, aggressively filtering out sentence pairs that may detrimentally affect the quality of our system's output. Some of our data are collected from human translations and some are synthetically generated. A part of the synthetic data is generated using an LLM, and we find that it increases the translation capability of our system significantly.

## 1 Introduction

We describe our submission to the 2024 WMT general translation task. Large Language Models (LLMs) have become near-ubiquitous in the field of Natural Language Processing (NLP) in the last couple of years. They have shown remarkable translation capabilities (see e.g. Xu et al., 2024a), but require significantly larger computational resources than previous neural MT (NMT) models, both for training and inference. Most openly available LLMs are primarily trained on English texts and may therefore need further training in order to be able to translate from or into less-resourced languages, such as Icelandic.

The ALMA models (Xu et al., 2024a) are LLM-based translation models, built on LLaMA-2. They have been trained to translate ten directions, including English↔Icelandic. We explore the capabilities of some of these models, the 7B and 13B parameter versions of ALMA-R (Xu et al., 2024b), and find that they generate very competitive translations as measured against the English–Icelandic WMT21 test sets (Akhbardeh et al., 2021), especially from Icelandic into English. Unfortunately, using our settings the translation speed was quite

slow (approximately one sentence per second) on an NVIDIA A100 GPU card.

We are interested in building faster models so we use the more traditional encoder-decoder Transformer architecture described in Vaswani et al. (2017). We collect all parallel data available to us for our language pair, generate additional synthetic pairs using the ALMA-R 13B parameter model and apply iterative back-translation using our own models. We apply filters to remove sentence pairs that may have detrimental effects on the models output.

We train four Transformer models[1] of varying sizes and let each model generate five translation candidates. A spelling and grammar checking model is then applied to the translations to generate "corrected" versions of the sentences. Finally the best candidate is selected from the pool of translations, corrected or not, using a reranking model.

We evaluate our models and approaches on the WMT21 test set for English→Icelandic.

## 2 Related Work

We only submit a system for the English→Icelandic translation direction. This language pair was previously one of the pairs for the WMT General Translation shared task in 2021 but prior to that, limited work had been published on MT for Icelandic. Brandt et al. (2011) describe a rule-based system for translating Icelandic→English, based on Apertium (Forcada et al., 2011). Jónsson et al. (2020) was the first published work describing SMT and NMT for Icelandic. Since 2021 the WMT21 evaluation data, as well as various parallel corpora projects, have made it more accessible to train and evaluate MT systems translating to or from Icelandic, and with that the language has been included in various research projects. We believe this is an indicator of the importance of evaluation campaigns, such

---

[1]Models available at `https://huggingface.co/arnastofnun`.

as the ones run in association with the WMT conferences, for less prominent languages.

Our approach uses an ensemble of four different translation models and a reranking model to select the best candidate. This is a common approach, motivated by the intuition that different systems may have different strengths. In recent work, Toral et al. (2023) use this approach in their experiments with literary translations. In their work on bidirectional reranking, Imamura and Sumita (2017) discuss reranking and ensembling for MT in some detail. Examples from the period of statistical MT include the work of Olteanu et al. (2006) and Wang et al. (2007), describing language model-based reranking on hypotheses generated by phrase-based SMT systems.

## 3 Data Selection and Filtering

Various parallel data are available for the English–Icelandic language pair. ParIce (Barkarson and Steingrímsson, 2019) is partly a collection of parallel corpora available elsewhere, which has been realigned and refiltered, and partly data compiled for that project, the largest source being regulatory texts published in relation with the European Economic Area (EEA) agreement. Data for the English–Icelandic language pair were collected within the Paracrawl project (Bañón et al., 2020), CCMatrix (Schwenk et al., 2021), MaCoCu (Bañón et al., 2022) and HPLT (Aulamo et al., 2023). Data for the language pair are also available from multiple smaller datasets distributed on OPUS (Tiedemann and Thottingal, 2020). We utilize all these datasets in training our models.

We also use synthetic data: Backtranslations made available by Jónsson et al. (2022), translations generated using the ALMA-R 13B parameter model and backtranslations generated by our trained models. We describe these in more detail in Section 3.3.

Khayrallah and Koehn (2018) show that incorrect translations, untranslated target text, misalignments, and other noisy segments in training data can have a detrimental effect on the quality of translations generated by NMT systems trained on that data. By filtering our training data rather aggressively, we try to minimize such noise.

### 3.1 ParIce

Even though care has been taken to realign and refilter data for the ParIce corpus, Steingrímsson et al.

(2023) show that it still contains noise, such as misalignments and mistranslations, that may be detrimental when training NMT systems. They refilter the data using a combination of approaches: Shallow filters based on simple heuristics, by using Bicleaner (Sánchez-Cartagena et al., 2018; Ramírez-Sánchez et al., 2020) and by employing classifiers (support vector machine-based ones (Cortes and Vapnik, 1995) had the best outcome) with a combination of scoring mechanisms, including LASER (Artetxe and Schwenk, 2019), LaBSE (Feng et al., 2022), NMTScore (Vamvas and Sennrich, 2022) using the M2M100 multilingual translation model (Fan et al., 2021), and WAScore, a word alignment-based score devised to measure word-level parallelism, introduced in Steingrímsson et al. (2021). In Steingrímsson (2023) these data are processed further by realigning the EEA texts in the ParIce corpus using SentAlign (Steingrímsson et al., 2023).

As the basis for our training we use the ParIce dataset, processed as described above, as well as parallel data extracted from Wikipedia using the comparable corpora mining approach described in (Steingrímsson et al., 2021) and sentence pairs extracted from version 9 of Paracrawl using the filtering approaches described above and in Steingrímsson et al. (2023).

### 3.2 Filtering the OPUS Datasets

An overview of the data for Icelandic-English parallel texts sourced from the OPUS catalog is provided in Appendix A. This data, accounting for redundant sentence pairs, amounts to 21.167.708[2] sentence pairs. At face value, this is a substantial amount of available data. However, the quality of these parallel texts is not reliable, with noisy and incorrect pairs being prevalent throughout most individual datasets in the catalog. To remedy this, and thus ensure that the data sourced via OPUS can be used effectively in our project, we applied an aggressive, sequential filtering process, with the goal of whittling away the majority of the low-quality sentence pairs.

Our sequential filtering process consists of ten individual steps, most of which only remove sentences from the data without modifying the content of other sentences. The process is *sequential*, in that the input of a filtering step is the output of the previous filtering step. Furthermore, the order of

---

[2]This applies to the state of the OPUS catalog at the time of development, i.e., April 2024.

Figure 1: Each filtering step's effect on OPUS dataset size

these steps is decided to ensure optimal processing time of the filters so that computationally heavy filtering steps process the least amount of data, which minimizes run time. For a detailed overview of each filtering step, see Appendix B.

The effects of each filtering step on the data amount is shown in Fig. 1. To ensure that our filtering methods affected our implementation positively, we intermittently added the output of the filtering process to our training pipeline and evaluated the performance. In particular, we used this approach to dial in the optimal LaBSE and NMT score cutoffs in our filters.

The final output of our filtering process produces a relatively high-quality data set of 2.056.704 English-Icelandic sentence pairs (roughly 9.71% of the original 21.167.708 raw sentence pairs sourced from the OPUS catalog), which we then add to our training data.

### 3.3 Synthetic Data

The dataset made available by Jónsson et al. (2022) contains translations from Europarl, Newscrawl, Wikipedia and the IGC. We perform a filtering step similar to the one used applied on the OPUS data, consisting of a length filter, removing all sentences that have fewer than four word tokens and more than 150, an overlap filter, removing all sentence pairs that share 40% or more of word tokens, and

a symbol filter removing all sentence pairs where more than 20% of characters in one of the sentences is non-alphabetical. Furthermore we use two scoring mechanisms for filtering, LaBSE, using a score threshold of 0.8, and NMTScore with a threshold of 0.4. These scores are selected based on the evaluation in (Steingrímsson et al., 2023). After filtering, we are left with 4.4M sentence pairs from this dataset.

We use the 13B parameter ALMA-R model to translate English sentences from Newscrawl to Icelandic and Icelandic texts from the Icelandic Gigaword Corpus (IGC) (Steingrímsson et al., 2018) to English. The Icelandic texts are sampled from three different subcorpora of the IGC, comprising news, scholarly journals, and literary texts. For each source sentence we generate five translations and use LaBSE to select the two best ones, granted that they exceed a threshold of a LaBSE score of 0.8 and pass through the three shallow filters described above: length, overlap and symbol filters. Our final set contains 8.9M sentence pairs translated from Icelandic to English and 700K sentence pairs translated from English to Icelandic.

Finally, we do iterative back-translation. We use the same training data as described above to train models to translate texts from the IGC to English. For the back-translations we use Transformer$_{\text{BIG}}$

| model | $d_{model}$ | $d_{ff}$ | $h$ | $N_{enc}$ | $N_{dec}$ |
|---|---|---|---|---|---|
| $Base$ | 512 | 2048 | 8 | 6 | 6 |
| $Base_{deep}$ | 512 | 2048 | 8 | 36 | 12 |
| $Big$ | 1024 | 4096 | 16 | 6 | 6 |
| $Big_{deep}$ | 1024 | 4096 | 16 | 36 | 12 |

Table 1: Model dimensions, heads and number of layers.

models (Vaswani et al., 2017), as described in Table 1. We use the same approach as before, generate five translations for each sentence and use LaBSE to select the two best ones, as long as they exceed the threshold of 0.8 and are not filtered out by the other filters. We do two iterations of translating and training models in both translation directions using backtranslated data. This results in a total of approximately 60M sentence pairs.

### 3.4 Other Data

To decide which datasets to use, we trained Transformer$_{BASE}$ models as described in Vaswani et al. (2017) and evaluated the models using the test set from WMT21. We started by training a baseline system using the dataset described in Section 3.1. We then added different datasets to the baseline data, trained new systems and evaluated them. If the new dataset seemed to improve the output we used that for our final system. In addition to previously described datasets we tried generating backtranslations using SMT and to add data from a bilingual lexicon using token-pair training as described by Jones et al. (2023). Table 2 shows chrF scores (Popović, 2015) for our different experi-

| Dataset | chrF |
|---|---|
| Baseline | 50.4 |
| Baseline+lexicon | 50.4 |
| Baseline+OPUS | 53.7 |
| Baseline+Jónsson | 53.5 |
| Baseline+Jónsson+SMT | 53.2 |
| Baseline+Jónsson+ALMA | 54.7 |
| Baseline+Jónsson+ALMA+OPUS | 55.1 |
| Baseline+Jónsson+ALMA+OPUS+BT1 | 56.4 |
| Baseline+Jónsson+ALMA+OPUS+BT2 | 56.8 |

Table 2: The table shows that when most of the datasets in our experiments are added to the training data the quality, as measured by chrF, increases. Exceptions to that are the experiments with adding token-pairs from an English-Icelandic lexicon and with using backtranslations generated by an SMT system. These two datasets are therefore not used in our final systems.

| Dataset | Sentence Pairs |
|---|---|
| Base | 2,277,023 |
| OPUS-filtered | 2,056,704 |
| Miðeind-BT | 2,559,806 |
| Miðeind-FT | 1,837,945 |
| ALMA-BT | 8,927,720 |
| ALMA-FT | 700,253 |
| IGC-BT-1 | 27,794,398 |
| IGC-BT-2 | 33,465,175 |

Table 3: Datasets used for training and number of sentence pairs in each dataset.

iments.

The total number of sentence pairs used for training is shown in Table 3

## 4 System Description

Our motivation for using multiple models is twofold: First, we want to use models that are computationally inexpensive to run and so we train models that can run on one consumer grade GPU. Second, systems of different sizes may have complementary strengths and so training multiple systems and reranking the results may give us better results than any one model.

We train four encoder-decoder Transformer models, all of which play a part in the translation pipeline. Two of the models follow the exact architecture described in Vaswani et al. (2017), i.e. the 'base' and 'big' versions of the original Transformer model, while the other two are deeper, using 36 encoder layers and 12 decoder layers instead of six. The difference between the four models is shown in Table 1.

The outputs from the translation models undergo two post-processing steps. First, they are run through a grammatical error correction model, a version of the byte-level sequence-to-sequence model ByT5 (Xue et al., 2022) that has been fine-tuned by Ingólfsdóttir et al. (2023) to correct spelling errors in Icelandic as well as handling more complex grammatical, semantic and stylistic issues. Second, we fix punctuation errors which translation models are prone to making when translating into Icelandic (mostly to do with quotation marks, which are different in Icelandic and English) as well as some that might be unique to our system, such as their incapability to translate emojis. As the grammatical error correction model proved too aggressive for our purposes, merging and splitting

| model | chrF |
|---|---|
| $Base$ | 56.8 |
| $Base_{deep}$ | 57.1 |
| $Big$ | 57.7 |
| $Big_{deep}$ | 57.7 |
| Ensemble+COMETKIWI | 58.3 |
| Ensemble+error correction +COMETKIWI | **58.4** |
| ALMA-R 7B | 52.2 |
| ALMA-R 13B | 53.4 |

Table 4: chrF scores for each of our models, compared with scores for the model ensembles and for the ALMA-R models. The scores are calculated on the WMT21 evaluation set.

| model | Selected | Unique |
|---|---|---|
| $Base$ | 293 | 158 |
| $Base_{deep}$ | 347 | 186 |
| $Big$ | 287 | 163 |
| $Big_{deep}$ | 419 | 246 |

Table 5: The number of sentences generated by each model selected for the final output when translating the WMT21 test set.

some sentences, normalizing informal language usage and hashtags, etc., we also revert some of the changes it introduced.

Using the WMT21 test set we experiment with an ensemble approach, using COMETKIWI-DA-22 (Rei et al., 2022) to select the best sentence out of 20 hypotheses made by the four models (each model generates five hypotheses using beam search with beam size 12). This raises the chrF score to 58.3 for our evaluation set. On top of this we add the spelling and grammar error correction, which gives us a very modest increase in quality as measured by chrF, shown in Table 4.

We investigate whether the COMETKIWI-DA-22 model prefers the output from some of the translation models over the others. Table 5 shows which translation models generated the translations ultimately chosen by the scoring model when experimenting on the WMT21 evaluation set of 1000 sentences. While translations by the deeper model are more likely to be selected, it is evident that all models are contributing, with the final selection containing 753 translation generated by only one model, and of these all models contribute over 150 translations each. 247 of the selected translations were generated by more than one model (non-unique translations). An ensemble approach thus seems to be likely to improve overall translation quality.

### 4.1 The pipeline

Basing our system on the most succesful approach in our experiments, our translation pipeline consists of three steps: First, using each of our four models, we generate five translation hypotheses using beam search for all source paragraphs, resulting in a total of 20 candidates.

Furthermore, each paragraph is segmented into sentences, $s_1, \ldots, s_n$. For each sentence, every model produces five hypotheses. These hypotheses are evaluated using COMETKIWI-DA-22, and the highest-scoring hypothesis is selected for each sentence. The selected hypotheses are concatenated to form a new paragraph. Finally, a single paragraph is created by combining the best translation of each sentence, leaving us with 25 translation candidates.

Each of these candidates is then corrected with regard to grammar, spelling and style using the ByT5 model described above.

These two steps, translating the source text and correcting the translations, result in a total of 50 translation candidates. In order to find the best candidate we use COMETKIWI-DA-22 to score all candidates. The highest scoring one is the selected translation of our system.

## 5 Results

We evaluate our system on the test data from WMT21. As expected, the bigger models perform better, but the best results are achieved by selecting translations from an ensemble of differently trained Transformer models. We use COMETKIWI-DA-22 to select the best translation out of 20 hypotheses made by the four models, five hypotheses by each using beam search with beam size 12. This raises the chrF score to 58.3 and when we add error correction on top, the score is slightly higher, 58.4, as shown in Table 4.

In the WMT24 general translation task, systems were evaluated using two automatic metrics, MetricX-23-XL (Juraska et al., 2023) and COMETKIWI-DA-XL (Rei et al., 2023), as well as by human evaluation. According to the automatic metrics, reported in Kocmi et al. (2024), our model is competitive among the open systems, although four closed systems achieve better scores. Results

| System Name | Type | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ |
|---|---|---|---|---|
| Unbabel-Tower70B | Closed | 1.0 | 2.5 | 0.740 |
| Claude-3.5 | Closed | 2.3 | 3.6 | 0.697 |
| Dubformer | Closed | 2.5 | 3.4 | 0.685 |
| IKUN | Open | 3.2 | 4.3 | 0.666 |
| GPT-4 | Closed | 3.4 | 4.7 | 0.673 |
| **AMI** | **Open** | **3.7** | **4.9** | **0.663** |
| IKUN-C | Constrained | 3.7 | 4.9 | 0.657 |
| TranssionMT | Closed | 4.2 | 5.5 | 0.653 |
| ONLINE-B | Closed | 4.2 | 5.5 | 0.652 |
| IOL-Research | Open | 4.3 | 5.7 | 0.655 |
| ONLINE-A | Closed | 5.5 | 6.4 | 0.603 |
| Llama3-70B | Open | 6.7 | 8.0 | 0.586 |
| ONLINE-G | Closed | 6.9 | 7.9 | 0.573 |
| CommandR-plus | Closed | 9.8 | 10.6 | 0.487 |
| Mistral-Large | Closed | 10.4 | 10.9 | 0.465 |
| Aya23 | Open | 15.2 | 14.9 | 0.311 |
| Phi-3-Medium | Closed | 16.2 | 15.7 | 0.278 |
| ONLINE-W | Closed | 18.1 | 19.5 | 0.296 |
| TSU-HITs | Constrained | 19.2 | 18.4 | 0.192 |
| CycleL | Constrained | 21.0 | 20.2 | 0.148 |

Table 6: Preliminary WMT24 General MT automatic ranking for English-Icelandic. Our system is in bold.

for the automatic metrics are shown in Table 6.

# 6 Conclusions and Future Work

We show that while Large Language Models have become nearly ubiquitous in Natural Language Processing, traditional encoder-decoder Transformer models remain a viable approach to machine translation, particularly when computational efficiency is a priority.

Nevertheless, our findings also reveal that integrating LLMs can be advantageous during the training process. Specifically, ALMA-R 13B proved to be an important part of our training pipeline, as the synthetic data it generated increased the quality of our translation systems.

Furthermore, our results indicate that while more training data usually result in a better translation system, low-quality data, such as the backtranslations generated with an SMT system, can have a detrimental impact on performance. Similarly, our experiments with a bilingual lexicon using token-pair training negatively affected the system's output. This may be due to a variety of reasons. Our SMT system could probably be improved as well as our approach to include data from a bilingual lexicon in the training data. This warrants further investigation.

Our filtering method, as described in Sections 3.2, 3.3 and Appendix B, has proven effective, even though it may be argued that it is still somewhat crude and more work into minimizing the loss of useful sentence pairs and more effectively remove detrimental sentence pairs would very likely improve the training data and in turn the translation models. For example, while we use LaBSE, LASER and NMT to evaluate sentence pairs, we apply individual cutoff values for each score. A better approach could entail using a classifier to combine all metrics for an optimal result.

Although currently impractical at production-scale, genetic algorithms, as shown by Jon and Bojar (2023) and Jon et al. (2023), show promising results in generating translation candidates. Given larger computational resources, similar approaches might prove useful and await future study.

# References

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai,

Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88, Online. Association for Computational Linguistics.

Mikel Artetxe and Holger Schwenk. 2019. Massively Multilingual Sentence Embeddings for Zero-Shot Cross-Lingual Transfer and Beyond. *Transactions of the Association for Computational Linguistics*, 7:597–610.

Mikko Aulamo, Nikolay Bogoychev, Shaoxiong Ji, Graeme Nail, Gema Ramírez-Sánchez, Jörg Tiedemann, Jelmer van der Linde, and Jaume Zaragoza. 2023. HPLT: High performance language technologies. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 517–518, Tampere, Finland. European Association for Machine Translation.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-Scale Acquisition of Parallel Corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

Marta Bañón, Miquel Esplà-Gomis, Mikel L. Forcada, Cristian García-Romero, Taja Kuzman, Nikola Ljubešić, Rik van Noord, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Peter Rupnik, Vít Suchomel, Antonio Toral, Tobias van der Werff, and Jaume Zaragoza. 2022. MaCoCu: Massive collection and curation of monolingual and bilingual data: focus on under-resourced languages. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 303–304, Ghent, Belgium. European Association for Machine Translation.

Starkaður Barkarson and Steinþór Steingrímsson. 2019. Compiling and Filtering ParIce: An English-Icelandic Parallel Corpus. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, pages 140–145, Turku, Finland. Linköping University Electronic Press.

Martha Dís Brandt, Hrafh Loftsson, Hlynur Sigurþórsson, and Francis M. Tyers. 2011. Apertium-IceNLP: A rule-based Icelandic to English machine translation system. In *Proceedings of the 15th Annual conference of the European Association for Machine Translation*, pages 217–224, Leuven, Belgium. European Association for Machine Translation.

Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20(3):273–297.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2021. Beyond English-Centric Multilingual Machine Translation. *J. Mach. Learn. Res.*, 22(1).

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT Sentence Embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Mikel Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis Tyers. 2011. Apertium: A free/open-source platform for rule-based machine translation. *Machine Translation*, 25:127–144.

Kenji Imamura and Eiichiro Sumita. 2017. Ensemble and reranking: Using multiple models in the NICT-2 neural machine translation system at WAT2017. In *Proceedings of the 4th Workshop on Asian Translation (WAT2017)*, pages 127–134, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Svanhvít Lilja Ingólfsdóttir, Petur Ragnarsson, Haukur Jónsson, Haukur Simonarson, Vilhjalmur Thorsteinsson, and Vésteinn Snæbjarnarson. 2023. Byte-level grammatical error correction using synthetic and curated corpora. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7299–7316, Toronto, Canada. Association for Computational Linguistics.

Josef Jon and Ondřej Bojar. 2023. Breeding machine translations: Evolutionary approach to survive and thrive in the world of automated evaluation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2191–2212, Toronto, Canada. Association for Computational Linguistics.

Josef Jon, Martin Popel, and Ondřej Bojar. 2023. CUNI at WMT23 general translation task: MT and a genetic algorithm. In *Proceedings of the Eighth Conference on Machine Translation*, pages 119–127, Singapore. Association for Computational Linguistics.

Alexander Jones, Isaac Caswell, Orhan Firat, and Ishank Saxena. 2023. GATITOS: Using a new multilingual lexicon for low-resource machine translation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 371–405, Singapore. Association for Computational Linguistics.

Haukur Páll Jónsson, Haukur Barri Símonarson, Pétur Orri Ragnarsson, Svanhvít Lilja Ingólfsdóttir, Vilhjálmur Þorsteinsson, and Vésteinn Snæbjarnarson. 2022. Long context synthetic translation pairs for english and icelandic (22.09). CLARIN-IS.

Haukur Páll Jónsson, Haukur Barri Símonarson, Vésteinn Snæbjarnarson, Steinþór Steingrímsson, and Hrafn Loftsson. 2020. Experimenting with Different Machine Translation Models in Medium-Resource Settings. In *Text, Speech, and Dialogue - 23rd International Conference, TSD 2020, Brno, Czech Republic, September 8-11, 2020, Proceedings*, volume 12284 of *Lecture Notes in Computer Science*, pages 95–103. Springer.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Huda Khayrallah and Philipp Koehn. 2018. On the Impact of Various Types of Noise on Neural Machine Translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83, Melbourne, Australia. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Preliminary WMT24 Ranking of General MT Systems and LLMs. *ArXiv*, abs/2407.19884.

Shuyo Nakatani. 2010. Language detection library for java.

Marian Olteanu, Pasin Suriyentrakorn, and Dan Moldovan. 2006. Language models and reranking for machine translation. In *Proceedings on the Workshop on Statistical Machine Translation*, pages 150–153, New York City. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Dario Amodei, Daniela Amodei, Jack Clark, Miles Brundage, and Ilya Sutskever. 2019. Better language models and their implications. *OpenAI blog*, 1(2).

Gema Ramírez-Sánchez, Jaume Zaragoza-Bernabeu, Marta Bañón, and Sergio Ortiz Rojas. 2020. Bifixer and Bicleaner: two open-source tools to clean your parallel data. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 291–298, Lisboa, Portugal. European Association for Machine Translation.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Víctor M. Sánchez-Cartagena, Marta Bañón, Sergio Ortiz-Rojas, and Gema Ramírez. 2018. Prompsit's submission to WMT 2018 Parallel Corpus Filtering shared task. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 955–962, Belgium, Brussels. Association for Computational Linguistics.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021. CCMatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500, Online. Association for Computational Linguistics.

Peter M. Stahl. 2024. Lingua - an accurate natural language detection library for short and mixed-language text. https://github.com/pemistahl/lingua-py. Accessed: 2024-08-21.

Steinþór Steingrímsson, Sigrún Helgadóttir, Eiríkur Rögnvaldsson, Starkaður Barkarson, and Jón Guðnason. 2018. Risamálheild: A Very Large Icelandic Text Corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*, LREC 2018, pages 4361–4366, Miyazaki, Japan.

Steinþór Steingrímsson, Hrafn Loftsson, and Andy Way. 2023. Filtering matters: Experiments in filtering training sets for machine translation. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 588–600, Tórshavn, Faroe Islands. University of Tartu Library.

Steinþór Steingrímsson, Pintu Lohar, Hrafn Loftsson, and Andy Way. 2021. Effective Bitext Extraction From Comparable Corpora Using a Combination of Three Different Approaches. In *Proceedings of the 14th Workshop on Building and Using Comparable Corpora (BUCC 2021)*, pages 8–17, Online (Virtual Mode). INCOMA Ltd.

Steinþór Steingrímsson. 2023. *Effectively compiling parallel corpora for machine translation in resource-scarce conditions*. Ph.D. thesis, Reykjavik University.

Steinþór Steingrímsson, Hrafn Loftsson, and Andy Way. 2023. SentAlign: Accurate and scalable sentence alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 256–263, Singapore. Association for Computational Linguistics.

Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-MT – building open translation services for the world. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 479–480, Lisboa, Portugal. European Association for Machine Translation.

Antonio Toral, Andreas Cranenburgh, and Tia Nutters. 2023. Literary-adapted machine translation in a well-resourced language pair. In Andrew Rothwell, Andy Way, and Roy Youdale, editors, *Computer-Assisted Literary Translation*, pages 27–52. Routledge.

Jannis Vamvas and Rico Sennrich. 2022. NMTScore: A multilingual analysis of translation-based text similarity measures. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 198–213, Abu Dhabi, United Arab Emirates.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, pages 5999–6009, Long Beach, California.

Wen Wang, Andreas Stolcke, and Jing Zheng. 2007. Reranking machine translation hypotheses with structured and web-based language models. In *2007 IEEE Workshop on Automatic Speech Recognition Understanding (ASRU)*, pages 159–164.

Titus Wormer. 2024. Franc - a natural language detection library. https://github.com/wooorm/franc. Accessed: 2024-08-21.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024a. A Paradigm Shift in Machine Translation: Boosting Translation Performance of Large Language Models. In *The Twelfth International Conference on Learning Representations*.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024b. Contrastive Preference Optimization: Pushing the Boundaries of LLM Performance in Machine Translation. *ArXiv*, abs/2401.08417.

Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. ByT5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306.

## A  OPUS Texts

The parallel texts we sourced from the OPUS catalog are listed in this section. The format of the list is as follows:

*Index*. **Name**; *version*; sentence pairs

For brevity, the *ELRC* parallel text names are abbreviated after the first entry in the list, with the *ditto* symbol ('"') replacing the 'ELRC' part of the name.

1. **CCAligned**; *v1*;                1,192,542
2. **CCMatrix**; *v1*;                8,723,145
3. **ECDC**; *v2016-03-16*;                2,512
4. **ELRC-2718-EMEA**; *v1*;                542,624
5. **"-3206-antibiotic**; *v1*;                816
6. **"-4295-www.malfong.is**; *v1*;                12,634
7. **"-4324-Government_Offices_I**; *v1*;                18,185
8. **"-4327-Government_Offices_I**; *v1*;                36,290
9. **"-4334-Rkiskaup_2020**; *v1*;                10,236
10. **"-4338-University_Iceland**; *v1*;                10,164
11. **"-502-Icelandic_Financial_**; *v1*;                1,525
12. **"-504-www.iceida.is**; *v1*;                1,055
13. **"-505-www.pfs.is**; *v1*;                2,866
14. **"-506-www.lanamal.is**; *v1*;                1,140
15. **"-5067-SciPar**; *v1*;                110,831
16. **"-508-Tilde_Statistics_Ice**; *v1*;                2,427
17. **"-509-Gallery_Iceland**; *v1*;                577
18. **"-510-Harpa_Reykjavik_Conc**; *v1*;                1,197
19. **"-511-bokmenntaborgin_is**; *v1*;                330
20. **"-516-Icelandic_Medicines**; *v1*;                711
21. **"-517-Icelandic_Directorat**; *v1*;                1,536
22. **"-597-www.nordisketax.net**; *v1*;                1,065
23. **"-718-Statistics_Iceland**; *v1*;                2,361
24. **"-728-www.norden.org**; *v1*;                41,073
25. **"-EMEA**; *v1*;                542,624
26. **"-antibiotic**; *v1*;                816
27. **"-www.norden.org**; *v1*;                41,073
28. **"-www.nordisketax.net**; *v1*;                1,065
29. **EUbookshop**; *v2*;                9,783
30. **GNOME**; *v1*;                28,776
31. **HPLT**; *v1*;                2,148,876
32. **KDE4**; *v2*;                98,989
33. **MaCoCu**; *v2*;                267,366

34. **MultiCCAligned**; *v1*;                1,192,537
35. **MultiHPLT**; *v1*;                     2,148,855
36. **MultiMaCoCu**; *v2*;                     267,366
37. **MultiParaCrawl**; *v7.1*;              2,392,423
38. **NLLB**; *v1*;                          8,723,145
39. **OpenSubtitles**; *v1*;                     7,138
40. **OpenSubtitles**; *v2016*;              1,359,224
41. **OpenSubtitles**; *v2018*;              1,569,189
42. **ParIce**; *v1*;                        2,097,022
43. **ParaCrawl**; *v7.1*;                   2,392,422
44. **ParaCrawl**; *v8*;                     5,724,373
45. **ParaCrawl**; *v9*;                     2,967,579
46. **QED**; *v2.0a*;                           27,611
47. **TED2020**; *v1*;                           2,430
48. **Tatoeba**; *v2*;                           8,139
49. **Tatoeba**; *v20190709*;                    9,436
50. **Tatoeba**; *v2020-05-31*;                  9,438
51. **Tatoeba**; *v2020-11-09*;                  9,440
52. **Tatoeba**; *v2021-03-10*;                  9,443
53. **Tatoeba**; *v2021-07-22*;                  9,443
54. **Tatoeba**; *v2022-03-03*;                  9,522
55. **Tatoeba**; *v2023-04-12*;                  9,600
56. **TildeMODEL**; *v2018*;                   420,712
57. **Ubuntu**; *v14.10*;                        2,155
58. **WikiMatrix**; *v1*;                       85,992
59. **WikiTitles**; *v3*;                       50,176
60. **XLEnt**; *v1*;                           962,661
61. **XLEnt**; *v1.1*;                         962,661
62. **XLEnt**; *v1.2*;                         962,661
63. **bible-uedin**; *v1*;                      62,163
64. **wikimedia**; *v20190628*;                    581
65. **wikimedia**; *v20210402*;                  2,625
66. **wikimedia**; *v20230407*;                  4,471

# B   Filtering steps

### Filter 1. Sentence length
Sentences should contain at minimum four characters and at maximum 150 characters.

### Filter 2. High inter-pair content overlap
Sentence pairs where the content of the source and target sentences are highly similar should be removed from the dataset.

### Filter 3. Character symbol filtering
All characters in the English and Icelandic alphabets (along with punctuation and numbers) designated as a set of allowed characters. Sentences containing less than 60% of these characters removed from the data and all characters outside the allowed

set removed from the remaining sentences.[3]

### Filter 4. LaBSE scoring
We use score each sentence pair using LaBSE (Feng et al., 2022) and remove all sentences with a score lower than 0.8[4].

### Filter 5. Language detection
We use various language detection software to gauge whether both the source and target sentences are in the correct language. The software we used was *fasttext* (Joulin et al., 2016), *franc* (Wormer, 2024), *lingua* (Stahl, 2024) and *langdetect* (Nakatani, 2010).

### Filter 6. Similar dataset pairs
As a safeguard, we remove any duplicate entries of our dataset if, for any reason, there remain duplicate instances after the previous filters. In our final experiment, this was rendered redundant, but was required in previous iterations and may prove useful in future iterations.

### Filter 7. Near-duplicate dataset pairs
Sentences are compared by removing content-specific words that are likely proper names and dates, etc., and comparing the remainder.

### Filter 8. Likely machine-translated target sentences
A GPT-2 (Radford et al., 2019) classifier is used to evaluate whether a given target sentence is machine-translated, based on a 10.000 sentence hand-evaluated reference set. If this is true for the target sentence, that pair is removed from the dataset.

### Filter 9. Existing datasets
As a final safeguard check, we remove any sentence pair that we already have on file in other datasets, as touched on in section 3.2.

### Filter 10. NMTScore cross-likelyhood 0.4
Finally, we use a translation cross-likelyhood NMTScore (Vamvas and Sennrich, 2022) to determine the translation quality of a given sentence pair. This step is computationally heavy and was therefore saved for last. Our experiments suggest that 0.4 is a suitable cutoff for our dataset.

---

[3]This is the last filtering step that inherently modifies the content inside individual sentences.

[4]This is a higher cutoff than the original LaBSE authors suggest to use, but our experiments suggests it better suits our data.

# IKUN for WMT24 General MT Task:
# LLMs Are here for Multilingual Machine Translation

**Baohao Liao**[1,2]    **Christian Herold**[2]    **Shahram Khadivi**[2]    **Christof Monz**[1]

[1]Language Technology Lab, University of Amsterdam

[2]eBay Inc., Aachen, Germany

b.liao@uva.nl

## Abstract

This paper introduces two multilingual systems, *IKUN* and *IKUN-C*, developed for the general machine translation task in WMT24. IKUN and IKUN-C represent an *open system* and a *constrained system*, respectively, built on Llama-3-8b and Mistral-7B-v0.3. Both systems are designed to handle all 11 language directions using a single model. According to automatic evaluation metrics, **IKUN-C achieved 6 first-place and 3 second-place finishes among all constrained systems, while IKUN secured 1 first-place and 2 second-place finishes across both open and constrained systems**. These encouraging results suggest that large language models (LLMs) are nearing the level of proficiency required for effective multilingual machine translation. The systems are based on a two-stage approach: first, continuous pre-training on monolingual data in 10 languages, followed by fine-tuning on high-quality parallel data for 11 language directions. The primary difference between IKUN and IKUN-C lies in their monolingual pre-training strategy. IKUN-C is pre-trained using constrained monolingual data, whereas IKUN leverages monolingual data from the OSCAR dataset. In the second phase, both systems are fine-tuned on parallel data sourced from NTREX, Flores, and WMT16-23 for all 11 language pairs.[1]

## 1 Introduction

Large language models (LLMs) (Touvron et al., 2023; Dubey et al., 2024; Jiang et al., 2023; OpenAI, 2023) serve as a crucial foundation for a wide range of applications. One significant advantage of LLMs is their ability to be applied across various tasks, thereby simplifying deployment processes. However, the application of LLMs to multilingual machine translation (MT) presents several challenges:

- Most LLMs are pre-trained on one or a few dominant languages, making direct fine-tuning on multilingual data insufficient for ensuring optimal performance, particularly for low-resource languages, which are often underrepresented in the training data.

- It remains unclear whether these LLMs, primarily pre-trained on a limited number of languages, effectively facilitate transfer learning across different languages (Tan et al., 2024).

- The large-scale nature of most LLMs presents significant challenges for efficient fine-tuning, particularly for researchers and practitioners with limited computational resources.

In the WMT24 general MT task (Kocmi et al., 2024a), our objective is to assess the capability of LLMs for multilingual MT, as an alternative to training bilingual systems from scratch (Wu et al., 2023). This paper provides a detailed account of how we developed our final multilingual system using LLMs for both the constrained and open tracks.

Firstly, we identified that certain LLMs exhibit inefficiencies in tokenizing sentences from languages that are underrepresented in the pre-training data. To address this, we extended the existing vocabulary to reduce the tokenized sentence length, thereby enhancing training efficiency. Secondly, we enriched the LLMs with knowledge across the 10 target languages through continued pre-training. This step is particularly crucial for underrepresented languages, as it facilitates transfer learning. Finally, we fine-tuned the models using high-quality parallel datasets across all 11 pairs.

Through this streamlined approach, our constrained multilingual system, IKUN-C, secured 6 first-place and 3 second-place rankings in the automatic evaluation. Our open multilingual system, IKUN, achieved 1 first-place and 2 second-place rankings across the open and constrained tracks.

---

[1]Please read our newest verion at https://arxiv.org/abs/2408.11512

These encouraging results demonstrate that LLMs can be effectively adapted for multilingual MT, broadening access to speakers of diverse languages.

## 2 Pre-trained LLM

LLMs are pre-trained on extensive web-scale data, encompassing a vast repository of general knowledge applicable to various tasks. Previous studies (Xu et al., 2024a,b) have demonstrated that LLMs can substantially enhance the performance of multilingual MT. Building on this insight, we adopt a pre-trained LLM as the foundation for our system.

**IKUN** is an open system developed with meticulous consideration of available resources and system capabilities. For this purpose, we selected Llama-3 (Dubey et al., 2024), one of the most advanced open-source LLMs available at the time of this competition. Due to constraints on computational resources, we opted for the 8B version[2] instead of the 70B version. A significant factor in our choice of Llama-3 was its strong support for multilingual applications, as evidenced by the efficiency with which its tokenizer handles all 11 languages involved in this competition (See §3). We also tried the instruct version, but it is worse than the pre-trained version.

**IKUN-C** is a constrained system based on Mistral-7B-v0.3 (Jiang et al., 2023), one of the three LLMs permitted for the constrained track. Prior to selecting Mistral-7B-v0.3[3], we conducted continuous pre-training on all three allowed LLMs — namely, Llama-2-7B, Llama-2-13B (Touvron et al., 2023), and Mistral-7B — using a subset of our monolingual data (approximately 1B tokens). The pre-training loss demonstrated that Mistral-7B outperformed Llama-2-7B and performed comparably to Llama-2-13B, leading us to select it as our architecture of choice.

## 3 Tokenizer Efficiency

A significant challenge in applying LLMs to multilingual MT lies in the efficiency of their tokenizers. These models are typically pre-trained on one or a few dominant languages, and when their tokenizers are applied to low-resource languages, they produce disproportionately long sequences of sub-words. This inefficiency leads to excessive GPU memory consumption during training.

To evaluate the efficiency of the tokenizer, we focus on comparing the length differences between tokenized English sentences and their corresponding non-English counterparts. Specifically, we define the length ratio as:

$$\text{length ratio} = \frac{\text{len(tokenizer}(x))}{\text{len(tokenizer}(y))}$$

where $y$ represents the English sentence, and $x$ denotes the paired non-English sentence. A smaller length ratio (close to 1) is desired, since it means that the tokenizer can encode the non-English sentence as efficient as the English sentence.

To facilitate a comparison of length ratios across different languages, English-centric multilingual data is essential. Fortunately, the FLoRes-200 dataset (Costa-jussà et al., 2022) possesses this characteristic. In the devtest and test sets of FLoRes-200, every English sentence is paired with translations in various other languages. We concatenate all sentences from the devtest set and compute the length ratio for each language, as illustrated in Figure 1. We also include NLLB's tokenizer (Costa-jussà et al., 2022) for a comparison.

We can observe: (1) NLLB consistently exhibits the smallest length ratio across all languages, likely due to its extensive optimization for hundreds of languages, thus serving as a lower bound in this context. (2) Mistral-v0.3 and Llama-3 demonstrate a notably high length ratio for Hindi, suggesting that Hindi is underrepresented in the pre-training data. (3) Compared to NLLB, the tokenizer of Mistral-v0.3 is significantly less efficient for Chinese, Japanese, Hindi, and Icelandic.

We opted to expand the vocabulary by incorporating new sub-words to reduce the length of tokenized sentences, thereby enhancing training efficiency. However, this approach introduces a trade-off between the addition of new sub-words and training performance. The embeddings for the newly introduced sub-words are initially untrained, and a substantial increase in sub-words may necessitate additional iterations of continuous pre-training. Consequently, our strategy for adding sub-words prioritizes those from languages with higher length ratios.

For our open system, IKUN, we didn't modify its tokenizer, since Llama-3 tokenizer is already efficient enough, only except for Hindi. For our constrained system, IKUN-C, we expanded its vocabulary for Chinese, Japanese, Hindi, and Icelandic through the following steps: (1) Generate

---

Figure 1: Tokenizer efficiency for various LLMs and languages. The larger the length ratio is, the less efficient the tokenizer is. We add new sub-words, from Chinese, Japanese, Hindi and Icelandic, to the Mistral-v0.3 vocabulary to construct the IKUN-C vocabulary. IKUN uses the Llama-3 tokenizer without any modification.

| Language pair | Num. | Language pair | Num. |
|---|---|---|---|
| cs-uk | 8768 | uk-cs | 8768 |
| ja-zh | 12858 | zh-ja | 12858 |
| en-zh | 36647 | zh-en | 34650 |
| en-cs | 30120 | cs-en | 28123 |
| en-de | 35564 | de-en | 33567 |
| en-hi | 11020 | hi-en | 9023 |
| en-is | 8010 | is-en | 6013 |
| en-ja | 18113 | ja-en | 16116 |
| en-ru | 32840 | ru-en | 30843 |
| en-es | 13808 | es-en | 11811 |
| en-uk | 11961 | uk-en | 9964 |
| en-fr | 4006 | fr-en | 4006 |
| **Total:** | | | 429457 |

Table 1: Number of parallel sentences.

a new vocabulary of 12K sub-words using monolingual data from the past two years for these four languages from News Crawl[4]; and (2) Merge the new vocabulary with the original. The efficiency of the resulted IKUN-C tokenizer is shown in Figure 1, demonstrating more efficiency for these four languages, especially for Hindi.

## 4 Experiments

We mainly follow the pipeline from Xu et al. (2024a), i.e. continuous pre-training on monolingual data for all 10 languages, and followed by fine-tuning on parallel data for all 11 language pairs.

### 4.1 Continuous Pre-training

Given that our selected LLMs, specifically Llama-3-8B and Mistral-7B-v0.3, are primarily pre-trained on English, it is necessary to incorporate knowledge from other languages through further pre-training, with particular emphasis on low-resource languages. Additionally, the word em-

beddings for the newly introduced sub-words in IKUN-C must also undergo training.

For the open system IKUN, we utilize monolingual data from the Oscar dataset (Suárez et al., 2020), covering all 10 target languages. We adopt the sampling strategy outlined by Xu et al. (2024a), described as:

$$P(l) \propto \left(\frac{D_l}{\sum_{l' \in L} D_{l'}}\right)^{\frac{1}{T}} \quad s.t. \quad \sum_{l' \in L} P(l') = \frac{9}{10}$$

where $D_l$ represents the number of words in language $l$[5], $T$ is the temperature parameter (set to 6), and $L$ denotes the set of all languages except English. The sampling probability for English is fixed at 1/10. The experimental settings for continuous pre-training are detailed in Table 2. We approximately pre-trained IKUN on an additional 8B tokens.

In the constrained system, only the provided data sources are permitted for use[6]. The IKUN-C system utilizes monolingual data from the News Crawl dataset for 9 languages, with the exception of Spanish, as the use of Spanish monolingual data from News Crawl is restricted. For Spanish, we incorporate monolingual data from the Leipzig Corpora (Goldhahn et al., 2012). Additionally, for Hindi, we augment the dataset with monolingual data from News Commentary due to the relatively limited amount of Hindi data available in the News Crawl dataset. This adjustment is crucial because Hindi is underrepresented in the pre-training of Mistral-7B-v0.3. Our experimental settings closely align with those detailed in Table 2.

---

[4]https://data.statmt.org/news-crawl/

[5]https://huggingface.co/datasets/oscar-corpus/OSCAR-2301

[6]https://www2.statmt.org/wmt24/mtdata/

| Hyper-parameter | Continuous pre-training | Finetuning |
|---|---|---|
| sampling probability | cs,de,en,es,hi,is,ja,ru,uk,zh = 0.1,0.13,0.1,0.13,0.08,0.05,0.08,0.13,0.08,0.12 | |
| duration | 60K steps | 1 epoch |
| batch size | 64 | 128 |
| sequence length | 2048 | max source length=512, max target length=512 |
| learning rate (lr) | 2e-5 | 2e-4 |
| warmup ratio | 0 | 0.01 |
| weight decay | 0.01 | 0.01 |
| lr scheduler | cosine | inverse_sqrt |
| training type | full finetuning | LoRA $r = 64$ for all layers |

Table 2: Experimental setting for continuous pre-training and subsequent finetuning.

| System | Metric | cs-uk | en-cs | en-de | en-es | en-hi | en-is | en-ja | en-ru | en-uk | en-zh | ja-zh |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| IKUN-C | MetricX ↓ | 2.4 | 4.3 | 2.0 | 3.5 | 7.1 | 4.9 | 4.3 | 4.7 | 4.7 | 4.2 | 6.2 |
| | CometKiwi ↑ | 0.648 | 0.618 | 0.641 | 0.666 | 0.499 | 0.657 | 0.669 | 0.649 | 0.622 | 0.624 | 0.519 |
| | AutoRank ↓ | 3.0 | 4.7 | 3.8 | 3.4 | 5.5 | 3.7 | 3.9 | 3.9 | 3.9 | 3.5 | 5.5 |
| | Place in constrained ↓ | 2 | 5 | 1 | 1 | 1 | 1 | 3 | 1 | 1 | 2 | 2 |
| IKUN | MetricX ↓ | 1.6 | 3.7 | 1.8 | 3.3 | 9.4 | 4.3 | 3.7 | 4.1 | 3.7 | 4.0 | 5.4 |
| | CometKiwi ↑ | 0.664 | 0.638 | 0.668 | 0.687 | 0.428 | 0.666 | 0.696 | 0.675 | 0.661 | 0.646 | 0.544 |
| | AutoRank ↓ | 2.3 | 3.9 | 3.0 | 2.8 | 7.7 | 3.2 | 3.1 | 3.2 | 2.8 | 3.1 | 4.4 |
| | Place in constrained&open ↓ | 2 | 5 | 4 | 3 | 5 | 1 | 6 | 3 | 2 | 5 | 6 |

Table 3: Preliminary results of our systems on the WMT24 test sets, taken from Kocmi et al. (2024b). The final human evaluation results are not released yet. "-" here means →. I.e. cs-uk is cs→uk.

## 4.2 Subsequent Fine-tuning

Previous studies (Wu et al., 2024; Zhou et al., 2023) have demonstrated that the quality of fine-tuning data is a critical factor in achieving optimal performance. Liao et al. (2021) further indicates that increasing the amount of back-translation data does not necessarily lead to better outcomes. In light of this, we exclusively utilize high-quality parallel data for the fine-tuning phase.

The high-quality parallel data is primarily sourced from FloRes-200 (Costa-jussà et al., 2022), NTREX-128 (Federmann et al., 2022), and previous WMT16-23 general MT/news tasks (Kocmi et al., 2023, 2022; Akhbardeh et al., 2021; Barrault et al., 2020, 2019, 2018; Bojar et al., 2017, 2016).

**FloRes-200:** As the FloRes-200 dataset provides parallel sentences across multiple languages, we leverage all 11 translation directions from the devtest and test sets. Importantly, our fine-tuning approach is not limited to the required directions listed in Table 3; instead, we fine-tune the model on both translation directions, e.g., en↔de, to facilitate broader applicability of the final model.

**NTREX-128:** We also incorporate parallel sentences from NTREX-128 for from-English translation directions, i.e. en→XX. In accordance with Federmann et al. (2022), which recommends using the en→XX translation direction, our fine-tuning

is confined to these directions rather than adopting a bidirectional approach. An exception is made for the en-fr pair, where bidirectional fine-tuning is applied due to the limited availability of parallel data for this pair (absent in previous WMTs).

**Past WMTs:** Additionally, we extract parallel sentences from the development and test sets of the WMT16-23 general MT tasks, provided they contain the necessary translation directions. For these sentences, we employ a bidirectional fine-tuning strategy.

The statistics for all parallel sentences are presented in Table 1. Notably, all systems are fine-tuned at the sentence level. Given that WMT development and test sets are at the document level, models could alternatively be fine-tuned at the document level or reformatted into a conversational structure for fine-tuning. This latter approach might be more effective for context-aware translation, as WMT24 applies context-based human evaluations. We reserve this exploration for future work. The fine-tuning setting is listed in Table 2.

## 4.3 Results

We present the preliminary results reported by Kocmi et al. (2024b) in Table 3, which includes four evaluation metrics. Both MetricX-23-XL (Juraska et al., 2023) and CometKiwi-DA-XL (Rei et al., 2023) have demonstrated strong cor-

relations with human evaluation (Freitag et al., 2023). AutoRank (Kocmi et al., 2024b), a normalized composite metric derived from MetricX and CometKiwi, scales the scores of each metric linearly to span the range from 1 to the total number of systems in a given language pair. The final automatic ranking is obtained by averaging these normalized scores. AutoRank can thus be considered a measure of the overall rank across all systems and tracks. Additionally, we report the rankings of our systems across various tracks.

It is noteworthy that both of our systems are multilingual, designed to handle all language pairs. The IKUN-C system, in particular, demonstrates promising performance in the constrained track, achieving 6 first-place and 3 second-place finishes. In both the open and constrained tracks, IKUN maintains strong performance, securing 1 first-place and 2 second-place positions, even when compared to systems that may leverage additional open-source data or specialize in a limited set of language pairs.

## 5 Conclusion

In this paper, we present a methodology for effectively adapting pre-trained LLMs to the task of multilingual machine translation. Our approach involves three primary steps: (1) expanding the vocabulary to accommodate languages that are underrepresented in the pre-training data, when necessary; (2) continuing pre-training the LLM on monolingual data to enhance its knowledge of underrepresented languages and to train the embeddings of newly introduced sub-words; and (3) fine-tuning the LLM on high-quality parallel data. Our experimental results demonstrate the efficacy of this straightforward pipeline, with IKUN-C securing 6 first-place finishes in the constrained track, and IKUN achieving 1 first-place ranking in both the open and constrained tracks.

## Acknowledgments

## References

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondrej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussà, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation, WMT@EMNLP 2021, Online Event, November 10-11, 2021*, pages 1–88. Association for Computational Linguistics.

Loïc Barrault, Magdalena Biesialska, Ondrej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubesic, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 conference on machine translation (WMT20). In *Proceedings of the Fifth Conference on Machine Translation, WMT@EMNLP 2020, Online, November 19-20, 2020*, pages 1–55. Association for Computational Linguistics.

Loïc Barrault, Ondrej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. Findings of the 2019 conference on machine translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation, WMT 2019, Florence, Italy, August 1-2, 2019 - Volume 2: Shared Task Papers, Day 1*, pages 1–61. Association for Computational Linguistics.

Loïc Barrault, Fethi Bougares, Lucia Specia, Chiraag Lala, Desmond Elliott, and Stella Frank. 2018. Findings of the third shared task on multimodal machine translation. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers, WMT 2018, Belgium, Brussels, October 31 - November 1, 2018*, pages 304–323. Association for Computational Linguistics.

Ondrej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. Findings of the 2017 conference on machine translation (WMT17). In *Proceedings of the Second Conference on Machine Translation, WMT 2017, Copenhagen, Denmark, September 7-8, 2017*, pages 169–214. Association for Computational Linguistics.

Ondrej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, An-

tonio Jimeno-Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana L. Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation, WMT 2016, colocated with ACL 2016, August 11-12, Berlin, Germany*, pages 131–198. The Association for Computer Linguistics.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, and et al. 2022. No language left behind: Scaling human-centered machine translation. *CoRR*, abs/2207.04672.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and et al. 2024. The llama 3 herd of models.

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online. Association for Computational Linguistics.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frédéric Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George F. Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation, WMT 2023, Singapore, December 6-7, 2023*, pages 578–628. Association for Computational Linguistics.

Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. 2012. Building large monolingual dictionaries at the leipzig corpora collection: From 100 to 200 languages. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation, LREC 2012, Istanbul, Turkey, May 23-25, 2012*, pages 759–765. European Language Resources Association (ELRA).

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, and et al. 2023. Mistral 7b. *CoRR*, abs/2310.06825.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin

Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024a. Findings of the WMT24 general machine translation shared task:
the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024b. Preliminary wmt24 ranking of general mt systems and llms.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popovic, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): llms are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation, WMT 2023, Singapore, December 6-7, 2023*, pages 1–42. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popovic. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation, WMT 2022, Abu Dhabi, United Arab Emirates (Hybrid), December 7-8, 2022*, pages 1–45. Association for Computational Linguistics.

Baohao Liao, Shahram Khadivi, and Sanjika Hewavitharana. 2021. Back-translation for large-scale multilingual machine translation. In *Proceedings of the Sixth Conference on Machine Translation, WMT@EMNLP 2021, Online Event, November 10-11, 2021*, pages 418–424. Association for Computational Linguistics.

OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.

Ricardo Rei, Nuno Miguel Guerreiro, José Pombal, Daan van Stigt, Marcos V. Treviso, Luísa Coheur, José G. C. de Souza, and André F. T. Martins. 2023. Scaling up cometkiwi: Unbabel-ist 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation, WMT 2023, Singapore, December 6-7, 2023*, pages 841–848. Association for Computational Linguistics.

Pedro Javier Ortiz Suárez, Laurent Romary, and Benoît Sagot. 2020. A monolingual approach to contextualized word embeddings for mid-resource languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 1703–1714. Association for Computational Linguistics.

Shaomu Tan, Di Wu, and Christof Monz. 2024. Neuron specialization: Leveraging intrinsic task modularity for multilingual machine translation. *CoRR*, abs/2404.11201.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, and et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288.

Di Wu, Shaomu Tan, Yan Meng, David Stap, and Christof Monz. 2024. How far can 100 samples go? unlocking overall zero-shot multilingual translation via tiny multi-parallel data. *CoRR*, abs/2401.12413.

Di Wu, Shaomu Tan, David Stap, Ali Araabi, and Christof Monz. 2023. Uva-mt's participation in the WMT23 general translation shared task. *CoRR*, abs/2310.09946.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024a. A paradigm shift in machine translation: Boosting translation performance of large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024b. Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. *CoRR*, abs/2401.08417.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. LIMA: less is more for alignment. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

# NTTSU at WMT2024 General Translation Task

**Minato Kondo**♠, **Ryo Fukuda**◇, **Xiaotian Wang**♠, **Katsuki Chousa**◇,
**Masato Nishimura**♠, **Kosei Buma**♠, **Takatomo Kano**◇, **Takehito Utsuro**♠
♠University of Tsukuba ◇NTT Communication Science Laboratories

## Abstract

The NTTSU team's submission leverages several large language models developed through a training procedure that includes continual pre-training and supervised fine-tuning. For paragraph-level translation, we generated synthetic paragraph-aligned data and used these data for training. In the task of translating Japanese to Chinese, we focused on speech domain translation. Specifically, we built Whisper models for Japanese automatic speech recognition (ASR). Since the dataset used for Whisper training contains many noisy data pairs, we combined the Whisper outputs using ROVER (Fiscus, 1997) to refine the transcriptions. Furthermore, we employed forward translation from audio as data augmentation, using both ASR models and a base translation model. To select the best translation from multiple hypotheses of the models, we applied Minimum Bayes Risk decoding after Quality Estimation (Fernandes et al., 2022), incorporating scores such as COMET-QE, COMET, and cosine similarity by LaBSE. We explored three different reranking strategies to handle two types of candidates from sentence- and paragraph-level translation and employed a fusion method that integrates all three.

## 1 Introduction

This paper provides a system description of the NTTSU team's submissions to WMT 2024. We took part in the General Translation Task (Kocmi et al., 2024a) for English-to-Japanese (En-Ja) and Japanese-to-Chinese (Ja-Zh). This task has three tracks with different constraints on the use of training data and pre-trained models. For En-Ja, we participated in the constrained track, which provides sets specifically allow training data and pre-trained models for use in traning the translation models. Additionally, for Ja-Zh, we participated in the open track, which allows the use of software and data under any open-source license.

Our team's submission leveraged several large language models developed through a training procedure (Guo et al., 2024; Kondo et al., 2024) that includes continual pre-training and supervised fine-tuning. For paragraph-level translation, we generated synthetic paragraph-aligned data and used these data for training.

In the task of translating Japanese to Chinese, we focused on speech domain translation. Specifically, we built Whisper models (Radford et al., 2022) for Japanese automatic speech recognition (ASR). We used the YODAS dataset (Li et al., 2024) for Whisper training. Since these data contained many noisy data pairs, we combined the Whisper outputs using ROVER (Fiscus, 1997) to refine the transcriptions. Furthermore, to enhance the robustness of the translation model against errors in the transcriptions, we performed data augmentation by forward translation from audio, using both ASR and base translation models.

To select the best translation from multiple hypotheses of the models, we applied Minimum Bayes Risk decoding after quality estimation (Fernandes et al., 2022), incorporating scores such as COMET-QE, COMET, and cosine similarity by LaBSE. We explored three different reranking strategies to handle two types of candidates from sentence- and paragraph-level translation and employed a fusion method that integrates all three.

## 2 System Overview

Our system had three main components: automatic speech recognition (ASR) models, translation models, and a reranking.

This year, speech domain translation was newly incorporated in the above task, and audio data, along with the organizer's transcription, were provided as input data. We were interested in the feasibility of speech translation from Japanese, so we created an ASR model for the Ja-Zh and used its transcription as the additional source text. More-

over, we used ROVER to refine the transcriptions.

For the translation model's architecture, we employed and trained the Transformer model and LLMs. To train the LLms, we carried out monolingual/parallel continual pre-training and supervised fine-tuning. The evaluation for this year was conducted at the paragraph level. To address this, we created sentence- and paragraph-level parallel data and utilized these data to build translation models for each level.

During the inference step, we used the translation models to independently translate at the sentence and the paragraph level, generating multiple candidates. We then selected the best translation candidate using a reranking that combines sentence- and paragraph-level reranking with MBR decoding after quality estimation.

## 3 Automatic Speech Recognition

For Ja-Zh speech translation, we fine-tuned various Whisper-based ASR models for the Japanese ASR task. We used the Japanese subset (ja100) of the YODAS dataset, which consists of approximately 3,000 hours of speech and transcriptions.

### 3.1 Dataset

During the dataset review, we found that the YODAS dataset contained many incorrect transcriptions (e.g., music and non-Japanese speech samples). To mitigate the negative impact of these incorrect samples, we refined the YODAS dataset. We integrated transcriptions of multiple hypotheses transcription generated from multiple ASR models to create a tuning dataset. Specifically, the following procedure was used to generate tuning data.

1. **Generation** We performed beam search decoding with multiple ASR models to generate multiple ASR hypotheses for each speech sample in ja100. This process yielded a set of hypotheses equal to the number of ASR models multiplied by the beam size. We set a beam size of 4.

2. **Language-based Filtering** We applied multistep filtering for the YODAS dataset. First, we used Whisper to transcribe the speech; then, we applied the Compact Language Detector v3 (CLD3) [1] to filter non-Japanese language. Next, we excluded the transcriptions that did not contain Japanese-specific characters (i.e.,

*hiragana* or *katakana*). After language-based filtering, we filterd out uncertain transcription that contained repetition. Specifically, texts with bi-grams appearing more than six times were excluded.

3. **Combination** After filtering, we combined multiple ASR hypotheses using the Recognizer Output Voting Error Reduction (ROVER) (Fiscus, 1997).

4. **CER-based Filtering** To filter uncertain samples of ROVER results, we applied accuracy-based filtering. We measured the character error rate (CER) between the ROVER results and the original subtitles in YODAS. A high CER indicates that either one or both may be significantly inaccurate. For the ASR training, we constructed a development set of 2k samples of CER $\leq 0.3$ data. No CER filtering was applied to the training set because no positive effect was observed in preliminary experiments. Finally, all ROVER results except the development set (1,614,110 segments) were treated as the training set. For the training of MT using the ASR data (described in §4.3), samples with CER $\leq 0.3$ (693,304 segments) were used.

To compare the quality of the original subtitles and the ROVER results, we subjectively evaluated the two corresponding transcriptions of 100 randomly selected samples. As a result, we determined that the ROVER results were of higher quality.

### 3.2 Model

To create the tuning data, we used two pre-trained ASR models: Whisper large-v3[2] and kotoba-whisper-v1.1[3], a Japanese-specific ASR model.

### 3.3 Training

Using the tuning data created through the above procedure with the two ASR models, we separately fine-tuned each of these models. The training of the model was conducted using the AdamW optimizer, with parameters set as $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1e-8$. We employed a linear decay learning rate scheduler and set the warmup steps to 500. The model's parameters were saved every 4000 steps.

---

[1] https://github.com/google/cld3

[2] https://huggingface.co/openai/whisper-large-v3

[3] https://huggingface.co/kotoba-tech/kotoba-whisper-v1.1

The training was carried out with a batch size of 32 samples over a single epoch. We selected the best model based on the loss in the dev set.

## 3.4 Inference

During inference, we performed a beam search with a beam size of 4 and combined these four hypotheses using ROVER. For the post-processing of the ASR stage, we integrated punctuation and sentence segmentation into the transcription. We used the fine-tuned version of xlm-roberta[4] and Bunkai (Hayashibe and Mitsuzawa, 2020)[5] for punctuation insertion and sentence segmentation, respectively. Finally, the two types of hypotheses from the two ASR models were passed to MT.

In the data generation process for MT training (§4.3), ROVER was not performed and the top-1 hypothesis of the beam search was used.

## 4 Primary Translation Model

### 4.1 Dataset

We used two types of text corpora: monolingual and parallel data. Monolingual data are used for monolingual continual pre-training, while parallel data are used for parallel continual pre-training, sentence-level supervised fine-tuning (SFT), and paragraph-level SFT.

**En-Ja** We used the following monolingual corpora: Common Crawl (Kocmi et al., 2022), Leipzig Corpora (Goldhahn et al., 2012), News Crawl, and News Commentary (Kocmi et al., 2023). We also used JParaCrawl v3.0 (Morishita et al., 2022), News Commentary (Kocmi et al., 2023), the Kyoto Free Translation Task Corpus (KFTT) (Neubig, 2011), TED Talks (Barrault et al., 2020), and past WMT test data as the parallel data. Since JParaCrawl v3.0 is automatically created and contains a certain amount of noisy data, we filtered the corpus based on sentence embeddings. We employed LaBSE (Feng et al., 2022) to embed the source and target sentences and then filtered out the sentence pairs in which the similarities are not between 0.4 and 0.9.

**Ja-Zh** We used the following monolingual corpora: Leipzig Corpora (Goldhahn et al., 2012), News Crawl, and News Commentary (Kocmi et al.,

2023). In order to obtain parallel data for continual pre-training, we used JParaCrawl Chinese v2.0 (Nagata et al., 2024). Since this corpus also contains noisy data, we filtered them using the same method as in the En-Ja task. For sentence-level SFT, we used ASPEC-JC (Nakazawa et al., 2016) and Flores-200 (NLLB Team et al., 2022) as training and development sets. In addition to the data for sentence-level SFT, we used News Commentary, WIT3 (Cettolo et al., 2012), Global Voice, and Neulab TedTalks (Tiedemann, 2012) as parallel corpora with context information for paragraph-level SFT.

### 4.2 Model Selection

For the En-Ja task, we used the largest available LLM in the constrained track, `Llama-2-13b`[6] (Touvron et al., 2023). For the Ja-Zh task, we used `TowerBase-13B-v0.1` [7] (Alves et al., 2024), a model based on Llama-2-13b that has been continually pre-trained with monolingual and parallel data.

Additionally, we developed and deployed a Transformer (Vaswani et al., 2017) model trained from scratch. As training data, we used JParaCrawl v3.0 for the En-Ja task and JParaCrawl Chinese v2.0 for the Ja-Zh task. The model configuration and hyperparameters are detailed in Table 1.

### 4.3 LLM Training Procedure

We conducted a three-stage training process based on research conducted on translation models using LLMs (Guo et al., 2024; Kondo et al., 2024). In the first stage, we performed continual pre-training with monolingual data. In the second stage, we conducted continual pre-training with parallel data. Finally, in the third stage, we carried out supervised fine-tuning. The detailed model configuration and hyperparameters are given in Table 1.

**Monolingual Continual Pre-Training** It has been reported that LLMs primarily pre-trained in English, such as Llama-2, have lower translation accuracy for languages other than English (Xu et al., 2024). Therefore, we performed continual pre-training using monolingual data to enhance the

---

[4]https://huggingface.co/1-800-BAD-CODE/
xlm-roberta_punctuation_fullstop_truecase
[5]https://github.com/megagonlabs/bunkai

[6]https://huggingface.co/meta-llama/
Llama-2-13b-hf
[7]https://huggingface.co/Unbabel/
TowerBase-13B-v0.1
[8]https://github.com/facebookresearch/fairseq
[9]https://github.com/huggingface/transformers

| Transformer Enc-Dec model | |
|---|---|
| Subword Size | 32,000 |
| Architecture | Transformer (big) |
| Optimizer | Adam ($\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1e-8$) |
| LR Scheduler | Inverse Square root decay |
| Warmup Steps | 4,000 |
| Max Learning Rate | 1e-3 |
| Dropout | 0.3 |
| Gradient Clipping | 1.0 |
| Label Smoothing | 0.1 |
| Batch Size | 512,000 tokens |
| Number of Updates | 50,000 steps |
| Implementation | fairseq[8] (Ott et al., 2019) |
| Common Settings for All LLMs Training Phases | |
| Warmup Ratio | 1% |
| Gradient Clipping | 1.0 |
| Weight Decay | 1.0 |
| Implementation | transformers[9] (Wolf et al., 2020) |
| Continual Pre-Training Settings | |
| Optimizer | AdamW ($\beta_1 = 0.9, \beta_2 = 0.95, \epsilon = 1e-5$) |
| LR Scheduler | Cosine |
| Max Learning Rate (full / LoRA) | 1.5e-4 / 2.0e-4 |
| Batch Size | 1,024 samples |
| Epoch | 1 |
| Context Length | 2,048 |
| Supervised Fine-tuning Settings | |
| Optimizer | AdamW ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1e-8$) |
| LR Scheduler | Inverse Square root decay |
| Max Learning Rate | 2.0e-4 |
| Batch Size | 1,024 samples |
| Epoch | 3 |
| LoRA Settings | |
| Rank / Alpha | 16 / 32 |
| Dropout | 0.05 |
| Target Modules | QKVO, FFN |

Table 1: Model configuration and hyperparameters.

generation capabilities in languages other than English.

We used randomly sampled data from the monolingual corpora described in §4.1. For the En-Ja task, we created two models, ver1 and ver2, and trained them using randomly sampled data of 1B and 4B tokens, respectively. In contrast, for the Ja-Zh task, we trained only a single model with randomly sampled data of 1B tokens due to the lack of time and GPU resources.

**Parallel Continual Pre-Training** After completing monolingual continual pre-training, we performed continual pre-training using parallel data. Based on the findings of (Kondo et al., 2024), we

used data where the source text is followed by its translation.

For the En-Ja task, the ver1 model was trained using LoRA (Hu et al., 2022), while the ver2 model was trained with full weights. Additionally, ver1 was trained using only the sentence-level parallel data from JParaCrawl v3.0, whereas ver2 utilized JParaCrawl v3.0 along with TED and News Commentary as pseudo-paragraph data.

**Supervised Fine-Tuning** After completing continual pre-training in monolingual and parallel data, we performed supervised fine-tuning using LoRA. The prompts applied to the training data were the same as those used in ALMA (Xu et al., 2024), and the same prompts were used during inference. Note that loss in the prompt outputs was excluded during training (Xu et al., 2024; Kondo et al., 2024).

Additionally, for domain adaptation, we performed SFT using data from each specific domain. For the En-Ja task, the ver1 model was fine-tuned using TED Talks, KFTT, and past WMT test data. In contrast, the ver2 model was fine-tuned with the same three datasets as ver1, plus two additional settings: using only the news domain data and using only the social domain data each from past WMT test data. Note that the past WMT test data used for SFT training consisted of the WMT20 development and test data, with the other test data from WMT21 to WMT23. For WMT21, both En-Ja and Ja-En directions were included, while WMT22 and WMT23 were composed only of the Ja-En direction. Additionally, the development data for all SFT were the WMT22 En-Ja data. As a result, we obtained a total of eight fine-tuned models for En-Ja. For Ja-Zh, we also performed SFT with synthetic data to enhance robustness against errors in the transcription for the speech domain. These data were constructed by forward translation from audio data using ASR and Transformer models.

## 5 Reranking

To enhance translation quality, we applied reranking to the candidate sentences. We conducted a comparative analysis of various methods and strategies, as described in §5.1 and §5.3, on the candidate generated by the methods described in §5.2.

### 5.1 Methods

The reranking approach is used to obtain the final output $\hat{y}$ from the set of candidate sentences $\mathcal{C}$ generated by the methods described in §5.2.

**Quality Estimation (QE)** This approach involves evaluating the candidates using reference-free quality estimation techniques, such as COMET-QE (Rei et al., 2021, 2022, 2023) and sentence embedding-based similarity, and subsequently selecting the candidate with the highest score, as follows:

$$\hat{y} = \underset{c \in \mathcal{C}}{\arg\max} \sum_{i=1}^{m} \lambda_i \text{QE}_i(x, c), \quad (1)$$

where $\text{QE}_i(\cdot, \cdot)$ is a reference-free quality estimation function and $\lambda_i$ represents its weight, subject to $\sum_{i=1}^{m} \lambda_i = 1$.

**Minimum Bayes Risk (MBR) decoding** MBR decoding (Fernandes et al., 2022) employs reference-based metrics to rank translation candidates. It aims to identify the translation that maximizes expected utility while equivalently minimizing the risk (Meister et al., 2020; Eikema and Aziz, 2020) as follows:

$$\hat{y}_i = \underset{c_i \in \mathcal{C}}{\arg\max} \frac{1}{|\mathcal{C}|} \sum_{j=1}^{|\mathcal{C}|} \text{RefMetric}(c_i, c_j), \quad (2)$$

where $\text{RefMetric}(\cdot, \cdot)$ is a reference-based metric. Note that MBR decoding scores the candidate using reference-based metrics by treating all candidates as reference texts without using an actual reference text.

**MBR after QE (QE→MBR)** This approach integrates QE with MBR decoding (Fernandes et al., 2022). The scores produced by the quality estimation procedure determined the top-n sample set from candidate set $\mathcal{C}$ as $\mathcal{C}_{\text{top-n}}$. Subsequently, MBR is applied to $\mathcal{C}_{\text{top-n}}$.

### 5.2 Candidate Generation

We generated five candidates for each model by varying the sampling methods during generation. For the speech domain in Ja-Zh, we had two extra transcriptions from our ASR models in addition to the official one. As a result, we generated five candidates for these two transcriptions and LLM models in the same manner. For models based on Llama-2-13b and TowerBase-13B-v0.1, the five methods were as follows: 1. greedy decoding (no sampling), 2. beam search with a beam size of 4, 3. temperature of 0.9, 4. temperature of 0.5, and 5. temperature of 0.3. For methods 3, 4, and 5, parameters other than temperature were set with

top_p at 0.6 and top_k at 50. We also used the top-5 candidates from beam search for the Transformer with a beam size of 6. As a result, a total of 45 candidate sentences were generated for the En-Ja task using the eight SFT models described in §4.3, along with the Transformer model, making a total of nine models.

Furthermore, for each SFT model, we employed two approaches to generate candidates.

**Sentence-Level Generation** First, we used pySBD[10] (Sadvilkar and Neumann, 2020) to split the original paragraph-level test data into sentences, and then we performed sentence-level inference to generate sentence candidates $\mathcal{C}_{sent}$.

**Paragraph-Level Generation** We used the paragraph data directly as model input for generating paragraph candidates $\mathcal{C}_{para}$.

### 5.3 Reranking System

For the two types of candidates mentioned in §5.2, we used three reranking strategies and one fusion method that integrates all three.

**Synthesized Paragraph Reranking** In each sentence-level inference result, we concatenated the sentences that originally belonged to the same paragraph in order and then performed reranking on the synthesized paragraph.

**Individual Sentence Reranking** We performed sentence-level reranking on the sentence candidates $\mathcal{C}_{sent}$ and then reconstructed the paragraphs from the final reranked results.

**Full Paragraph Reranking** The paragraph candidates $\mathcal{C}_{para}$ were used as the objects of reranking, directly generating paragraph-level results.

**Multi-Attribute Candidate Reranking** We established a larger set of multi-attribute candidates $\mathcal{C}_{mac}$ according to the three reranking strategies mentioned above:

- Synthesized paragraph candidates by concatenating the sentences in order from sentence candidates $\mathcal{C}_{sent}$.

- Paragraph data reconstructed on the results obtained by different reranking methods from Individual Sentence Reranking.

- Paragraph candidates $\mathcal{C}_{para}$ generated by Paragraph-Level Generation.

---

[10] https://github.com/nipunsadvilkar/pySBD

| | CER (YODAS) | COMET (WMT test) |
|---|---|---|
| Whiper large-v3 | 7.7 | 0.4598 |
| + FT | 4.8 | 0.4601 |
| kotoba-whisper-v1.1 | 12.6 | 0.4407 |
| + FT | 5.0 | 0.4518 |
| Official transcription | - | 0.7278 |

Table 2: ASR performances and their translation accuracies. Second column is CER results on the evaluation data of the YODAS dataset. Third column is COMET results on the speech domain of this year's WMT test set.

Then, $\mathcal{C}_{mac}$ was used for paragraph-level reranking.

# 6 Experiment and Analysis

## 6.1 Results of ASR

The second column of Table 2 shows the ASR results (with and without fine-tuning on the YODAS dataset) for the Ja-Zh speech translation. Note that this evaluation was not done in combination with the ROVER system. We confirmed that fine-tuning improved the recognition performance on the YODAS dataset. The third column of Table 2 shows the translation results[11] for the WMT test set. Fine-tuning resulted in a relative improvement of 2.5% for kotoba-whisper-v1.1, but no significant improvement was observed for Whisper-large-v3, even through it demonstrated high ASR performance before fine-tuning. Moreover, our models performed worse than the official transcriptions. We trained the ASR models using relatively short audio samples, whereas the audio samples in the test set were longer than 30 seconds. This gap between the training and test conditions likely contributed to the degradation in speech recognition accuracy. In addition, we prepared ASR models for a wide range of topics, domains, and noise levels for open-domain speech input. For this purpose, we used the YODAS dataset instead of datasets such as TED, CSJ, and Libri, which contain clean speech with human transcriptions. However, this strategy did not turn out to be suitable for the WMT test set. In fact, when we listened to the speech from the test set, the SNR was high and clean. This gap may have also contributed to the degradation. These findings will be leveraged for future improvements.

[11]We used `wmt22-comet-da`. During this evaluation, we used the official transcription as the source text for all hypotheses because it would be the most accurate transcription. `https://huggingface.co/Unbabel/wmt22-comet-da`

| Model | Input | COMET22 |
|---|---|---|
| Ver1 | Sentence | **0.8218** |
| | Paragraph | 0.7666 |
| Ver2 | Sentence | **0.8352** |
| | Paragraph | 0.8349 |

Table 3: COMET Scores of Sentence-Level and Paragraph-Level SFT on WMT23 En-Ja test data

| Scoring Function | COMET22 |
|---|---|
| LaBSE-cos | 0.8364 |
| Comet-QE20 | 0.8797 |
| Comet-QE21 | 0.8837 |
| CometKiwi22 | 0.8821 |
| CometKiwi23-xl | 0.8819 |
| 0.5×Comet-QE20 + 0.5×LaBSE-cos | 0.8835 |
| 0.8×Comet-QE21 + 0.2×LaBSE-cos | **0.8856** |
| 0.9×CometKiwi22 + 0.1×LaBSE-cos | 0.8824 |
| 0.9×CometKiwi23-xl + 0.1×LaBSE-cos | 0.8830 |

| MBR ratio | COMET22 |
|---|---|
| QE (Top 10%) | 0.8911 |
| QE (Top 20%) | 0.8940 |
| QE (Top 30%) | 0.8949 |
| QE (Top 40%) | 0.8950 |
| QE (Top 50%) | **0.8955** |
| QE (Top 60%) | **0.8955** |
| QE (Top 70%) | 0.8954 |
| QE (Top 80%) | 0.8953 |
| QE (Top 90%) | 0.8953 |
| 100% | 0.8953 |

Table 4: COMET Scores of QE and MBR decoding on WMT23 En-Ja test data. The 45 candidates used were generated by the methods in §5.2. MBR decoding was performed after QE with the best scoring function, `0.8×Comet-QE21 + 0.2×LaBSE-cos`.

## 6.2 Sentence-Level versus Paragraph-Level in SFT

In the SFT experiments using past WMT test data, we evaluated whether sentence-level or paragraph-level source texts achieved better accuracy by assessing them with COMET (`wmt22-comet-da`) on the WMT23 En-Ja test data. For paragraph-level training, the data were reconstructed from sentence-level to paragraph-level based on the .xml files provided by WMT. Table 3 shows the results, indicating that sentence-level inputs achieved higher accuracy than those of paragraph-level inputs. Therefore, for subsequent SFT, we used only sentence-level inputs.

## 6.3 Results of Quality Estimation

To identify the scoring function in Eq.(1) that yields the highest translation accuracy, we compared ten

| ID | System | MetricX ↓ | CometKiwi ↑ |
|----|--------|-----------|-------------|
| (a) | Synthesized Para | 2.8830 | 0.7260 |
| (b) | Individual Sent | 2.8100 | 0.7273 |
| (c) | Full Para | 2.7263 | 0.7260 |
| (d) | Multi-Attribute | **2.6321** | **0.7310** |

Table 5: Results of Reranking Systems on WMT24 En-Ja test data. Systems (a)~(c) used 45 candidates, while System (d) used 100 candidates, consisting of 45 from $\mathcal{C}_{sent}$, 45 from $\mathcal{C}_{para}$, and 10 results obtained by Individual Sentence Reranking using the 10 methods listed in Table 4. All of the system results are based on Top 50% MBR decoding after QE with the best scoring function, `0.8×Comet-QE21 + 0.2×LaBSE-cos`.

different scoring functions based on the findings in the paper. We used COMET-QE and LaBSE cosine similarity for scoring functions and evaluated them with COMET on the WMT23 En-Ja test data. Since the WMT23 test data are sentence-level, we used the 45 candidate sentences generated through paragraph-level generation, where each sentence was directly input, as described in §5.2. Additionally, the reranking system utilized Full Paragraph Reranking, as described in §5.3. Table 4 shows the results, indicating that `0.8×wmt21-comet-qe + 0.2×LaBSE-cos` achieved the highest accuracy. Therefore, this scoring function was adopted for subsequent experiments and finally the submitted system.

### 6.4 Resluts of MBR after QE

We investigated the proportion of MBR that achieved the highest accuracy under the same conditions as in §6.3. Table 4 shows the results, indicating that accuracy was maximized at 50%. Therefore, in subsequent experiments and the submitted system, the proportion of MBR was set to 50%.

### 6.5 Results of Reranking Systems

Table 5 shows the results of the reranking system on WMT24 En-Ja. We used MetricX-23-XL (Juraska et al., 2023) and CometKiwi-DA-XL (Rei et al., 2023) as evaluation metrics, consistent with the WMT24 preliminary report (Kocmi et al., 2024b). From these results, it was found that the Multi-Attribute Candidate Reranking achieved the highest accuracy. Therefore, we adopted Multi-Attribute Candidate Reranking for the submitted system.

## 7 Conclusion

In this paper, we described our system for the WMT'24 General Translation Task. We developed

ASR models for the speech domain in Ja-Zh and used Transformer and LLMs for the translation models. We trained LLMs using a three-stage training process: Monolingual Continual Pre-training, Parallel Continual Pre-Training, and Supervised Fine-Tuning. Finally, we applied reranking method and strategies to the translation candidates generated by the translation models. Our analyses confirmed the effectiveness of our reranking method and strategies for paragraph-level translation.

## References

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv:2402.17733*.

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 conference on machine translation (WMT20). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55, Online. Association for Computational Linguistics.

Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. WIT3: Web inventory of transcribed and translated talks. In *Proceedings of the 16th Annual Conference of the European Association for Machine Translation*, pages 261–268, Trento, Italy. European Association for Machine Translation.

Bryan Eikema and Wilker Aziz. 2020. Is MAP decoding all you need? the inadequacy of the mode in neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic

BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

J.G. Fiscus. 1997. A post-processing system to yield reduced word error rates: Recognizer output voting error reduction (rover). In *1997 IEEE Workshop on Automatic Speech Recognition and Understanding Proceedings*, pages 347–354.

Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. 2012. Building large monolingual dictionaries at the Leipzig corpora collection: From 100 to 200 languages. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 759–765, Istanbul, Turkey. European Language Resources Association (ELRA).

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 639–649, Mexico City, Mexico. Association for Computational Linguistics.

Yuta Hayashibe and Kensuke Mitsuzawa. 2020. Sentence boundary detection on line breaks in Japanese. In *Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)*, pages 71–75, Online. Association for Computational Linguistics.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór

Steingrímsson, and Vilém Zouhar. 2024a. Findings of the WMT24 general machine translation shared task:
the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024b. Preliminary wmt24 ranking of general mt systems and llms. *arXiv:2407.19884*.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Minato Kondo, Takehito Utsuro, and Masaaki Nagata. 2024. Enhancing translation accuracy of large language models through continual pre-training on parallel data. In *Proceedings of the 21st International Conference on Spoken Language Translation (IWSLT 2024)*, pages 203–220, Bangkok, Thailand (in-person and online). Association for Computational Linguistics.

Xinjian Li, Shinnosuke Takamichi, Takaaki Saeki, William Chen, Sayaka Shiota, and Shinji Watanabe. 2024. Yodas: Youtube-oriented dataset for audio and speech. *arXiv:2406.00899*.

Clara Meister, Ryan Cotterell, and Tim Vieira. 2020. If beam search is the answer, what was the question? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP),*

pages 2173–2185, Online. Association for Computational Linguistics.

Makoto Morishita, Katsuki Chousa, Jun Suzuki, and Masaaki Nagata. 2022. JParaCrawl v3.0: A large-scale English-Japanese parallel corpus. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6704–6710, Marseille, France. European Language Resources Association.

Masaaki Nagata, Makoto Morishita, Katsuki Chousa, and Norihito Yasuda. 2024. A japanese-chinese parallel corpus using crowdsourcing for web mining. *arXiv:2405.09017*.

Toshiaki Nakazawa, Manabu Yaguchi, Kiyotaka Uchimoto, Masao Utiyama, Eiichiro Sumita, Sadao Kurohashi, and Hitoshi Isahara. 2016. ASPEC: Asian scientific paper excerpt corpus. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 2204–2208, Portorož, Slovenia. European Language Resources Association (ELRA).

Graham Neubig. 2011. The Kyoto free translation task. http://www.phontron.com/kftt.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *arXiv:1902.01382*.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision.

Ricardo Rei, Ana C Farinha, Chrysoula Zerva, Daan van Stigt, Craig Stewart, Pedro Ramos, Taisiya Glushkova, André F. T. Martins, and Alon Lavie. 2021. Are references really needed? unbabel-IST 2021 submission for the metrics shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 1030–1040, Online. Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Nipun Sadvilkar and Mark Neumann. 2020. PySBD: Pragmatic sentence boundary disambiguation. In *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, pages 110–114, Online. Association for Computational Linguistics.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pier-

ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-
icz, Joe Davison, Sam Shleifer, Patrick von Platen,
Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,
Teven Le Scao, Sylvain Gugger, Mariama Drame,
Quentin Lhoest, and Alexander Rush. 2020. Trans-
formers: State-of-the-art natural language processing.
In *Proceedings of the 2020 Conference on Empirical
Methods in Natural Language Processing: System
Demonstrations*, pages 38–45, Online. Association
for Computational Linguistics.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Has-
san Awadalla. 2024. A paradigm shift in machine
translation: Boosting translation performance of
large language models. In *The Twelfth International
Conference on Learning Representations*.

# SCIR-MT's Submission for WMT24 General Machine Translation Task

**Baohang Li, Zekai Ye, Yichong Huang, Xiaocheng Feng, Bing Qin**
Harbin Institute of Technology
{baohangli,zkye,ychuang,xcfeng,qinb}@ir.hit.edu.cn

## Abstract

This paper introduces the submission of SCIR research center of Harbin Institute of Technology participating in the WMT24 machine translation evaluation task of constrained track for English to Czech. Our approach involved a rigorous process of cleaning and deduplicating both monolingual and bilingual data, followed by a three-stage model training recipe. During the testing phase, we used the beam serach decoding method to generate a large number of candidate translations. Furthermore, we employed COMET-MBR decoding to identify optimal translations.

## 1 Introduction

This paper presents the submission from the SCIR-MT in the WMT24 machine translation evaluation task, focusing on the **constrained** track of English to Czech translation. In the field of machine translation, the quality of translation systems has been improved with the development of large language models and the increase in data volume. However, achieving high-quality translation outputs under limited conditions remains a challenging task due to resource and computational constraints (Freitag and Al-Onaizan, 2017).

Our team has adopted a series of innovative methods to address this challenge. Initially, we conducted a rigorous cleaning and deduplication process for both monolingual and bilingual data to ensure the quality of the training dataset. Subsequently, we implemented a three-stage model training strategy, including monolingual continual pre-training, bilingual continual pre-training, and translation-specific supervised intruction tuning. During the testing phase, we utilized the beam search decoding method to generate a multitude of candidate translations and applied the COMET-MBR (Fernandes et al., 2022) decoding strategy to identify the optimal translations.

The structure of this paper is as follows: we first provide a detailed description of the data preprocessing steps and strategies; then, we outline our foundational model selection and training strategy; next, we introduce the decoding algorithms used in the testing phase; and finally, we present the COMET-MBR decoding method and report our experimental results on the wmttest2023 dataset. These methods have led to excellent performance in terms of both BLEU (Post, 2018) and COMET (Rei et al., 2020) scores, demonstrating the effectiveness of our approach.

## 2 Data Preprocessing

### 2.1 Provided Data

**Bilingual Corpus** We used all the provided bitext corpora: Europarl v10, ParaCrawl v9 (Bañón et al., 2020), Common Crawl, News Commentary v18.1, Wiki Titles v3, WikiMatrix, Tilde MODEL corpus, and TED Talks (Cettolo et al., 2012).

**Monolingual Corpus** We also used the following provided monolingual data: News Crawl, Europarl v10, News Commentary, Common Crawl, and Leipzig Corpora (Biemann et al., 2007).

### 2.2 Data Cleaning

Data cleaning played a pivotal role in improving the quality of our training dataset. During this stage, we implemented several key steps to ensure the quality of the bilingual data and monolingual data, respctively.

#### 2.2.1 Bilingual Corpus

Given that a significant portion of the training dataset is synthetically-aligned, we need to use a comprehensive data preprocessing pipeline to ensure good translation quality. In particular, we sequentially performed heuristic-based, statistics-based, and embedding-based methods to filter our data.

280

**Heuristic-based** The following heuristic-based filters are used before applying the others:

- **Language Detection** We use fasttext[1] (Joulin et al., 2016) to filter out sentence pairs mismatching the English-Czech direction.

- **Numerical Matching** If one sentence in a pair has a number (ordinal, date, etc.), we also checked the other sentence if a matching number is present. If a match is not detected, the pair is removed.

**Statistics-based** We employed statistics-based filters on sentence pairs following (Cruz, 2023). We first tokenized then applied the following statistics-based filters:

- **Length Filter** We removed pairs containing sentences with more than 50 characters.

- **Pair Length Ratio** We removed pairs where the ratio of the string lengths between the source and target sentences is greater than 1.2.

- **Symbol Token Ratio** We removed any sentence pairs in which either the source or target sentence appears more than 5 times.

- **Messy Token Ratio** We removed pairs where the number of messy characters in the sentences exceeds 2.

- **Most Frequent Words Gap** We measured the symmetry of bilingual text pairs by calculating the difference in the occurrence counts of the most frequent words in each text, and removed pairs where this difference exceeded 5.

**Embedding-based** Finally, we experimented with the use of sentence embedding models to compute the cross-lingual embedding similarity between the sentence pair. We used LaBSE (Feng et al., 2020) models to embed both the source and target sentences then computed a cosine similarity score between the two. The pair must have a similarity score $0.95 < s \le 1$ to be kept.

After rigorous data cleaning, we filtered the bilingual training data from 56,288,239 pairs to 2,725,848 pairs, retaining only 4.8% of the highest quality data for continual pre-training.

### 2.2.2 Monolingual Corpus

For incremental pre-training of large language models(Wu et al., 2024), we employed the Data-Juicer[2] (Chen et al., 2024b) to filter monolingual data in English and Czech. The filtering part includes the following filters: 1) Number of words, 2) Character repetition ratio, 3) Word repetition ratio, 4) Special character ratio, 5) Stop word ratio, 6) Flagged word ratio, 7) Language identification confidence, 8) Perplexity score, 9) Document length (number of characters), 10) Number of lines, 11) Short line length ratio, 12) Short line ratio.

To address the challenge of assessing the quality of the Czech data, we assumed that the Czech data provided by the competition organizers was of generally acceptable quality, reflecting a reasonable approximation of Czech syntax and expression. To further enhance data quality and improve model performance, we applied the Interquartile Range (IQR) (Whaley III, 2005) statistical method to establish a threshold for data filtering. The IQR method is particularly advantageous because it allows for the objective identification of outliers samples without making specific assumptions about the data distribution.

We calculated the IQR for the Czech dataset to define a reasonable range for data quality. Any samples falling outside this range were deemed potential outliers and excluded from the training data. By evaluating each data pair against these quality filtering criteria, we ensured that only samples within the acceptable range contributed to the training process. This approach enabled us to retain the most representative, high-quality samples, thereby enhancing the overall performance of the translation model. Table 1 presents the number of rows in each dataset with/without filtering.

| Corpus | w/o. Filtering | w. Filtering |
|---|---|---|
| Common Crawl corpus | 333 M | 37 M |
| News Crawl | 12M | 4.6M |
| Leipzig Corpora | 4 M | 1.9 M |
| Europarl v10 | 669 K | 391K |
| News Commentary v18.1 | 283 K | 138 K |

Table 1: Czech Corpus Statistics. Line counts are listed before and after filtering.

---

## 3 Translation Model Training

This Section describes our foundation model selection and model training strategy.

### 3.1 Model Configuration

We adopted LLaMA-2-13B as our foundation model considering its impressive performance on most English benchmarks after pre-training on 1.4T tokens (Touvron et al., 2023). Specifically, our Translation Model was initialized from the LLaMA-2-13B model to reduce the computational cost and continues to train on massive Czech and parallel corpus.

### 3.2 Training Strategy

In pre-trained models such as LLaMA-2, which are primarily trained on English data, integrating monolingual data during continual pre-trainingalongside parallel data has been shown to substantially enhance performance (Guo et al., 2024; Alves et al., 2024). Leveraging this approach, we improved our translation model by first incorporating monolingual data during the continual pre-training phase of models initially trained in English. This was followed by further continual pre-training using parallel data. Finally, we conducted instruction fine-tuning with a limited amount of bilingual data. Our models were developed using the LLaMA-Factory framework (Zheng et al., 2024), which facilitated this comprehensive training process.

#### 3.2.1 Stage-1: Monolingual Continual Pre-training

In the initial phase of our training approach, we conducted secondary pre-training on the large language model (LLM) utilizing the carefully-curated monolingual dataset (shown in 2.2.2). The core objective of this stage is to enrich the LLM's understanding and generation capabilities in non-English languages.

We aimed to strengthen the LLM's multilingual capabilities by exposing it to a diverse monolingual corpus. Although this step was related to machine translation, it was designed to lay a solid foundation for the model's language proficiency, which was critical for the subsequent stages focusing on translation tasks.

**Hyperparameters** We used the AdamW optimizer, with $\beta_1 = 0.9$, $\beta_2 = 0.95$, and $\epsilon =$ $1.0 \times 10^{-8}$. The context length is 2048, and training is conducted for 1 epoch. We performed validation every 100 training steps. We used a cosine learning rate schedule with a warmup ratio of 1% and a peak learning rate of $2 \times 10^{-5}$. We applied a weight decay of 0.1 and gradient clipping of 1.0. We utilized eight NVIDIA RTX A800 GPUs, processing 1 batch on each GPU with a gradient accumulation step of 32, achieving an effective batch size of 256. During training, Flash-Attenion(Dao et al., 2022), bfloat16 precision, gradient checkpointing, and DeepSpeed ZeRO Stage 2(Rasley et al., 2020) were employed. With these configurations, the training process was completed in 5 days, which accelerates the overall training duration.

#### 3.2.2 Stage-2: Bilingual Continual Pre-training

Bilingual Continual Pre-training is a methodology that involves ongoing training on bilingual datasets to improve the model's alignment between languages. This approach facilitates the model's ability to capture detailed syntactic and semantic correspondences across languages. Such fine-grained alignment is helpful for machine translation, as it enhances the accuracy of encoding source language information and improves the quality of the generated translations, thereby producing more precise and fluent translation outcomes.

**Hyperparameters** We performed continual pre-training on the model that achieves the minimum validation. We adopted the AdamW optimizers parameter used in Section 3.2.1. Weight decay and gradient clipping remained the same as in Section 3.2.1. We used a cosine learning rate schedule without warmup and a peak learning rate of $1 \times 10^{-5}$. We conduct validation every 10% of the total training steps for Continual Pre-training with Sentence-aligned Parallel Data, with 1 epoch and a batch size of 256.

#### 3.2.3 Stage-3: Translation-specific Supervised Fine-Tuning

During the instruction fine-tuning stage, we constructed bilingual translation data in a question and answer format, where the instruction language was the source language for translation. We also employed full-scale parameter training. As highlighted in previous research (Xu et al., 2023), instruction fine-tuning of large language models

(LLMs) benefits from limited yet high-quality datasets. To ensure the optimal quality of data during fine-tuning, we followed previous research practices and used translation fine-tuning datasets constructed from the WMT validation data. These datasets, which underwent rigorous quality control measures, were ideal for fine-tuning purposes.

**Hyperparameters** We adjusted the AdamW optimizers parameters as used in Section 3.2.1. Weight decay and gradient clipping remained the same as in Section 3.2.1. The peak learning rate was set to $9.0 \times 10^{-6}$ for full fine-tuning, without warmup, using an inverse square schedule. We conducted validation every 10% of the total training steps for SFT, with 3 epochs and a batch size of 64.

## 4 Decoding algorithms

In the test stage, we first generated multiple candidate translations for the given source sentence. Then, we performed MBR to determine the final translation.

### 4.1 Candidate Generation

During the testing phase, we produced 42 high-quality candidate translations. To enhance the diversity of these results, we employed In-Context Learning (ICL) techniques alongside the beam search algorithm. Specifically, we began by sampling various translation examples to serve as demonstrations, which contributed to greater result diversity. We then applied beam search with a beam width of 5 to generate the final set of 42 top hypotheses. This approach effectively integrates context learning and diversity sampling, thereby optimizing both the coverage and quality of the translations.

### 4.2 COMET-MBR

COMET-MBR (Fernandes et al., 2022) employs Minimum Bayes Risk (MBR) decoding (Kumar and Byrne, 2004; Eikema and Aziz, 2020) with a COMET model (Rei et al., 2020) that has been trained on direct assessments. Typically, a translation $\hat{y}^{MAP} \in \mathcal{V}_Y^{|y^*|}$ is generated using Maximum-A-Posteriori (MAP) decoding, defined as:

$$\hat{y}^{MAP} = \underset{y \in \mathcal{Y}}{argmax} \ \log p(y|x), \qquad (1)$$

where $\mathcal{Y} \subseteq \bigcup_{i=1}^{\infty} \mathcal{V}_Y^i$ represents the search space of target sentences. Unlike MAP decoding, MBR decoding aims to identify the translation that minimizes the Bayes risk:

$$\hat{y}^{MBR} = \underset{h \in \mathcal{Y}}{argmax} \underbrace{E_{y' \sim p(y|x)}[u(y', h)]}_{\approx \frac{1}{m} \sum_{j=1}^m u(y^{(j)}, h)}, \qquad (2)$$

where $\bar{\mathcal{Y}} \subseteq \mathcal{Y}$ denotes a set of translation hypotheses, and $u : \mathcal{Y} \times \mathcal{Y} \rightarrow R$ is the utility function. In our study, we utilize COMET[3] (Rei et al., 2020) as the utility function $u$. It is important to note that the hypotheses set $\bar{\mathcal{Y}}$ and the sample set used for expectation estimation, $\{y^{(1)}, \ldots, y^{(m)}\}$, are shared, except for $h$, i.e., $\{y^{(1)}, \ldots, y^{(m)}\} = \bar{\mathcal{Y}} \setminus \{h\}$. Consequently, given a candidate set, the computational complexity of MBR decoding is on the order of $\mathcal{O}(m^2)$, which leads to slower inference speeds as $m$ increases.

## 5 Experimental Results

We evaluated the translation performance of our system on the WMTTest2023 dataset (Tom et al., 2023) and the Flores-200 benchmark (Costa-jussà et al., 2022). To assess translation quality, we employed both BLEU and COMET scores, utilizing the COMET model Unbabel/wmt22-comet-da[3]. Table 2 provides a comparative analysis with two existing commercial translation systems, Baidu[4] and Google[5]. In this table, "Stage1" "Stage2" and "Stage3" refer to the respective stages of our model training process. The performance labeled as "COMET-MBR" corresponds to the results of applying our MBR decoding approach to the candidate translations.

## 6 Conclusion

In this paper, we describe the materials we submitted for the general translation task at WMT2024. We participated in a constrained track: En→Cs. We trained a machine translation model based on LLaMA, utilizing a comprehensive data pipeline for filtering and curation. This pipeline integrates embedding-based, heuristic-based, and statistics-based filters. Subsequently, we employed a three-stage training method to enhance the translation capabilities of the model. Additionally, we utilized minimum Bayes risk decoding to refine the translation candidates. On two benchmark

---

[3]https://huggingface.co/Unbabel/wmt22-comet-da
[4]https://fanyi.baidu.com
[5]https://translate.google.com

| Methods | Flores | | WMT23 | |
|---|---|---|---|---|
| | BLEU | COMET | BLEU | COMET |
| *Existing Systems* | | | | |
| *Baidu* | 31.43 | 89.26 | 35.34 | 86.04 |
| *Google* | 36.81 | 91.51 | 50.25 | 89.90 |
| *Ours(Based on LLaMA2-13B)* | | | | |
| *Baseline* | 23.74 | 86.44 | 22.08 | 79.71 |
| *+Stage1* | 25.75 | 88.84 | 26.72 | 84.18 |
| *+Stage2* | 31.60 | 89.83 | 33.29 | 85.09 |
| *+Stage3* | 32.95 | 89.51 | 35.60 | 87.76 |
| *+COMET-MBR* | 33.44 | 92.27 | 36.61 | 89.09 |

Table 2: Comparison of translation performance using BLEU and COMET scores. We use LLaMA-2-13B as our base model.

datasets, our system outperformed Baidu and exhibited performance comparable to Google, both of which are unconstrained business systems with significantly more training data.

**Future Directions.** In the future, we aim to investigate how to prevent the catastrophic forgetting problem in the general capabilities of LLMs caused by continual pre-training on non-English data, which will help models benefit from effective translation-specific prompting techniques (Huang et al., 2024a; He et al., 2024; Chen et al., 2024a). Additionally, it is promising to train multiple translation systems based on different pre-training language models and combine their outputs with the ensemble learning strategies (Huang et al., 2024b; Jiang et al., 2023).

# References

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*, 2024.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, et al. Paracrawl: Web-scale acquisition of parallel corpora. Association for Computational Linguistics (ACL), 2020.

Chris Biemann, Gerhard Heyer, Uwe Quasthoff, and Matthias Richter. The leipzig corpora collection-monolingual corpora of standard size. *Proceedings of Corpus Linguistic*, 2007, 2007.

Mauro Cettolo, Christian Girardi, and Marcello Federico. Wit3: Web inventory of transcribed and translated talks. In *Proceedings of the Conference of European Association for Machine Translation (EAMT)*, pages 261–268, 2012.

Andong Chen, Lianzhang Lou, Kehai Chen, Xuefeng Bai, Yang Xiang, Muyun Yang, Tiejun Zhao, and Min Zhang. DUAL-REFLECT: Enhancing large language models for reflective translation through dual learning feedback mechanisms. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 693–704, Bangkok, Thailand, August 2024a. Association for Computational Linguistics. URL https://aclanthology.org/2024.acl-short.64.

Daoyuan Chen, Yilun Huang, Zhijian Ma, Hesen Chen, Xuchen Pan, Ce Ge, Dawei Gao, Yuexiang Xie, Zhaoyang Liu, Jinyang Gao, et al. Data-juicer: A one-stop data processing system for large language models. In *Companion of the 2024 International Conference on Management of Data*, pages 120–134, 2024b.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*, 2022.

Jan Christian Blaise Cruz. Samsung r&d institute philippines at wmt 2023. *arXiv preprint arXiv:2310.16322*, 2023.

Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.

Bryan Eikema and Wilker Aziz. Is map decoding all you need? the inadequacy of the mode in neural machine translation. *arXiv preprint arXiv:2005.10283*, 2020.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. Language-agnostic bert sentence embedding. *arXiv preprint arXiv:2007.01852*, 2020.

Patrick Fernandes, António Farinhas, Ricardo Rei, José GC de Souza, Perez Ogayo, Graham Neubig, and André FT Martins. Quality-aware decoding for neural machine translation. *arXiv preprint arXiv:2205.00978*, 2022.

Markus Freitag and Yaser Al-Onaizan. Beam search strategies for neural machine translation. *arXiv preprint arXiv:1702.01806*, 2017.

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. A novel paradigm boosting translation capabilities of large language models. *arXiv preprint arXiv:2403.11430*, 2024.

Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. Exploring human-like translation strategy with large language models. *Transactions of the Association for Computational Linguistics*, 12:229–246, 2024. doi: 10.1162/tacl_a_00642. URL https://aclanthology.org/2024.tacl-1.13.

Yichong Huang, Xiaocheng Feng, Baohang Li, Chengpeng Fu, Wenshuai Huo, Ting Liu, and Bing Qin. Aligning translation-specific understanding to general understanding in large language models, 2024a. URL https://arxiv.org/abs/2401.05072.

Yichong Huang, Xiaocheng Feng, Baohang Li, Yang Xiang, Hui Wang, Bing Qin, and Ting Liu. Ensemble learning for heterogeneous large language models with deep parallel collaboration, 2024b. URL https://arxiv.org/abs/2404.12715.

Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. LLM-blender: Ensembling large language models with pairwise ranking and generative fusion. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14165–14178, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.792. URL https://aclanthology.org/2023.acl-long.792.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hérve Jégou, and Tomas Mikolov. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*, 2016.

Shankar Kumar and Bill Byrne. Minimum bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 169–176, 2004.

Matt Post. A call for clarity in reporting bleu scores. *arXiv preprint arXiv:1804.08771*, 2018.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. Comet: A neural framework for mt evaluation. *arXiv preprint arXiv:2009.09025*, 2020.

Kocmi Tom, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. Findings of the 2023 conference on machine translation (wmt23):

Llms are here but not quite there yet. In *WMT23-Eighth Conference on Machine Translation*, pages 198–216, 2023.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Dewey Lonzo Whaley III. The interquartile range: Theory and estimation. Master's thesis, East Tennessee State University, 2005.

Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. Continual learning for large language models: A survey. *arXiv preprint arXiv:2402.01364*, 2024.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*, 2023.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyan Luo. Llamafactory: Unified efficient fine-tuning of 100+ language models. *arXiv preprint arXiv:2403.13372*, 2024.

# AIST AIRC Systems for the WMT 2024 Shared Tasks

**Matīss Rikters**[1]
[1]Artificial Intelligence
Research Center (AIRC)
National Institute of Advanced
Industrial Science and Technology
matiss.rikters@aist.go.jp

**Makoto Miwa**[1,2]
[2]Toyota Technological
Institute, Japan
makoto-miwa@toyota-ti.ac.jp

## Abstract

This paper describes the development process of NMT systems that were submitted to the WMT 2024 General Translation and Biomedical shared tasks by the team of AIST AIRC. At WMT 2024 AIST AIRC participated in the General Machine Translation shared task and the Biomedical Translation task. We trained constrained track models for translation between English, German, and Japanese. Before training our models, we first filtered the parallel data, then performed iterative back-translation and additional filtering. We experimented with training baseline Transformer models, Mega models, and fine-tuning open-source T5 and Gemma model checkpoints using the filtered parallel data. Our primary submissions contain translations from ensembles of two Mega model checkpoints and our contrastive submissions are generated by our fine-tuned T5 model checkpoints.

## 1 Introduction

We describe the machine translation (MT) systems submitted to the WMT 2024 General Translation and Biomedical Translation tasks developed by the team of AIST AIRC. We experimented with data quality control by filtering out noisy examples from parallel and monolingual data sets before training, and corpora selection. We also compared several modeling approaches by contrasting our previous year's best constrained submission (Rikters and Miwa, 2023) – the Mega model (Ma et al., 2023) to open track approaches of fine-tuning T5 (Raffel et al., 2020) and Gemma (Mesnard et al., 2024) model open-source checkpoints. When fine-tuning T5 and Gemma models, we experimented with adding named entity (NE) annotations (Rikters and Miwa, 2024) to improve rare word translation, since struggling to correctly translate less common NEs was one of the most common errors identified in human evaluations of our WMT 2023 submissions.

## 2 Data

In the General Translation task we only participated in the constrained track, so our data selection was limited to only the parallel corpora provided by the shared task organizers, which for German and Japanese was unchanged from the previous year. For the Biomedical Translation task we used a combination of General Translation task data and Biomedical Translation task data.

All parallel training data and monolingual data for back-translation were filtered before starting any training, which has been proven very effective in previous WMT shared tasks (Pinnis et al., 2018). The filtering process we used is detailed by Rikters (2018). We did not perform any parallel data distillation for our submissions this year.

For the system development process in the General Translation task, we selected News Test sets from the WMT 2022 shared task as development data and test sets from WMT 2023 as evaluation data. Statistics of the data we used are shown in Table 1. For the Biomedical Translation task we used the same combination of 2022 and 2023 development / evaluation data sets.

### 2.1 Data Selection

To not overwhelm the full combined training data set with lower-quality web-crawled data, we 1) limited the English-German Paracrawl to 50 million parallel sentences; and 2) up-scaled all data from other sources to match the amount of the Paracrawl data after filtering by doubling for English-German and tripling for English-Japanese.

### 2.2 Filtering

Even though all training data need not always be perfect and methods like back-translation intentionally generate somewhat noisy additional training data, some types of noise are more harmful than others. Since most training corpora are produced

| Corpus / Filtering | | EN-DE | EN-JA |
|---|---|---|---|
| All other | Before | 16,752,302 | 8,076,155 |
| | After | 13,737,028 | 7,076,869 |
| Paracrawl | Before | 50,000,000 | 21,891,738 |
| | After | 44,533,635 | 21,088,689 |
| | Combined | 72,007,691 | 42,319,296 |
| Medline | | 45,796 | - |
| UFAL Medical | | 3,036,581 | - |

| Corpus / Filtering | Monolingual | |
|---|---|---|
| | Before | After |
| DE | 43,613,631 | 37,110,981 |
| JA | 22,193,545 | 21,558,123 |
| EN | 47,333,840 | 36,756,542 |

Table 1: Training data statistics for all other parallel data without Paracrawl, a subset of Paracrawl, combined development and evaluation data from the past WMT shared tasks, and monolingual data. Sentence counts are listed before and after filtering.

partially or fully automatically, errors such as misalignments between source and target sentences or direct copies of source to target can occur, as well as some amounts of third language data in seemingly bilingual data sets.

To avoid such problems, we used data cleaning and pre-processing methods described by Rikters (2018). The filtering part includes the following filters: 1) unique parallel sentence filter; 2) equal source-target filter; 3) multiple sources - one target and multiple targets - one source filters; 4) non-alphabetical filters; 5) repeating token filter; and 6) correct language filter. We also perform pre-processing consisting of the standard Moses (Koehn et al., 2007) scripts for punctuation normalization, cleaning, and Sentencepiece (Kudo and Richardson, 2018) for splitting into subword units for training MEGA models, and the default tokenizers for T5 and Gemma. The filters were applied to the given parallel sentences, monolingual news sentences before performing back-translation, and both sets of synthetic parallel sentences resulted from back-translating the monolingual news.

## 2.3 Back-translation

Increasing the amount of in-domain training data with synthetic back-translated corpora (Sennrich et al., 2016) is a common practice in cases with considerable amounts of in-domain monolingual data. However, since the task recently shifted from

'news' to 'general' text translation, the definition of what would be considered in-domain data became less clear. Furthermore, for the constrained track the selection of provided monolingual data from the organizers is still limited to news and web-crawled data. No other monolingual data that would be considered more similar to what the 'general' test sets may include, such as user generated (social media), conversational, and e-commerce data are provided in the task. For our experiments we continued to assume that a significant portion of the test data would still be from the news domain. Therefore, we chose to only use the provided monolingual News crawl, News discussions, and News Commentary corpora for back-translation.

## 2.4 Post-processing

In post-processing of the model output we aimed to mitigate some of the most commonly noticeable mistakes that the models were generating. We mainly noticed two often occurring problems in output from all models: 1) difficulties in translating emoji symbols; and 2) occasional repetitions of words or phrases.

While all English and German alphabet letters and even Japanese characters are covered in the large training data corpora, the Unicode emoji were mostly formed and clearly defined only in the past decade, and new emoji are still added every year or two with the next release planned for late 2024[1]. Emoji are also not often present in MT training data, therefore full emoji coverage is absent from model vocabularies, which leads to occasional <unk> tokens being generated as output if emoji were present in the input. In order to keep using the models without re-training, we replaced any <unk> tokens in the output using a dictionary of any emojis appearing in the input.

Furthermore, the occasional hiccuping or hallucinating of models on less common input sequences seems to still be present, sometimes generating repetitions of tokens or phrases. We replaced any consecutive repeating n-grams with a single n-gram. The same was applied to repeating n-grams that have a preposition between them, i.e., *the victim of the victim*.

Both post-processing approaches gave BLEU score improvements of around 0.1 - 0.2.

---

[1] https://emojipedia.org/unicode-16.0

## 3 Model Configurations

While preparing our submissions we experimented with three main model types between the constrained and open system tracks. For our primary submission we chose the constrained Mega models similar to our last year's primary submission (Rikters and Miwa, 2023), and for contrastive submissions we used T5 models (Raffel et al., 2020) fine-tuned on NE-annotated General Translation task data, and Gemma models (Mesnard et al., 2024) tuned on General Translation task data.

### 3.1 Mega

Ma et al. (2023) proposed a moving average equipped gated attention mechanism (MEGA) - a single-head gated attention mechanism equipped with exponential moving average to incorporate inductive bias of position-aware local dependencies into the position-agnostic attention mechanism. Compared to the Transformer model, MEGA has a single-head gated attention mechanism instead of multi-head attention, which enables gains in efficiency while not sacrificing on performance.

For training our Mega models we used the implementation[2] provided by the authors, which is based on FairSeq (Ott et al., 2019).

### 3.2 T5

We experiment with multi-task training and fine-tuning the T5 model (Raffel et al., 2020) for translation between English → German, as well as its multilingual counterpart mT5 (Xue et al., 2021) for English → Japanese translation. We compare the results with non-modified versions of T5, Flan-T5, and the multilingual mT5.

We combine and shuffle all training data for the tasks, and experiment fine-tuning the large versions ( 1B parameters) of the T5 models using a random subset of 10M parallel sentences. We base this choice on observations from preliminary experiments where the small versions of T5 models often converged before reaching 1M examples and base models converged before seeing 10M, since the pre-trained checkpoints are already quite capable as is.

We used the Adafactor optimizer (Shazeer and Stern, 2018) with FP16 training, effective batch sizes of 256 or 512 depending on the model size, evaluation every 1000 steps, and early stopping set to 10 checkpoints of evaluation loss not improving.

[2] https://github.com/facebookresearch/mega

We set learning rate to 0.0001, weight decay to 0.01, and train each model on a single machine with eight NVIDIA A100 GPUs.

### 3.3 Gemma

We experimented with adapting 7B and 9B parameter sizes of the 1.1 and 2 version Gemma models (Mesnard et al., 2024) using the in-domain data provided for the General Translation shared task. We used the same random subset of 10M training examples as we did for training T5 models.

## 4 Results

### 4.1 General Translation Task

We include the official preliminary automatic ranking results provided by the organizers in Tables 2 and 3. Our primary submissions rank 2nd and 4th among the constrained track (with a white background) for EN-DE and EN-JA respectively. Sadly, they were both not selected for human evaluation by the task organizers due to a large number of submissions and budget constraints this year. References had also not been released as of writing the final submission, therefore, additional metrics or manual assessment of the translations could not be performed.

### 4.2 Biomedical Translation Task

For the Biomedical Translation task we compared our best models trained for the General Translation task with ones fine-tuned on the biomedical training data, as well as dedicated models trained on the biomedical data from the start. Table 4 shows our preliminary results from developing Mega models for the English↔German tracks of the Biomedical Translation task. We only used different configurations of the MEGA models and compared them with the baseline model submitted to the general translation task. Our best configuration was an ensemble of three separate model checkpoints trained on a mixture of biomedical training data and general data, and fine-tuned on biomedical data.

Table 5 lists the preliminary official results of the Biomedical Translation task provided by the task organizers. According to the BLEU scores, our models seem to be ranked 2nd in both translation directions, overtaken only by the submissions from Unbabel, which are 70B parameter large language models. Similarly to the General Translation task, references for these had also not been released as of writing the final submission, therefore, additional

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation |
|---|---|---|---|---|
| IOL-Research | 2.3 | 1.6 | 0.692 | ✓ |
| Llama3-70B § | 2.5 | 1.7 | 0.686 | ✓ |
| Aya23 | 2.7 | 1.8 | 0.680 | ✓ |
| IKUN | 3.0 | 1.8 | 0.668 | ✓ |
| IKUN-C | 3.8 | 2.0 | 0.641 | ✓ |
| CUNI-NL | 4.2 | 2.1 | 0.624 | |
| **AIST-AIRC** | 7.2 | 3.3 | 0.551 | |
| Occiglot | 8.2 | 3.8 | 0.539 | |
| MSLC | 11.9 | 4.4 | 0.390 | |
| TSU-HITs | 13.3 | 5.6 | 0.395 | |
| CycleL2 | 27.0 | 11.5 | 0.091 | |
| CycleL | 27.0 | 11.5 | 0.091 | |

Table 2: Preliminary WMT24 General MT automatic ranking for English→German (excluding closed systems).

| System Name | AutoRank ↓ | MetricX ↓ | CometKiwi ↑ | Human evaluation |
|---|---|---|---|---|
| Team-J | 1.9 | 2.9 | 0.740 | ✓ |
| NTTSU | 1.9 | 2.6 | 0.731 | ✓ |
| IOL-Research | 2.3 | 3.1 | 0.724 | ✓ |
| Aya23 | 2.3 | 3.1 | 0.719 | ✓ |
| Llama3-70B § | 2.6 | 3.5 | 0.714 | ✓ |
| IKUN | 3.1 | 3.7 | 0.696 | |
| IKUN-C | 3.9 | 4.3 | 0.669 | ✓ |
| **AIST-AIRC** | 6.6 | 6.5 | 0.583 | |
| CycleL | 24.0 | 22.4 | 0.101 | |

Table 3: Preliminary WMT24 General MT automatic ranking for English→Japanese (excluding closed systems).

metrics or manual assessment of the translations could not be performed.

## 5 Conclusion

In this paper we described the development process of the AIST AIRC's NMT systems that were submitted for the WMT 2024 shared tasks on general domain text translation and biomedical translation. We compared training MEGA models to fine-tuning T5 and Gemma model architectures in search of the best decoding approach for improving upon output quality. Our results showed that the MEGA model architecture remains highly competitive even in the modern world of large language models, and fine-tuning LLMs with NE-annotated data does not necessarily lead to higher automatic evaluation scores. Especially in the Biomedical Translation task our 100M parameter models demonstrated high competitiveness with the leading 70B parameter models, falling only

0.42 BLEU points behind for EN→DE.

In total, output from four primary systems was submitted to the two shared tasks by AIST AIRC for the English↔German and English→Japanese translation directions.

In future work, we plan to experiment with incorporating document-level training data and modeling longer sequences with appropriate available training data. In terms of data, we intend to increase vocabulary coverage by adding all known unicode emoji symbols to the vocabulary even if they are not present in the training data, as well as additionally sample Paracrawl data where emoji are present.

## Acknowledgements

| Configuration | EN→DE | DE→EN |
|---|---|---|
| Baseline General model | 27.23 | 35.00 |
| General BT model | 26.47 | 33.90 |
| Bio trained/adapted | 31.33 | 40.21 |
| Bio-Baseline ensemble | 30.95 | 39.14 |
| Bio-best-last | 31.33 | 40.14 |
| Bio-ens-15 | 31.23 | 40.12 |
| Bio-ens-14 | 31.21 | 39.80 |
| Bio-ens-14-15 | 31.44 | 40.17 |
| Bio-ens-14-15-2 | **31.47** | **40.45** |

Table 4: Biomedical task development BLEU score results evaluated on the 2023 Biomedical Translation task test set. The top 3 rows are single model results from the baseline model of the General Translation task, the model after back-translation (BT), and the models specifically trained and adapted on the biomedical (Bio) task data. All remaining rows are combinations of ensembles consisting of best, last, and other checkpoints from the baseline and biomedical specific models.

| System Name | EN→DE | DE→EN |
|---|---|---|
| ADAPT | 30.16 | 36.93 |
| **AIST-AIRC** | 33.80 | 45.92 |
| DCUGenNLP | 16.46 | 32.60 |
| HW-TSC | 28.77 | 45.79 |
| Unbabel | 34.22 | 49.05 |

Table 5: Preliminary WMT24 Biomedical Translation Task BLEU score results.

## Ethics Statement

Our work fully complies with the ACL Code of Ethics[3]. We use only publicly available datasets and relatively low compute amounts while conducting our experiments to enable reproducibility. We do not perform any studies on other humans or animals in this research.

## References

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. 2023. Mega: Moving average equipped gated attention. In The Eleventh International Conference on Learning Representations.

Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open models based on gemini research and technology.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations), pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

---

[3]https://www.aclweb.org/portal/content/acl-code-ethics

290

Mārcis Pinnis, Matīss Rikters, and Rihards Krišlauks. 2018. Tilde's machine translation systems for WMT 2018. In Proceedings of the Third Conference on Machine Translation: Shared Task Papers, pages 473–481, Belgium, Brussels. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(140):1–67.

Matīss Rikters. 2018. Impact of Corpora Quality on Neural Machine Translation. In Proceedings of the 8th Conference Human Language Technologies - The Baltic Perspective (Baltic HLT 2018), Tartu, Estonia.

Matiss Rikters and Makoto Miwa. 2023. AIST AIRC submissions to the WMT23 shared task. In Proceedings of the Eighth Conference on Machine Translation, pages 155–161, Singapore. Association for Computational Linguistics.

Matiss Rikters and Makoto Miwa. 2024. Entity-aware multi-task training helps rare word machine translation. In Proceedings of the 17th International Natural Language Generation Conference, pages 47–54, Tokyo, Japan. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In International Conference on Machine Learning, pages 4596–4604. PMLR.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 483–498, Online. Association for Computational Linguistics.

# Occiglot at WMT24: European Open-source Large Language Models evaluated on Translation

**Eleftherios Avramidis**[1], **Annika Grützner-Zahn**[1], **Manuel Brack**[1,2],
**Patrick Schramowski**[1,2,3], **Pedro Ortiz Suarez**[4], **Malte Ostendorff**[5], **Fabio Barth**[1],
**Shushen Manakhimova**[1], **Vivien Macketanz**[1], **Georg Rehm**[1], **Kristian Kersting**[1,2,3,6]

[1]German Research Center for Artificial Intelligence (DFKI), Germany
[2]Computer Science Department, TU Darmstadt [3]Hessian.AI [4]Common Crawl Foundation
[5]Occiglot [6]Centre for Cognitive Science, TU Darmstadt

## Abstract

This document describes the submission of the very first version of the Occiglot open-source large language model to the General MT Shared Task of the 9th Conference of Machine Translation (WMT24). Occiglot is an open-source, community-based LLM based on Mistral-7B, which went through language-specific continual pre-training and subsequent instruction tuning, including instructions relevant to machine translation. We examine the automatic metric scores for translating the WMT24 test set and provide a detailed linguistically-motivated analysis. Despite Occiglot performing worse than many of the other system submissions, we observe that it performs better than Mistral7B, which has been based upon, which indicates the positive effect of the language specific continual-pretraining and instruction tuning. We see the submission of this very early version of the model as a motivation to unite community forces and pursue future LLM research on the translation task.

## 1 Introduction

Occiglot, initiated in March 2024, is a community-based open-source initiative for "Polyglot Language Models for the Occident". We believe that our dedicated language modeling solutions will not only maintain Europe's academic and economic competitiveness and AI sovereignty, but also have a profound Impact on the preservation of linguistic diversity, multilingualism, and cultural richness. Occiglot is an academic, non-profit research collective committed to open science and open-source LLM development.

Although Occiglot is in the early stages of development, it entails a significant amount of work for large-scale data collection, model pre-training and tuning, and multi-faceted evaluation. Since LLMs can be used in various use cases, targeted evaluation, starting in the first stages, is important for revealing strengths and weaknesses. The shared task of the 9th Conference of Machine Translation (WMT24; Kocmi et al., 2024a) provides the opportunity for testing the performance of the LLM in a translation task.

First, this paper reviews some indicative items of related work (section 2). Then, in section 3 we present the details on the development of the Occiglot model (section 3.1), the training data related to translation (section 3.2) and the engineering towards machine translation and outline the issues and directions for further improvements. Section 4 presents the evaluation, whereas a conclusion is given in section 5.

## 2 Related work

Prompting LLMs for translation output has been successfully employed since the early years of LLMs (Brown et al., 2020), with the few-shot enhanced context approach indicating good results (Vilar et al., 2023). Later approaches suggested that an adaptive method of few-shot prompting may be even more beneficial (Agrawal et al., 2023; Zhang et al., 2023; Soudi et al., 2024). Enis and Hopkins (2024) deal with evaluating Claude 3 Opus, as compared to other LLMs, with regard to machine translation of low resource languages.

The motivation of Occiglot, to focus LLM development on languages other than English, is confirmed by Diandaru et al. (2024), who suggest that models centered around languages other than English could provide a more efficient foundation for multilingual applications. Zan et al. (2024) follow a similar approach to ours, including instruction tuning tailored to particular target languages. Stap et al. (2024) suggest that including monolingual data as part of the fine-tuning data, we can maintain the abilities while simultaneously enhancing overall translation quality.

292

## 3 The language model

### 3.1 Training

The submission at WMT24 is based on the current, first version (v0.1) of the Occiglot bilingual models for English-Spanish and English-German, released in March and April 2024 respectively. That version provides a broader LLM collection for the five largest European languages: English, German, French, Spanish, and Italian. Out of these languages, only German and Spanish are official language directions of the WMT24 shared task and, therefore, the respective bilingual models are chosen for this submission.

The models are based on the Mistral-7B, which was pre-trained for English. In addition, bilingual continual pre-training and subsequent instruction tuning for each language were performed. Both models include the dataset Open-Hermes-2B[1], which contains content in English language and code. The German model `occiglot-7b-de-en-instruct` was trained on 180M tokens of additional multilingual and code instructions, including the German subsets of DiscoLM (which includes the publicly available `germanrag` dataset), Open Assistant Conversations Dataset v2 (OASST-2; Köpf et al., 2023) and Aya-Dataset (Singh et al., 2024). The Spanish model `occiglot-7b-es-en-instruct` was trained on 160M tokens of additional multilingual and code instructions, including the datasets Mentor-ES, the Stanford Question Answering Dataset v2[2] (SQuAD; Carrino et al., 2020) and the Spanish subsets of OASST-2 and Aya-Dataset.

The full instruction fine-tuning took place on an H100 with 8 GPUs for 0.6–4 training epochs (depending on dataset sampling). We used the `axolotl` framework, maintaining a precision of bf16, a global batch size: 128 (with 8192 context length and Cosine Annealing with Warm-up). The tokenizer is unchanged from Mistral-7B-v0.1.

All pre-trained and instruction-tuned checkpoints are available on Hugging Face[3] under the Apache 2.0 license. Note that the model was not safety-aligned and might generate problematic outputs.

### 3.2 Translation data during training

Both the bilingual German and Spanish models were subjected to paired English translation data during continual pre-training. Specifically, the training data contains paired sentences from Tatoeba (Tiedemann, 2020) and Opus 100 (Zhang et al., 2020). The samples are presented as one coherent text using a diverse set of templates, like

```
Given the following passage:
<German sentence>
a good English translation is:
<English sentence>
```

About 470k and 380k similar translation examples were included during the continual pre-training of the bilingual German and Spanish model, respectively.

Additionally, the instruction tuning stage of both models also includes multilingual data. For the bilingual Spanish model, as mentioned above, parts of the instruction training set were taken from a translated version of the SQuAD, which contains Spanish questions about English literature, for example. More importantly for our task, the incorporated open-assistant OASST-2 dataset also includes about 100 samples of direct instructions for translations between English and Spanish. Similarly, the employed German instruction tuning dataset contains over 2000 dedicated translation examples.

### 3.3 Prompting translations

During the development of the model, we devised a system prompt instructing the model to perform as a dedicated translator and we found that this prompt is immensely helpful when employing the downstream model for translation tasks. Nevertheless, for the WMT submission we decided to use a prompting method which is similar to the way other LLMs are prompted, so that the results are comparable. Prompting was based on the 5-shot templates used by the organizers General Shared task of Machine Translation to prompt GPT-4[4]. The exact prompt used can be seen in Figure 1.

The suggested practice for MT prompting is multi-shot, where one provides first 4 source/translation samples and then only a source awaiting the translation. Occiglot was giving as an answer not only the translation, but was proceeding with generating more text, on the similar

---

[1] https://huggingface.co/teknium
[2] https://huggingface.co/datasets/ccasimiro/squad_es
[3] https://huggingface.co/collections/occiglot/occiglot-eu5-7b-v01-65dbed502a6348b052695e01

[4] https://github.com/wmt-conference/wmt23-news-systems/tree/master/tools/LLM-prompt

```
SYSTEM_PROMPT = "You are a very good
translator. Please translate the given
texts from English to 1.  target_lang
as precisely and accurately as possible
without changing the structure and answer
only with one translation."

PROMPT = "Please translate this into 1.
{target_lang}:

{source_seg}
1. {translation}"
```

Figure 1: Prompt used

| System Name | AutoRank↓ | MetricX↓ | Comet Kiwi↑ |
|---|---|---|---|
| Unbabel | 1.0 | 1.1 | 0.723 |
| Dubformer | 1.8 | 1.2 | 0.694 |
| ... | | | |
| GPT-4 | 1.8 | 1.4 | 0.700 |
| ... | | | |
| Mistral-Large | 2.0 | 1.5 | 0.694 |
| ... | | | |
| IKUN-C | 3.8 | 2.0 | 0.641 |
| ... | | | |
| CUNI-NL | 4.2 | 2.1 | 0.624 |
| AIST-AIRC | 7.2 | 3.3 | 0.551 |
| NVIDIA-NeMo † | 7.4 | 3.5 | 0.558 |
| Occiglot | 8.2 | 3.8 | 0.539 |
| MSLC | 11.9 | 4.4 | 0.390 |
| TSU-HITs | 13.3 | 5.6 | 0.395 |

Table 1: Indicative comparisons from the preliminary WMT24 General MT automatic ranking for English-German.

| System Name | Comet Kiwi ↑ |
|---|---|
| Occiglot | 0.539 |
| Mistral 7B v0.1 | 0.429 |

Table 2: Comparison between Occiglot and its pre-trained model Mistral7B on English-German

pattern, which was difficult to post-process. We had to write a post-processing script that isolates the translation from the additional superfluous text. Nevertheless, we suspect that this post-processing script may have not operated properly in all cases, as we have some hundreds of empty outputs.

The second issue we faced was the inference speed. We loaded the model locally on a python script in the GPU cluster and used the hugging-face `pipeline` command to prompt. The German model was too slow (2-7sec per segment), which made it very tight to meet the deadline. We therefore enabled multiple workers with batches (batch_size=64, num_workers=4) which gave indeed a big acceleration. The behavior of the model was a bit different in the batch mode, so we had to include a system prompt (which was not used for the Spanish model). The parameters of the request command with batches were also different (e.g. the limit max_new_tokens), so it is not sure if parallelizing gave the same results as the single worker mode would have given. The Spanish model was fast enough, and the Spanish test set significantly smaller, so we didn't have to parallelize.

Finally, the German model was going through memory spikes and was killed several times by the administrator rules of our GPU cluster. This may have to do with the test set, as the German test set contains a higher number of examples with more complex sequences. In the future, we have to modify our scripts to stream directly to a file and have the possibility to resume from a particular line in case of a crash.

## 4 Evaluation

### 4.1 Comparison with other WMT systems

The preliminary results (Kocmi et al., 2024b) of the General MT task, based on automatic measures Table 1, indicate a low performance of Occiglot as compared to other systems. We attribute these results to the fact that the development of our LLM is in the early stage and the model has undergone a relatively minimal optimization for translation. Additionally, we have strong indications that the post-processing script did not account for all possible cases. The fact that the model delivered some hundreds of empty outputs is also a matter that may have contributed to the low scores (although it needs to be noted that the parent model Mistral-Large, prompted by the WMT24 organizers, has delivered a higher number of empty outputs). Finally, we should note that the comparison is mostly done with LLMs with a higher number of parameters, as compared to our system. Therefore, this comparison should only be seen with a grain of salt.

### 4.2 Comparison with pre-trained model

Occiglot performs better in translating from English-German than the pre-trained model Mistral 7B v0.1, it has been based on. This indicates a

| category | items | acc |
|---|---|---|
| Ambiguity | 22 | 86.4 |
| Coordination & ellipsis | 124 | 60.5 |
| False friends | 40 | 92.5 |
| Function word | 40 | 75.0 |
| LDD & interrogatives | 207 | 76.3 |
| Lexical Morphology | 39 | 61.5 |
| MWE | 123 | 76.4 |
| Named entity & terminology | 112 | 77.7 |
| Negation | 18 | 66.7 |
| Non-verbal agreement | 109 | 87.2 |
| Punctuation | 37 | 51.4 |
| Subordination | 191 | 85.3 |
| Verb semantics | 23 | 60.9 |
| Verb tense/aspect/mood | 3249 | 71.9 |
| Verb valency | 114 | 65.8 |
| micro-average | 4448 | 72.8 |
| macro-average | 4448 | 73.0 |

Table 3: Performance of the Occiglot English-German model with regard to linguistically-motivated categories

success of the bilingual continual pre-training and subse- quent instruction tuning for this particular language direction.

### 4.3 Fine-grained linguistic analysis

Additionally to the automatic scores, we provide here some fine-grained analysis based on particular linguistic categories, based on a linguistically-motivated test suite (Macketanz et al., 2022, 2021; Avramidis et al., 2020). The results can be seen in Table 3 and a more detailed view of the phenomena is displayed in Table 4. The model is particularly strong in *false friends*, which typically refers to lexemes that are identical in their phonological or orthographic form across two languages but have different meanings. It also performs relatively well in handling *non-verbal agreement*, i.e. ensuring that nouns and pronouns agree in gender, number and sometimes case across the sentence (particularly *substitution and coreference*), as well as in *lexical ambiguity*, where a word changes its meaning depending on a context, and *subordination* (particularly *adverbial and subject clause*). *Subordination* refers to the relationship between clauses where one clause is syntactically dependent on the main clause. However, it performs poorly in *punctuation* and particularly quotation marks, which means the model fails to correctly mark direct speech, quotations, or special terms. The low accuracy in *negation* is also particularly concerning, given the semantic importance of this category.

## 5 Conclusion and further work

We presented an entry participation of a new open-source community-based LLM. Despite some efforts to improve our LLM performance towards translation, the resulting model performs poorly as compared to other systems. Nevertheless, the challenges served as a motivation to unite community forces and initiate research on a new LLM task, which may be further improved in the future. Aside from the automatic scores, by applying a linguistically motivated test suite, we could gain some insights into the linguistic categories which perform better or worse. Further work may include more optimization towards translation, improvement of the prompting and post-processing mechanism and addition of more languages. A more direct comparison with models of similar parameter size (7B) should also be considered in the future.

## Acknowledgements

| phenomenon | items | acc |
|---|---|---|
| Ambiguity | 22 | 86.4 |
| Lexical ambiguity | 22 | 86.4 |
| Coordination & ellipsis | 124 | 60.5 |
| Gapping | 20 | 25.0 |
| Pseudogapping | 19 | 73.7 |
| Right node raising | 18 | 88.9 |
| Sluicing | 20 | 75.0 |
| Stripping | 23 | 39.1 |
| VP-ellipsis | 24 | 66.7 |
| False friends | 40 | 92.5 |
| Function word | 40 | 75.0 |
| Focus particle | 23 | 78.3 |
| Question tag | 17 | 70.6 |

| phenomenon | items | acc |
| --- | --- | --- |
| LDD & interrogatives | 207 | 76.3 |
| Extraposition | 18 | 55.6 |
| Inversion | 27 | 77.8 |
| Multiple connectors | 20 | 80.0 |
| Negative inversion | 20 | 80.0 |
| Pied-piping | 19 | 73.7 |
| Polar question | 18 | 77.8 |
| Preposition stranding | 19 | 57.9 |
| Split infinitive | 19 | 94.7 |
| Topicalization | 20 | 80.0 |
| Wh-movement | 27 | 81.5 |
| Lexical Morphology | 39 | 61.5 |
| Functional shift | 17 | 70.6 |
| Noun formation (er) | 22 | 54.5 |
| MWE | 123 | 76.4 |
| Collocation | 20 | 90.0 |
| Compound | 16 | 87.5 |
| Idiom | 20 | 40.0 |
| Nominal MWE | 20 | 75.0 |
| Prepositional MWE | 18 | 83.3 |
| Verbal MWE | 29 | 82.8 |
| Named entity & terminology | 112 | 77.7 |
| Date | 19 | 73.7 |
| Domainspecific Term | 18 | 83.3 |
| Location | 19 | 84.2 |
| Measuring unit | 21 | 76.2 |
| Onomatopeia | 15 | 53.3 |
| Proper name | 20 | 90.0 |
| Negation | 18 | 66.7 |
| Non-verbal agreement | 109 | 87.2 |
| Coreference | 35 | 88.6 |
| Genitive | 18 | 83.3 |
| Personal Pronoun Coreference | 13 | 92.3 |
| Possession | 27 | 81.5 |
| Substitution | 16 | 93.8 |
| Punctuation | 37 | 51.4 |
| Quotation marks | 37 | 51.4 |
| Subordination | 191 | 85.3 |
| Adverbial clause | 19 | 94.7 |
| Cleft sentence | 17 | 76.5 |
| Contact clause | 22 | 72.7 |
| Indirect speech | 19 | 89.5 |
| Infinitive clause | 19 | 84.2 |
| Object clause | 20 | 95.0 |
| Pseudo-cleft sentence | 19 | 78.9 |
| Relative clause | 39 | 89.7 |
| Subject clause | 17 | 82.4 |
| Verb semantics | 23 | 60.9 |
| Verb tense/aspect/mood | 3249 | 71.9 |
| Conditional | 20 | 70.0 |
| Ditransitive - conditional I progressive | 53 | 71.7 |
| Ditransitive - conditional I simple | 55 | 76.4 |
| Ditransitive - conditional II progressive | 56 | 48.2 |
| Ditransitive - conditional II simple | 54 | 77.8 |
| Ditransitive - future I progressive | 52 | 86.5 |
| Ditransitive - future I simple | 110 | 70.0 |
| Ditransitive - future II progressive | 55 | 34.5 |
| Ditransitive - future II simple | 51 | 29.4 |
| Ditransitive - past perfect progressive | 56 | 62.5 |
| Ditransitive - past perfect simple | 55 | 67.3 |
| Ditransitive - past progressive | 57 | 77.2 |
| Ditransitive - present perfect progressive | 57 | 75.4 |
| Ditransitive - present perfect simple | 51 | 80.4 |
| Ditransitive - present progressive | 55 | 85.5 |
| Ditransitive - simple past | 76 | 85.5 |
| Ditransitive - simple present | 50 | 84.0 |
| Gerund | 25 | 80.0 |
| Imperative | 15 | 46.7 |
| Intransitive - conditional I progressive | 27 | 92.6 |
| Intransitive - conditional I simple | 28 | 96.4 |
| Intransitive - conditional II progressive | 27 | 66.7 |
| Intransitive - conditional II simple | 29 | 69.0 |
| Intransitive - future I progressive | 30 | 83.3 |
| Intransitive - future I simple | 68 | 91.2 |
| Intransitive - future II progressive | 28 | 53.6 |
| Intransitive - future II simple | 35 | 48.6 |
| Intransitive - past perfect progressive | 30 | 46.7 |
| Intransitive - past perfect simple | 35 | 71.4 |
| Intransitive - past progressive | 32 | 81.3 |
| Intransitive - present perfect progressive | 29 | 82.8 |
| Intransitive - present perfect simple | 29 | 72.4 |
| Intransitive - present progressive | 61 | 85.2 |
| Intransitive - simple past | 35 | 80.0 |
| Intransitive - simple present | 38 | 68.4 |
| Modal | 288 | 71.5 |
| Modal negated | 304 | 75.0 |
| Reflexive - conditional I progressive | 35 | 74.3 |
| Reflexive - conditional I simple | 34 | 64.7 |
| Reflexive - conditional II progressive | 34 | 58.8 |
| Reflexive - conditional II simple | 34 | 76.5 |
| Reflexive - future I progressive | 30 | 60.0 |
| Reflexive - future I simple | 68 | 54.4 |
| Reflexive - future II progressive | 34 | 41.2 |
| Reflexive - future II simple | 33 | 39.4 |
| Reflexive - past perfect progressive | 35 | 42.9 |
| Reflexive - past perfect simple | 34 | 67.6 |
| Reflexive - past progressive | 33 | 87.9 |
| Reflexive - present perfect progressive | 32 | 68.8 |
| Reflexive - present perfect simple | 34 | 79.4 |
| Reflexive - present progressive | 33 | 75.8 |
| Reflexive - simple past | 33 | 78.8 |
| Reflexive - simple present | 31 | 61.3 |
| Transitive - future II progressive | 30 | 36.7 |
| Transitive - conditional I progressive | 30 | 86.7 |
| Transitive - conditional I simple | 27 | 85.2 |
| Transitive - conditional II progressive | 28 | 89.3 |
| Transitive - conditional II simple | 25 | 80.0 |
| Transitive - future I progressive | 30 | 73.3 |
| Transitive - future I simple | 57 | 84.2 |
| Transitive - future II simple | 32 | 65.6 |
| Transitive - past perfect progressive | 28 | 89.3 |
| Transitive - past perfect simple | 28 | 71.4 |
| Transitive - past progressive | 44 | 70.5 |
| Transitive - present perfect progressive | 27 | 88.9 |
| Transitive - present perfect simple | 29 | 79.3 |
| Transitive - present progressive | 39 | 84.6 |
| Transitive - simple past | 38 | 89.5 |
| Transitive - simple present | 34 | 88.2 |
| Verb valency | 114 | 65.8 |
| Case government | 14 | 85.7 |
| Catenative verb | 18 | 83.3 |
| Mediopassive voice | 22 | 54.5 |
| Passive voice | 19 | 78.9 |
| Resultative | 19 | 63.2 |
| Semantic roles | 22 | 40.9 |
| micro-average | 4448 | 72.8 |
| phen. macro-average | 4448 | 73.2 |
| categ. macro-average | 4448 | 73.0 |

Table 4: Performance of the Occiglot English-German model with regard to linguistically-motivated phenomena

# References

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. In-context Examples Selection for Machine Translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8857–8873, Toronto, Canada. Association for Computational Linguistics.

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, Aljoscha Burchardt, and Sebastian Möller. 2020. Fine-grained linguistic evaluation for state-of-the-art machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 346–356, Online. Association for Computational Linguistics.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. ArXiv:2005.14165 [cs].

Casimiro Pio Carrino, Marta R. Costa-jussà, and José A. R. Fonollosa. 2020. Automatic Spanish Translation of SQuAD Dataset for Multi-lingual Question Answering. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5515–5523, Marseille, France. European Language Resources Association.

Ryandito Diandaru, Lucky Susanto, Zilu Tang, Ayu Purwarianti, and Derry Wijaya. 2024. Could We Have Had Better Multilingual LLMs If English Was Not the Central Language? ArXiv:2402.13917 [cs].

Maxim Enis and Mark Hopkins. 2024. From LLM to NMT: Advancing Low-Resource Machine Translation with Claude. ArXiv:2404.13813 [cs].

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024a. Findings of the WMT24 general machine translation shared task:
the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar.

2024b. Preliminary wmt24 ranking of general mt systems and llms.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. OpenAssistant conversations - democratizing large language model alignment. In *Advances in neural information processing systems*, volume 36, pages 47669–47681. Curran Associates, Inc.

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, and Hans Uszkoreit. 2022. A linguistically motivated test suite to semi-automatically evaluate German–English machine translation output. In *Proceedings of the thirteenth language resources and evaluation conference*, pages 936–947, Marseille, France. European Language Resources Association.

Vivien Macketanz, Eleftherios Avramidis, Shushen Manakhimova, and Sebastian Möller. 2021. Linguistic evaluation for the 2021 state-of-the-art machine translation systems for German to English and English to German. In *Proceedings of the Sixth Conference on Machine Translation*, pages 1059–1073, Online. Association for Computational Linguistics.

Shivalika Singh, Freddie Vargus, Daniel D'souza, Börje Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura O'Mahony, Mike Zhang, Ramith Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Moura, Dominik Krzemiński, Hakimeh Fadaei, Irem Ergun, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Chien, Sebastian Ruder, Surya Guthikonda, Emad Alghamdi, Sebastian Gehrmann, Niklas Muennighoff, Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. 2024. Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11521–11567, Bangkok, Thailand. Association for Computational Linguistics.

Abdelhadi Soudi, Mohamed Hannani, Kristof Van Laerhoven, and Eleftherios Avramidis. 2024. Exploring the potential of large language models in adaptive machine translation for generic text and subtitles. In *Proceedings of the 17th workshop on building and using comparable corpora (BUCC) @ LREC-COLING 2024*, pages 51–58, Torino, Italia. ELRA and ICCL.

David Stap, Eva Hasler, Bill Byrne, Christof Monz, and Ke Tran. 2024. The Fine-Tuning Paradox: Boosting Translation Quality Without Sacrificing LLM Abilities. ArXiv:2405.20089 [cs].

Jörg Tiedemann. 2020. The Tatoeba Translation Challenge – Realistic Data Sets for Low Resource and Multilingual MT. In *Proceedings of the Fifth Conference*

*on Machine Translation*, pages 1174–1182, Online. Association for Computational Linguistics.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for Translation: Assessing Strategies and Performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15406–15427, Toronto, Canada. Association for Computational Linguistics.

Changtong Zan, Liang Ding, Li Shen, Yibing Zhen, Weifeng Liu, and Dacheng Tao. 2024. Building Accurate Translation-Tailored LLMs with Language Aware Instruction Tuning. ArXiv:2403.14399 [cs].

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting Large Language Model for Machine Translation: A Case Study. In *Proceedings of the 40th International Conference on Machine Learning*, pages 41092–41110. PMLR. ISSN: 2640-3498.

Biao Zhang, Philip Williams, Ivan Titov, and Rico Sennrich. 2020. Improving Massively Multilingual Neural Machine Translation and Zero-Shot Translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1628–1639, Online. Association for Computational Linguistics.

# CoST of breaking the LLMs

**Ananya Mukherjee**\*, **Saumitra Yadav**\*, **Manish Shrivastava**
MT-NLP Lab, LTRC, KCIS, IIIT Hyderabad, India
ananya.mukherjee@research.iiit.ac.in
saumitra.yadav@research.iiit.ac.in
m.shrivastava@iiit.ac.in

## Abstract

This paper presents an evaluation of 16 machine translation systems submitted to the Shared Task of the 9th Conference of Machine Translation (WMT24) for the English-Hindi (en-hi) language pair using our Complex Structures Test (CoST) suite. Aligning with this year's test suite subtask theme, "Help us break LLMs", we curated a comprehensive test suite encompassing diverse datasets across various categories, including autobiography, poetry, legal, conversation, play, narration, technical, and mixed genres.

Our evaluation reveals that **all the systems struggle significantly with the archaic style of text like legal and technical writings or text with creative twist like conversation and poetry datasets**, highlighting their weaknesses in handling complex linguistic structures and stylistic nuances inherent in these text types. Our evaluation identifies the strengths and limitations of the submitted models, pointing to specific areas where further research and development are needed to enhance their performance. Our test suite is available at https://github.com/AnanyaCoder/CoST-WMT-24-Test-Suite-Task.

## 1 Introduction

Neural Machine Translation (NMT) has seen substantial progress in recent years, achieving impressive quality that benefits many everyday applications. The advent of large language models (LLMs) has further enhanced translation capabilities. However, despite these advancements, there remain challenges that generic evaluation methods often fail to address. While traditional evaluations using random text samples might show overall success,

they may not reveal subtle issues where MT systems struggle, such as handling complex linguistic structures, idiomatic expressions, and diverse text types like conversations, poetry, legal documents, and technical writing. These flaws can be obscured by average performance metrics or overlooked entirely. A more systematic method for identifying linguistic issues in translation outputs involves using test suites or challenge sets to evaluate the system's performance on specific tasks. (Manakhimova et al., 2023). Test suites offer a standardized approach to evaluating MT systems, revealing strengths and weaknesses in handling complex text types.

In this context, we present the results of using test suites to analyze state-of-the-art machine translation systems across various categories. These evaluations were conducted as part of the theme "Help Us Break LLMs" for the 9th Conference on Machine Translation (WMT24). The test suites were used to evaluate systems submitted for the English-Hindi language pair.

We have curated a unique test suite comprising sentences from 9 categories across 16 sources to evaluate how large language models (LLMs) perform. The diversity of these categories allows us to assess the LLMs' capabilities beyond the typical news or generic domains, which often focus on reporting or narrative writing styles. Details of our test suite are provided in Section 2.

We perform reference-free and reference-based evaluations of the Hindi translations of this test suite, produced by 16 different machine translation (MT) systems submitted to the General Translation Task at WMT24 (Kocmi et al., 2024a). For reference-less evaluation, we employ COMET-Kiwi (Rei et al., 2022), while (Papineni et al., 2002),

---

\* Authors contributed equally

chrF (Popović, 2015, 2017), MEE4 (Mukherjee et al., 2020; Mukherjee and Shrivastava, 2023), BERTScore (Zhang* et al., 2020), and COMET (Rei et al., 2020) are used for reference-based evaluation. Professional English-to-Hindi translators provide the reference translations. Our results indicate that, for the English-to-Hindi language pair, **LLMs show weaker performance on datasets related to poetry, legal, and conversational content**. Details of our evaluation experiments are discussed in Section 3, and our analysis is presented in Section 4.

## 2 CoST: Complex Structure Testsuite

Table 1 depicts the dataset categories and the distribution within our test suite. The "Original" column presents the initial count of selected sentences for each category, as gathered from the datasets. The last column, "CoST," displays the final count of sentences included in the test suite. Our test suite is designed to evaluate translations across

- Multiple Writing Style: Prose, Conversation, Autobiography, Legal Writing, Literary Narrative and Technical Documents.

- Lexical Choice: As we are sampling test suites from various domains, there is a decent mixture of domain-specific words, e.g. Legal Text, Technical Text, etc.

In total, 1,947 English sentences were selected based on criteria such as sentence length, depth of dependency tree, combination of noun phrases, verb phrases, named entities, etc. Ensuring a test suite containing sentences with good representation from simple to complex structures.

## 3 Evaluation Strategy

To evaluate the performance of the 16 submitted MT systems, we performed both automatic and manual evaluations.

### 3.1 Automatic Evaluation

In automatic evaluation, we leveraged both reference-less and reference-based metrics.

| Category | Dataset | Original | CoST |
|---|---|---|---|
| poetry | Kabir ke Dohe | 11 | 9 |
| | Amir Khusro | 9 | 9 |
| narration | ShortStories | 177 | 72 |
| | Post Office | 440 | 10 |
| | Glimpses of Bengal | 101 | 64 |
| | The Home and the World | 236 | 183 |
| | The gardener | 277 | 27 |
| | Abridged Merchant of Venice | 63 | 31 |
| | Christmas Carole | 923 | 308 |
| legal | Legal Text | 2862 | 638 |
| | IIT Bombay Jud | 167 | 83 |
| mix | IN22 | 570 | 241 |
| conversation | Friends | 77 | 53 |
| play | King Of Dark Chamber | 35 | 22 |
| autobiography | My Reminiscences | 109 | 110 |
| Technical | Technical Papers | 185 | 87 |
| | **Total** | **6242** | **1947** |

Table 1: Data Statistics of **CoST**.

### 3.1.1 Reference-less Evaluation

For the reference-less automatic evaluation, we utilize COMETKIWI (Rei et al., 2022) scores, which offer quality estimation scores derived from the source sentence and MT output.

### 3.1.2 Reference-based Evaluation

With the help of professional English-to-Hindi translators, we also provide one gold reference translation for each source sentence in the test suite. We evaluate the machine translation outputs against these references using BLEU (Papineni et al., 2002), chrF (Popović, 2015, 2017), MEE4 [1] (Mukherjee et al., 2020; Mukherjee and Shrivastava, 2023), BERTScore (Zhang* et al., 2020), and COMET (Rei et al., 2020).

### 3.2 Manual Evaluation

The manual analysis was done by professional native speakers. They were instructed to identify mistranslations and hallucinations and make note of other translation errors like wrong post positions to get more nuanced information regarding the performance of the systems.

## 4 Results and Analysis

The results of the automatic evaluation are reported in Table 2. Ranks are shown in parentheses for each metric, where (1) is the highest rank. It is clearly evident that evaluations from all the metrics rank TranssionMT as the best system, followed by ONLINE-B

---

[1]https://www.kaggle.com/ananyacoder/mee4-metric-run

and Claude-3.5. In contrast, CycleL is ranked the lowest, preceded by IKUN-C and IKUN. We also observe that according to the Preliminary WMT24 Ranking of General MT Systems and LLMs (Kocmi et al., 2024b), Unbabel-Tower70B is listed as the top performer. However, its performance decreases on CoST. For more category-wise informative results, we looked at the performance of systems for each category using lexical-based metric (Figure 2 and 3), embedding- based metric (Figure 4 and 5) and supervised metric (6) and (Figure 1). These results illustrate that **all systems underperform with poetry, legal, and conversation data**. In contrast, the systems consistently exhibit strong performance with autobiography, play, and mixed (IN22) data.

The analysis shows a clear trend, i.e., systems struggle with specific genres like poetry, legal, and conversation while excelling in narrative styles such as autobiography and play. This suggests that the training data for these systems may be heavily skewed towards narrative writing, hence strong performance in those areas. The sub-par performance in poetry, conversational and legal texts might reflect challenges in handling diverse linguistic and stylistic features that are less prevalent in the training data.

## 4.1 Qualitative Analysis

These manual assessments are carried out by professional Hindi speakers who hold graduate-level qualifications and possess good knowledge in the domains covered by our test suite.

### 4.1.1 Handling Named Entities

Source: *Labanya said to her sister in soothing tones : " Don't be upset about it , dear ; I will see what I can do to prevent it . "*

Most models successfully translated "Labanya" correctly, preserving the original name. However, the outputs from Claude-3.5, GPT-4, NVIDIA-NeMo, ONLINE-A, Unbabel-Tower70B, and ZMT show variations or distortions of the name, indicating potential issues with **name recognition or transliteration** in these models.

In another instance, IKUN-C, IKUN, Llama3-70B, NVIDIA-NeMo, ONLINE-A, Unbabel-Tower70B, and ZMT systems have translated 'Phoebe' as Phob, Phobey, Phoyeb, Phoyebe; surprisingly ONLINE-G has generated चाँद (meaning moon, as Phobe is one of the moons of Saturn).

### 4.1.2 Spelling and Typological Errors

Except for Llama3-70B, IOL_Research, and CommandR-plus, all other models tend to generate हूं instead of हूँ, indicating a recurring spelling error in their outputs.

### 4.1.3 Omissions

The Hindi translations produced by the IKUN and IKUN-C systems consistently suffer from **incompleteness**, often leaving out key parts of the original sentences, undermining the accuracy and reliability of the translations, making them less effective for conveying the full meaning of the source text.

### 4.1.4 Incorrect Lexical Word Choices

Choosing the right word in translation is crucial for preserving the essence, tone, and intention of the original sentence. For instance, Unbabel-Tower70B accurately translates "well," whereas all other systems translate it as "alright" or "okay." These alternatives do not fit the context as well, thereby **affecting the tone and overall quality** of the translation.

Source: *I'd be pulling up shoots of grass to use them to check the wind, and looking at maps of ports and piers and roads.*

However, Aya23 and IOL_Research translate it as "removing," while the remaining systems use "pull." These variations of "remove" and "pull" slightly **affect the accuracy and well-formedness** of the Hindi translation.

## 5 Conclusion

This paper evaluates translations from 16 MT systems submitted to the General Translation Shared Task WMT24 on **Complex Structures Test** suite which was designed to cover various writing styles and domains beyond the typical news and generic data, consisting 1,947 unique sentences selected for their lexical and structural diversity. We conducted automatic reference-based, automatic reference-free, and manual evaluations. Our thorough analysis reveals significant limitations in these LLMs,

Figure 1: System-wise plots of average COMET-KIWI Scores for each category.

| | reference-free | reference-based | | | | |
|---|---|---|---|---|---|---|
| **System** | **COMET-KIWI** | BLEU | chrF | MEE4 | BERTScore | COMET |
| **TranssionMT** | 0.815 (1) | 68.399 (1) | 81.577 (1) | 0.903 (1) | 0.942 (1) | 0.835 (1) |
| **Claude-3.5** | 0.815 (1) | 43.321 (3) | 66.385 (3) | 0.85 (3) | 0.898 (3) | 0.803 (3) |
| **ONLINE-B** | 0.814 (2) | 67.733 (2) | 80.768 (2) | 0.898 (2) | 0.933 (2) | 0.83 (2) |
| Unbabel-Tower70B | 0.809 (3) | 38.634 (6) | 62.811 (6) | 0.842 (5) | 0.886 (5) | 0.799 (4) |
| Llama3-70B | 0.791 (4) | 34.164 (9) | 58.612 (8) | 0.83 (6) | 0.874 (7) | 0.767 (5) |
| IOL_Research | 0.79 (5) | 32.991 (10) | 57.244 (10) | 0.825 (8) | 0.869 (8) | 0.765 (6) |
| ZMT | 0.785 (6) | 42.277 (5) | 65.614 (5) | 0.843 (4) | 0.893 (4) | 0.75 (9) |
| ONLINE-A | 0.785 (6) | 42.324 (4) | 65.637 (4) | 0.843 (4) | 0.893 (4) | 0.75 (9) |
| GPT-4 | 0.785 (6) | 31.795 (11) | 57.227 (11) | 0.826 (7) | 0.868 (9) | 0.755 (8) |
| CommandR-plus | 0.785 (6) | 29.088 (12) | 54.918 (12) | 0.816 (10) | 0.858 (10 | 0.757 (7) |
| Aya23 | 0.761 (7) | 27.938 (13) | 53.473 (13) | 0.81 (11) | 0.852 (11) | 0.728 (10) |
| ONLINE-G | 0.735 (8) | 35.952 (7) | 60.861 (7) | 0.825 (8) | 0.875 (6) | 0.669 (12) |
| NVIDIA-NeMo | 0.734 (9) | 34.635 (8) | 57.977 (9) | 0.821 (9) | 0.868 (9) | 0.689 (11) |
| IKUN-C | 0.658 (10) | 10.89 (15) | 38.711 (14) | 0.693 (12) | 0.752 (12) | 0.591 (13) |
| IKUN | 0.574 (11) | 12.181 (14) | 36.159 (15) | 0.657 (13) | 0.731 (13) | 0.546 (14) |
| CycleL | 0.366 (12) | 1.77 (16) | 16.476 (16) | 0.347 (14) | 0.665 (14) | 0.33 (15) |

Table 2: System-wise ranking based on reference-free and reference-based metrics. Top 3 are highlighted in bold. Ranks are mentioned in brackets. The rows are colour coded highlighting the top scores in green and low scores in red.



Figure 2: Category-wise plots of average BLEU Scores for all the submitted MT systems.

Figure 3: Category-wise plots of average chrF Scores for all the submitted MT systems.



Figure 4: Category-wise plots of average BERTScore Scores for all the submitted MT systems.

Figure 5: Category-wise plots of average MEE4 Scores for all the submitted MT systems.



Figure 6: Category-wise plots of average COMET Scores for all the submitted MT systems.

particularly in translating poetry, conversational, and legal texts. Additionally, our manual review uncovered issues such as incorrect word choices, spelling errors, and poor handling of named entities. Despite their advancements, these LLMs show notable weaknesses in handling diverse and complex linguistic contexts. This highlights the need for continued refinement and broader training data to improve their performance across a wider range of text types and domains.

# References

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024a. Findings of the WMT24 general machine translation shared task:
the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024b. Preliminary wmt24 ranking of general mt systems and llms.

Shushen Manakhimova, Eleftherios Avramidis, Vivien Macketanz, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian MÃ¶ller. 2023. Linguistically motivated evaluation of the 2023 state-of-the-art machine translation: Can chatgpt outperform nmt? In *Proceedings of the Eighth Conference on Machine Translation*, pages 224–245, Singapore. Association for Computational Linguistics.

Ananya Mukherjee, Hema Ala, Manish Shrivastava, and Dipti Misra Sharma. 2020. Mee: An automatic metric for evaluation using embeddings for machine translation. In *2020 IEEE 7th International Conference on Data Science and Advanced Analytics (DSAA)*, pages 292–299. IEEE.

Ananya Mukherjee and Manish Shrivastava. 2023. Mee4 and xlsim : Iiit hyd's submissions' for

wmt23 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 798–803, Singapore. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

# WMT24 Test Suite: Gender Resolution in Speaker-Listener Dialogue Roles

**Hillary Dawkins**  **Isar Nejadgholi**  **Chi-kiu Lo** 羅致翹
Digital Technologies Research Centre
National Research Council Canada (NRC-CNRC)
{hillary.dawkins, isar.nejadgholi, chikiu.lo}@nrc-cnrc.gc.ca

## Abstract

We assess the difficulty of gender resolution in literary-style dialogue settings and the influence of gender stereotypes. Instances of the test suite contain spoken dialogue interleaved with external meta-context about the characters and the manner of speaking. We find that character and manner stereotypes outside of the dialogue significantly impact the gender agreement of referents within the dialogue. https://github.com/hillary-dawkins/wmt24-gender-dialogue.

## 1 Introduction

Gender bias and gender effects in machine translation are prevalent in translation directions where gender relevancy increases from source to target language (Savoldi et al., 2021; Barclay and Sami, 2024; Savoldi et al., 2023). English has minimal morphological effects caused by natural gender, whereas many languages (e.g. French, Spanish, Czech, Icelandic, German) have grammatical gender cases for various parts of speech which sometimes need to align with natural gender for animate nouns. For example "I am happy" in the source language English has divergent translations in the target language French ("Je suis heureux/heureuse") depending on the natural gender of the speaker. The consequence is that gender-alignment errors can easily arise in such translation directions. Furthermore, stereotypes are known to drive gender agreement (e.g., systems may tend to prefer the translation "Je suis jolie" over "Je suis joli" for "I am pretty" despite incomplete gender context) (Sólmundsdóttir et al., 2022), and these stereotype effects can persist even when unambiguous gender information is provided (Stanovsky et al., 2019; Troles and Schmid, 2021; Kocmi et al., 2020).

Typically, these gender effects are studied in isolation or semantically-bleached settings (as in the above examples). There it is known that the internal characteristics of adjective words, such as

the gender stereotype, sentiment, and type (appearance or character), are significant factors influencing the choice of gender agreement in translation (Sólmundsdóttir et al., 2022). However, the need for gender agreement also occurs in more complex settings, such as over long ranges, and passages involving multiple potential referents.

Due to increasing interest in paragraph-level translation and literary domains, here we assess the challenge of speaker-listener role resolution in literary dialogue settings. In particular, the gender of the speaker and listener must be resolved correctly to obtain a correct translation, and we suppose that gender stereotype effects can further add to the task difficulty. We find that stereotypical character descriptions and manners of speaking are significant influences on the gender alignment, generally overshadowing the internal adjective traits.

## 2 Test Suite Description

This test suite measures the gender resolution tendencies of machine translation systems in literary-style dialogue settings. In this setting, spoken dialogue (in quotations or otherwise delimited) is interleaved with meta-context about the dialogue (e.g., the speaker, the listener(s), and character and environment descriptions). When spoken dialogue refers to a person, a challenge arises in resolving the referent given the meta-context. The test suite includes three target languages (Spanish, Czech, and Icelandic), where the gender of the referent affects the correct translation.

Here, we focus on two-person conversations, where adjectives are used within dialogue to describe either the speaker or the listener. Within a single source passage, both characters may take on both the speaker and listener roles at times. Since adjectives are gender-neutral in the source language (English), the gender of the adjective's referent must be determined from the meta-context, if possible. The test suite contains inputs where

the gender remains unknown given the complete context (termed gender-ambiguous cases), and inputs where the gender can be unambiguously resolved given the complete context (termed gender-determined cases).

The test suite contains a handful of template types (each detailed in Appendix A) to assess the influence of stereotype cues in the meta-context and the structural features of the passage. Stereotype features include character descriptions and the manner of speaking (controlled using adverbs). Structural features include the number of referents in a single passage, partial or complete gender information, first- or third-person speakers, and adjective repetition. Some challenging features of the templates include adjectives that appear before the referent is introduced, and repeated adjectives referring to different entities.

The templates use vocabulary sets for adjectives ($n = 350$), gender-stereotyped adverbs ($n = 29$), and gender-stereotyped occupation words ($n = 44$). Each adjective is labeled with its gender stereotype (M/F/neutral), sentiment (positive/negative/neutral), and type (character/appearance). The full vocabulary set with annotations is released as part of the test suite contribution.

## 3 Methodology

The adjective translations are extracted from the target languages and processed using dictionary searches[1] to obtain the gender agreement label. The advantage of using dictionary searches over automated morphological gender taggers is that irregular adjectives (e.g. "rosa" in Spanish) are correctly classified, and the use of different parts of speech or out-of-dictionary words can also be monitored. For example, the use of a gender-neutral noun phrase or direct substitution of an English word should be counted as a neutral label for our purposes. Only when a translated word is not found in any dictionary search, is it passed to auto gender-tagging based on its morphological features (e.g. an "o" vs. an "a" ending in Spanish). This second pass allows for (possibly hallucinated) out-of-dictionary words to be included in the analysis, but only if

they strongly resemble a regular adjective form (e.g. "víktur" in Icelandic may be derived from the English source word "victorious", but clearly a masculine adjective ending has been chosen in translation). A small portion of words remains unclassified after both passes are complete, meaning that they neither exist in the dictionary nor resemble a regular adjective in the target language. The fine-grained annotations for each extracted translation, in addition to the final gender label (one of M, F, N, or unclassified), are released with the test suite results for further analysis.

The scope of analysis in this paper is limited to the subset of M- and F-labeled translations. That is, when a gendered adjective form *is* chosen by a translation system, we are interested in the factors that influence this choice, and the corresponding translation errors that occur when an adjective form does not match the referent's gender. To this end, results throughout the paper are presented in three ways.

When the gender of a referent is unknown, we report the proportion of masculine and feminine adjective declensions to observe the system's tendencies in ambiguous settings. When the gender of a referent is known, we report the accuracy of the adjective declensions. Typically, the underlying effect (e.g., the influence of stereotypes) is the same in both cases. However, it is important to know that the effect persists even when unambiguous gender context is available. Both proportion and accuracy results are always reported using balanced subsets[2] of the relevant test suite subset.

Lastly, we wish to understand the relative importance of factors that influence the system's choice of gender agreement. To do so, we perform regression analyses where the dependent variable to predict is the gender declension of the translation, and independent variables include both internal adjective factors (the gender stereotype, sentiment, and type), and external factors that are introduced through the meta context (e.g. character descriptions). The regression coefficients are reported with significance levels.

## 4 Gender-Stereotyped Manner

Firstly, we observe that the manner of speaking in literary dialogue settings can significantly affect

*(a)* "I think I'm <u>courageous</u>," I said.

*(b)* "I think I'm <u>statečný</u> (**M↑↑**)," I said sauvely.

*(c)* "I think I'm <u>statečná</u> (**F↑↑**)," I said gently.

Figure 1: Gender-stereotyped adverbs outside of the dialogue affect the adjective's gender agreement with the speaker within the dialogue. Source sentences in English include instances without adverbs (a) and with stereotypically masculine (b) or feminine adverbs (c). When translated to the target language, adjectives tend to align with the stereotype (adjectives shown here in Czech).

the gender prediction of the speaker. Furthermore, this influence is susceptible to gender stereotypes. Refer to the example shown in Figure 1.

Within the **Stereo-Adverb** test suite subset, all adjectives refer to a first-person speaker (I), and therefore the natural gender of the adjective's referent is ambiguous in the source language. We report the proportion of male declensions on subsets (a) with no adverb, (b) a male-stereotyped adverb, and (c) a female-stereotyped adverb (full results in Appendix B). The majority of systems display a difference greater than 10% when the adverb switches from male- to female-stereotyped. The systems with the largest effects are shown in Table 1. Note that the most affected systems include those that defy the usual default-male agreement in ambiguous gender cases in the baseline setting (i.e., in the absence of any adverb). Here we see that the default-female agreement is unstable with respect to stereotype cues.

To compare the influence of speaking manner to the influence of internal adjective traits, we perform regression analysis where the dependent variable to predict is the gender declension. Independent variables are the gender stereotype label of the adverb, and the gender stereotype, the sentiment, and the type (appearance or character) of the adjective. The analysis shows that adverb influence is comparable or stronger than these internal adjective characteristics within this test suite (see Table 9).

## 5 Gender-Stereotyped Characters

Secondly, we observe that character descriptions that align with socially held stereotypes impact gender resolution within spoken dialogue. Refer to the examples shown in Figure 2.

Within the **Stereo-Character** test suite subset,

all adjectives refer to one of two characters that have been given some stereotypical descriptions using both occupations and attributive adjectives. Template variations include single-speaker dialogue, where adjectives refer to either the speaker (I) or listener (you) (see template 3), and two-speaker conversations where both participants are referenced by each speaker (see template 4).

In ambiguous gender cases (Figure 2a), we report the stereotype effect again by looking at the tendency of the system to choose either a female or male adjective declension depending on the referent stereotype (full results in Appendix B). Characters that are described by male-leaning gender stereotypes are very likely to receive a masculine adjective, whereas the use of feminine adjectives increases for female-stereotyped characters (pushing against the default-male baseline), as shown in Table 2 for the most affected systems.

Furthermore, we find that this effect persists in determined gender cases (Figure 2b) such that agreement accuracy can drop significantly when the actual gender opposes a socially-held stereotype. We report this observation as the difference in accuracy between the PRO and ANTI template subsets (full results in Appendix B). Approximately half of the tested systems are not robust to stereotype cues even when the correct, unambiguous gender context is provided. The most affected systems are shown in Table 3. As with the stereotyped adverb effect, we perform a regression analysis to probe the relative influence of stereotypical character descriptions compared to the internal adjective factors. We find that the character descriptions are much more impactful on the adjective's gender form than the internal adjective traits within this dialogue setting (see Table 15).

## 6 "Opposite" or "Same" Binary Gender Speaker Bias

Finally, in the absence of any gender stereotype effects, we assess the overall "vanilla" difficulty of this gender resolution task in non-challenge settings and the influence of different structural elements in the source input. In doing so, we observe that an "opposite" or same binary gender bias exists. That is, in dialogue settings between two speakers, some systems strongly predict one speaker to be male and the other female, while other systems strongly prefer same-gender pairs. This observation holds in both ambiguous and determined cases.

| System | $F$ | $M$ | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|--------|-----|-----|-------|-------|-------|-------|--------|
| CUNI-MH | 0.703 | 0.297 | 0.379 | 0.621 | 0.950 | 0.050 | **0.571** |
| ONLINE-W | 0.591 | 0.409 | 0.387 | 0.613 | 0.884 | 0.116 | **0.497** |
| CommandR-plus | 0.340 | 0.660 | 0.120 | 0.880 | 0.554 | 0.446 | **0.434** |
| Aya23 | 0.370 | 0.631 | 0.187500 | 0.813 | 0.612 | 0.388 | **0.425** |

Table 1: **Gender-Stereotyped Manner**: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Adverb** test suite subset for the most affected translation systems in the English to Czech translation direction. All adjectives self-refer to the speaker of unknown gender. Subscripts ($M$ and $F$) denote the use of gender-stereotyped adverbs to control the manner of speaking (e.g., $M_F$ denotes the proportion of adjectives with a male declension within instances using a stereotypically feminine adverb, as shown in Figure 1 example (c)). The unsubscripted results refer to no adverb (as shown in Figure 1 example (a)). The overall strength of the adverb effect is reported using the difference $\Delta M_{M-F}$.

The prudish secretary smiled. "I think I'm adj (**F↑↑**)," they said to the grizzled judge.

The prudish secretary smiled. "I think you're adj (**M↑↑**)," they said to the grizzled judge.

The pretty nurse smiled. "I think I'm adj (**F↑↑**) and you're adj (**M↑↑**)," they said.
The strong doctor laughed back. "No, you're adj (**F↑↑**) and I'm adj (**M↑↑**)," they replied.

(a) Ambiguous cases: Adjectives refer to characters of unknown gender.

**PRO** (natural gender aligns with stereotypical character): (**Accuracy↑↑**)

The chaste secretary smiled. "I think I'm adj (**F↑↑**)," she said to the handsome analyst.

**ANTI** (natural gender opposes the stereotypical character): (**Accuracy↓↓**)

The chaste secretary smiled. "I think I'm adj (**F↑**)," he said to the handsome analyst.

(b) Determined cases: Adjectives refer to characters of known gender. The known gender either aligns with the stereotype (PRO) or opposes the stereotype (ANTI).

Figure 2: Gender-stereotyped character descriptions outside of the dialogue affect the adjective's gender agreement.

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|--------|-------|-------|-------|-------|--------|
| Claude-3.5 | 0.000 | 1.000 | 0.391 | 0.609 | **0.391** |
| CommandR-plus | 0.012 | 0.988 | 0.401 | 0.598 | **0.390** |
| Aya23 | 0.122 | 0.878 | 0.429 | 0.571 | **0.307** |
| Unbabel-Tower70B | 0.058 | 0.942 | 0.359 | 0.640846 | **0.302** |
| GPT-4 | 0.000 | 1.000 | 0.274 | 0.726 | **0.274** |

Table 2: **Gender-Stereotyped Characters**: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Character-Amb** test suite subset (Figure 2a) for the most affected translation systems in the English to Spanish translation direction, partitioned by the referent's gender stereotype (denoted by subscripts). The true gender of the referent is unknown, but the choice of declension is affected by the stereotypical character description. The overall strength of the character description effect is reported by the difference $\Delta M_{M-F}$.

| System | Accuracy (PRO) | Accuracy (ANTI) | $\Delta$(PRO, ANTI) |
|---|---|---|---|
| ONLINE-W | 0.985 | 0.414 | **0.571** |
| GPT-4 | 0.990 | 0.527 | **0.463** |
| Aya23 | 1.000 | 0.655 | **0.345** |
| IKUN | 0.975 | 0.702 | **0.273** |

Table 3: **Gender-Stereotyped Characters**: The accuracy in gender-adjective agreement on the **Stereo-Character-Det** test suite subset (Figure 2b) for the most affected translation systems in the English to Spanish translation direction. The true gender of the character either aligns with (PRO) or opposes (ANTI) the stereotypical description. The presence of stereotypical character descriptions can significantly decrease the gender translation accuracy.

Refer to the examples shown in Figure 3.

In ambiguous gender cases, we can observe this effect as the proportion of adjective declension choices conditioned on the known gender of the second character in the conversation (Figure 3a). Note that adjectives may still either refer to the speaker or listener, and both types are affected by the presence of a second known gender. Full results are shown in Appendix B, and a summary of the most affected systems is shown in Table 4.

In determined gender cases, the tendency to assume either the same or opposite binary gender pairs manifests as decreased accuracy in cases that oppose this assumption. We report the accuracy in adjective agreement on test subsets where (a) only one gender is specified (Figure 3c), (b) both genders are specified and are opposite, (c) both genders are specified and are the same (Figure 3b). Subset (a) is usually easiest for most systems because the same or opposite gender effect is not possible. The difference in accuracy between subsets (b) and (c) indicates the strength and direction of this effect. Full results are shown in Appendix B and a summary is shown in Table 5.

We note that the observed decrease in accuracy on gender pairings that oppose the system's presupposition is being driven by two features within our templates: 1. Adjectives that occur before their referent (if reading left to right), and 2. A consistency effect. For example, refer to the two examples shown in Figure 3b. In both cases, the first adjective to translate occurs before it's referent, but after the gender of the speaker is known. Adjectives in this position are very likely to align with the same or opposite gender of the speaker in affected systems, depending on the effect direction. Following the incorrect translation of the first adjective, we observe that the adjective in the last position is likely to also be incorrect, possibly owing to a consistency effect since these refer to the same entity.

Using regression analysis, we predict the adjective declension conditioned on structural factors: the gender of the other speaker, whether the referent is the speaker (I) or the listener (you), the gender choice in preceding adjectives that refer to the same entity (consistency), and whether the adjective occurs before the referent is introduced ("look-ahead" position), as well as the internal traits of the adjective as always, and the true gender label for determined cases. Controlling for internal traits, the correct gender label, and the default masculine baseline, we observe that both "look-ahead" and referent role (listener) are influential structural factors affecting the task difficulty (refer to table 27).

# 7 Future Work

Here, the scope of analysis is limited to the cases where a translation system has chosen either a masculine or feminine adjective form, and ignores those cases where a neutral translation strategy was used instead. However, the labeling methodology as described in Section 3 does produce a test suite with annotated neutral labels as well. The observed neutral strategies vary by target language and include the use of adjectives with the same form for the female and male gender cases (e.g. regular adjectives ending in "e" in Spanish, or "í" in Czech), the use of the neuter gender case if it exists (as in Czech and Icelandic), direct substitution of the gender-neutral source (English) adjective, the use of alternative forms (e.g. translated adjectives ending "o/a" in Spanish or "(ur)" in Icelandic), and the use of noun phrases in place of adjectives, which may be gender-neutral depending on the target language. Some of these strategies may be considered to be more correct than others (i.e. applying the neuter gender case to a person is not grammatically correct, but may still be preferred to misgendering in ambiguous cases).

> I smiled. "I think I'm adj (F↑)," I said.
> He laughed back. "No, you're adj (F↑)," he replied.

> The man smiled. "I think you're adj (F↑↑↑)," he said.
> I laughed back. "No, I'm adj (F↑↑)," I replied.

> The man smiled. "I think I'm adj and you're adj (F↑↑↑)," he said.
> I laughed back. "No, you're adj and I'm adj (F↑↑)," I replied.

> I smiled. "I think I'm adj (F↑) and you're adj," I said.
> He laughed back. "No, you're adj (F↑) and I'm adj," he replied.

(a) Ambiguous cases: Adjectives refer to a character of unknown gender, while the gender of the second character in the conversation is known (male in these examples). Adjectives referring to the gender-ambiguous character are more likely to agree with the opposite gender of the speaker (i.e., take feminine forms in these examples).

> The man smiled. "I think you're adj (acc↓↓↓)," he said.
> He laughed back. "No, I'm adj (acc↓↓)," he replied.

> The man smiled. "I think I'm adj and you're adj (acc↓↓↓)," he said.
> He laughed back. "No, you're adj (acc↓↓) and I'm adj (acc↓)," he replied.

(b) Determined cases where the gender of both speakers is known. Accuracy decreases for same-gender pairs due to the opposite binary gender effect.

> I smiled. "I think you're adj (acc↑↑)," I said.
> She/he laughed back. "No, I'm adj (acc↑↑)," she/he replied.

> The wo/man smiled. "I think I'm adj (acc↑↑) and you're adj," she/he said.
> I laughed back. "No, you're adj (acc↑↑) and I'm adj," I replied.

(c) Determined cases where the gender of one speaker is known. Accuracy is generally high in the absence of a second gender (i.e., the opposite binary gender effect is not possible).

Figure 3: The opposite binary gender effect is present in both ambiguous (a) and determined (b) cases. Determined cases with a single known gender (c) are unchallenging despite having the same structural components (i.e. both speaker (I) and listener (you) resolutions, and need to "look ahead" in the text to find the adjective's referent). All effects are the same but flipped for systems that prefer same-gender speaker pairs.

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|
| Claude-3.5 | 0.419 | 0.581 | 0.074 | 0.926 | **-0.346** |
| CommandR-plus | 0.764 | 0.236 | 0.426 | 0.574 | **-0.338** |
| IKUN-C | 0.292 | 0.708 | 0.703 | 0.297 | **0.410** |
| IKUN | 0.256 | 0.744 | 0.726 | 0.274 | **0.470** |

Table 4: **Opposite or Same Binary Gender Effect**: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Structure-Amb** test suite subset (Figure 3a) for the most affected systems in the English to Spanish translation direction. All adjectives refer to someone of an unknown gender in conversation with someone of a known gender (where that known gender is denoted by the subscripts). Systems Claude-3.5 and CommandR-plus show the greatest tendency to assume opposite-gender speaker pairs ($\Delta M_{M-F} \ll 0$), and systems IOL-Research and IKUN show the greatest tendency to assume same-gender speaker pairs ($\Delta M_{M-F} \gg 0$).

| System | Acc (one gender) | Acc (same genders) | Acc (opp genders) | $\Delta$(same, opp) |
|---|---|---|---|---|
| CommandR-plus | 0.987 | 0.797 | 0.991 | **-0.194** |
| Llama3-70B | 0.957 | 0.806 | 0.977 | **-0.171** |
| ONLINE-A | 0.734 | 0.828 | 0.668 | **0.160** |
| ONLINE-G | 0.726 | 0.827 | 0.625 | **0.202** |

Table 5: **Opposite or Same Binary Gender Effect**: The accuracy in gender-adjective agreement on the **Structure-Det** test suite subset (Figures 3c and 3b) for the most affected systems in the English to Spanish translation direction. The second speaker in the conversation is either unknown (one gender subset), the same, or opposite to the adjective referent of known gender. Systems with an opposite binary gender effect suffer on the same-gender subset such that the difference in accuracy $\Delta$(same, opp) $\ll 0$, and systems with a same-gender preference suffer on the opposite-gender subset such that the difference in accuracy $\Delta$(same, opp) $\gg 0$.

Further analysis is needed to understand how often neutral strategies are used in both the ambiguous and determined gender cases, and what factors influence a translation system's choice or ability to use a neutral strategy (Savoldi et al., 2024; Piergentili et al., 2023; Lauscher et al., 2023).

# 8 Conclusion

In conclusion, this test suite provides an opportunity to study the challenging task of referent resolution within literary-style dialogue settings. When spoken dialogue refers to characters described outside of dialogue in the meta-context, it adds an extra layer of complexity to the gender agreement task. Here we focus on language directions that are prone to gender agreement errors due to greater gender relevancy in the target language than the source language. We find that stereotypical character descriptions and manners of speaking are significant influences for some translation systems. Furthermore, some systems strongly prefer to resolve two-person conversations as same- or opposite-gender pairs. All observed effects are present in both ambiguous and determined gender cases.

# Limitations

This test suite uses simple templates to study the influence of structural factors in a controlled manner. Although templates are varied and contain quite a few structure variables, they do not represent the diversity or complexity of real literary settings. Having identified the stereotype effects and challenge features within this test suite, future work could compile a real in-the-wild literary dialogue test suite by seeking out instances with these features of interest.

The primary limitation of this work is the focus on binary gender. All determined gender cases within the test suite are either male or female, and the analysis of chosen gender declensions is limited to masculine and feminine forms. This is partially due to the availability of known stereotypes for binary gender, and partially due to the binary nature of gender morphology in the target languages (even if neuter grammatical gender exists, it does not apply to animate nouns). Future work should investigate the use of neutral strategies when gender is unknown as a way to avoid misgendering non-binary referents.

# Ethics Statement

As discussed in the Limitations section, the focus on binary gender throughout the paper is a serious ethical concern, and we stress here that similar research questions are applicable to non-binary genders. We hope that the analysis presented here and the test suite results might encourage the inclusion of non-binary natural gender in future work.

# Acknowledgements

# References

Peter J Barclay and Ashkan Sami. 2024. Investigating markers and drivers of gender bias in machine translations.

Tom Kocmi, Tomasz Limisiewicz, and Gabriel Stanovsky. 2020. Gender coreference and bias evaluation at WMT 2020. In *Proceedings of the Fifth Conference on Machine Translation*, pages 357–364, Online. Association for Computational Linguistics.

Anne Lauscher, Debora Nozza, Ehm Miltersen, Archie Crowley, and Dirk Hovy. 2023. What about "em"? how commercial machine translation fails to handle (neo-)pronouns. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 377–392, Toronto, Canada. Association for Computational Linguistics.

Andrea Piergentili, Beatrice Savoldi, Dennis Fucci, Matteo Negri, and Luisa Bentivogli. 2023. Hi guys or hi folks? benchmarking gender-neutral machine translation with the GeNTE corpus. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14124–14140, Singapore. Association for Computational Linguistics.

Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. Gender bias in machine translation. *Transactions of the Association for Computational Linguistics*, 9:845–874.

Beatrice Savoldi, Marco Gaido, Matteo Negri, and Luisa Bentivogli. 2023. Test suites task: Evaluation of gender fairness in MT with MuST-SHE and INES. In *Proceedings of the Eighth Conference on Machine Translation*, pages 252–262, Singapore. Association for Computational Linguistics.

Beatrice Savoldi, Andrea Piergentili, Dennis Fucci, Matteo Negri, and Luisa Bentivogli. 2024. A prompt response to the demand for automatic gender-neutral translation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 256–267, St. Julian's, Malta. Association for Computational Linguistics.

Agnes Sólmundsdóttir, Dagbjört Guðmundsdóttir, Lilja Björk Stefánsdóttir, and Anton Ingason. 2022. Mean machine translations: On gender bias in Icelandic machine translations. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3113–3121, Marseille, France. European Language Resources Association.

Gabriel Stanovsky, Noah A. Smith, and Luke Zettlemoyer. 2019. Evaluating gender bias in machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1679–1684, Florence, Italy. Association for Computational Linguistics.

Jonas-Dario Troles and Ute Schmid. 2021. Extending challenge sets to uncover gender bias in machine translation: Impact of stereotypical verbs and adjectives. In *Proceedings of the Sixth Conference on Machine Translation*, pages 531–541, Online. Association for Computational Linguistics.

# A  Test Suite Templates

## A.1  Stereo-Adverb Templates

Examples in the Stereo-Adverb test suite subset take the form:

$$\text{"I think I'm } A\text{," I said } adverb. \qquad (1)$$

where $A$ ($n = 130$) denotes an adjective sampled from the full adjective set, and $adverb$ can be none, $M$-stereotyped ($n = 3$) or $F$-stereotyped ($n = 3$). In total, there are $N = 910$ source sentences in this subset $\big(N = 130 \times (1 + 3 + 3)\big)$.

## A.2  Stereo-Character Templates

All examples in the Stereo-Character test suite subset contain two characters that are introduced using gender-stereotyped descriptions. For simplicity, all character descriptions are in the form:

$$C_g = a_g occ_g \qquad (2)$$

where $a_g$ is gender-stereotyped adjective, and $occ_g$ is a matching gender-stereotyped occupation (e.g. "pretty nurse" or "strong doctor"). In each example, there is one female-stereotyped character ($n = 22$) and one male-stereotyped character ($n = 22$). We denote the character pairs as $(C_g, C_{\bar{g}})$.

Templates in this test suite subset come in both single-speaker and two-way conversation styles. In the single-speaker template, examples are of the form:

The $C_g$ smiled. "I think {I'm, you're} $A$,"
{he, she, they} said to the $C_{\bar{g}}$. $\qquad (3)$

where I'm+{he, she} combinations produce gender-determined referents, and you're+{he, she, they} and I'm+they combinations produce gender-ambiguous referents. There are 22 character pairs, 2 character orders, 2 referent pronoun variants, and 3 speaker pronoun variants, for a total of 264 base templates. Each base template is paired with 4 unique adjectives sampled from the full adjective set, for a total of $N = 1056 = 22 \times 2 \times 2 \times 3 \times 4$ source sentences (352 determined and 704 ambiguous).

In the two-way conversation template, examples are of the form:

The $C_g$ smiled. "I think I'm $A_1$ and you're $A_2$,"
they said.
The $C_{\bar{g}}$ laughed back. "No, you're $A_3$, but I'm $A_4$,"
they replied.

$$\qquad (4)$$

such that the gender of all adjective referents is ambiguous. To observe how the system handles repeated adjectives in the input that refer to different entities, 4 adjective equality variations are used:

$$
\begin{aligned}
&(A_1, A_2, A_3, A_4) \\
&(A_1, A_2, A_2, A_4) \\
&(A_1, A_2, A_3, A_1) \\
&(A_1, A_2, A_2, A_1).
\end{aligned}
\qquad (5)
$$

There are 22 character pairs, 2 character orders, and 4 adjective equality patterns, for a total of 176

base templates. For each base template, 5 unique adjective tuples $(A_1, A_2, A_3, A_4)$ are sampled from the full adjective set, for a total of $N = 880 = 22 \times 2 \times 4 \times 5$ source sentences. Note that each source sentence provides 4 adjective agreement samples.

### A.3 Structure Templates

The structure templates do not include any gender-stereotyped variables, and instead focus on structural variables in dialogue settings between two speakers. There are two template styles: one where all adjectives refer to the same entity, and one where both characters are referenced in equal measure. Both template styles have variations in the provided gender context: two speakers of known gender, such that each adjective's correct gender agreement is always determined, or one known gender and one unknown gender (first-person), such that the adjective's gender is either ambiguous or determined depending on the referent.

The first template style with complete gender context:

The {woman, man} smiled. "I think {I'm, you're} $A_1$," {she, he} said.
{He, She} laughed back. "No, [{you're, I'm} not $A_1$, but] {you are, I am} $A_2$," {he, she} replied. (6)

where the text contained by [...] denotes an optional chaining effect on $A_1$. There are 4 gender combinations for the two characters $\big((M, M), (F, F), (F, M), (M, F)\big)$, 2 pronoun referent variations (I, you), and 2 chaining variants (present or not), for 16 base templates. For each base template, 60 unique adjective tuples $(A_1, A_2)$ are sampled from the full adjective set, for a total of $N = 960 = 4 \times 2 \times 2 \times 2 \times 60$ source sentences.

The first template style with partial gender context:

{I, The wo/man} smiled. "I think {I'm, you're} $A_1$," {I, s/he} said.
{S/he, I} laughed back. "No, [{you're, I'm} not $A_1$, but] {you are, I am} $A_2$," {s/he, I} replied. (7)

As above, there are 4 gender combinations $\big((M, ?), (F, ?), (?, M), (?, F)\big)$, 2 pronoun referent variations, 2 chaining variations, and 60 unique adjective tuples, for a total of $N = 960$ source sentences. The structure variables split this subset in half between ambiguous and determined cases. When the unknown gender (first-person speaker, I) appears first and the first pronoun referent is "I", or the known gender speaker appears first and the first person referent is "you", all adjectives are gender-ambiguous ($n = 480$ source sentences, $n = 1200$ adjective instances). Otherwise, all adjectives are gender-determined ($n = 480$ source sentences, $n = 1200$ adjective instances). Note that each source sentence contains 2-3 adjective instances, depending on whether the optional chaining effect is included.

The second template style with complete gender context:

The {man, woman} smiled. "I think I'm $A_1$ and you're $A_2$," {he, she} said.
{He, She} laughed back. "No, you're $A_3$, but I'm $A_4$," {he, she} replied. (8)

where there are 4 possible gender combinations, 4 adjective equality patterns as described by equation (5), and 60 unique adjective tuples $(A_1, A_2, A_3, A_4)$, for a total of $N = 960$ source sentences with 4 determined adjective instances each.

The second template style with partial gender context:

{I, The wo/man} smiled. "I think I'm $A_1$ and you're $A_2$," {I, s/he} said.
{S/He, I} laughed back. "No, you're $A_3$, but I'm $A_4$," {s/he, I} replied. (9)

where again there are 4 possible gender combinations, 4 adjective equality patterns, and 60 unique adjective tuples, for a total of $N = 960$ source sentences. As with template (7), adjectives in this subset are split evenly between ambiguous and determined cases. However, unlike (7), both ambiguous and determined adjectives appear together (equally) in the same source passage. Note that all adjective positions are split evenly between determined and ambiguous cases, as determined by the variable position of the speakers.

## B Results for All Systems

315

| System | $F$ | $M$ | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|---|---|
| Aya23 | 0.30308 | 0.69692 | 0.14804 | 0.85196 | 0.53892 | 0.46108 | 0.39088 |
| Claude-3.5 | 0.08439 | 0.91561 | 0.00000 | 1.00000 | 0.14930 | 0.85070 | 0.14930 |
| CommandR-plus | 0.37172 | 0.62828 | 0.19682 | 0.80318 | 0.51624 | 0.48376 | 0.31942 |
| Dubformer | 0.08120 | 0.91880 | 0.06135 | 0.93865 | 0.10631 | 0.89369 | 0.04496 |
| GPT-4 | 0.01125 | 0.98875 | 0.00000 | 1.00000 | 0.02308 | 0.97692 | 0.02308 |
| IKUN | 0.60329 | 0.39671 | 0.49211 | 0.50789 | 0.61719 | 0.38281 | 0.12508 |
| IKUN-C | 0.49174 | 0.50826 | 0.47581 | 0.52419 | 0.45274 | 0.54726 | -0.02306 |
| IOL-Research | 0.08787 | 0.91213 | 0.03196 | 0.96804 | 0.12500 | 0.87500 | 0.09304 |
| Llama3-70B | 0.03901 | 0.96099 | 0.00000 | 1.00000 | 0.08274 | 0.91726 | 0.08274 |
| MSLC | 0.14011 | 0.85989 | 0.15891 | 0.84109 | 0.19874 | 0.80126 | 0.03983 |
| ONLINE-A | 0.09344 | 0.90656 | 0.05932 | 0.94068 | 0.11059 | 0.88941 | 0.05126 |
| ONLINE-B | 0.08717 | 0.91283 | 0.06970 | 0.93030 | 0.09690 | 0.90310 | 0.02721 |
| ONLINE-G | 0.14204 | 0.85796 | 0.14241 | 0.85759 | 0.16289 | 0.83711 | 0.02048 |
| ONLINE-W | 0.26113 | 0.73887 | 0.08274 | 0.91726 | 0.49821 | 0.50179 | 0.41548 |
| TranssionMT | 0.10359 | 0.89641 | 0.10000 | 0.90000 | 0.14516 | 0.85484 | 0.04516 |
| Unbabel-Tower70B | 0.32142 | 0.67858 | 0.17822 | 0.82178 | 0.42691 | 0.57309 | 0.24869 |

Table 6: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Adverb** test suite subset for all systems (English to **Spanish**). All adjectives self-refer to the speaker of unknown gender. In affected systems, the use of a male-stereotyped adverb to control the manner of speaking increases the use of male adjectives compared to female adjectives (see subscript $M$ denoting the use of male-stereotyped adverbs), and vice versa (see subscript $F$ denoting the use of female-stereotyped adverbs). The baseline, non-subscripted results refer to the proportions of male and female adjective use in the absence of any adverb. The overall strength of the adverb effect can be captured by the difference $\Delta M_{M-F}$.

| System | $F$ | $M$ | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|---|---|
| Aya23 | 0.36944 | 0.63056 | 0.18750 | 0.81250 | 0.61244 | 0.38756 | 0.42494 |
| CUNI-DocTransformer | 0.36044 | 0.63956 | 0.27865 | 0.72135 | 0.52107 | 0.47893 | 0.24241 |
| CUNI-GA | 0.41955 | 0.58045 | 0.35779 | 0.64221 | 0.47974 | 0.52026 | 0.12195 |
| CUNI-MH | 0.70343 | 0.29657 | 0.37886 | 0.62114 | 0.95018 | 0.04982 | 0.57132 |
| CUNI-Transformer | 0.40925 | 0.59075 | 0.37895 | 0.62105 | 0.42209 | 0.57791 | 0.04314 |
| Claude-3.5 | 0.19281 | 0.80719 | 0.00769 | 0.99231 | 0.37334 | 0.62666 | 0.36564 |
| CommandR-plus | 0.33985 | 0.66015 | 0.11950 | 0.88050 | 0.55371 | 0.44629 | 0.43421 |
| GPT-4 | 0.05730 | 0.94270 | 0.00000 | 1.00000 | 0.11644 | 0.88356 | 0.11644 |
| IKUN | 0.26889 | 0.73111 | 0.14492 | 0.85508 | 0.32364 | 0.67636 | 0.17872 |
| IKUN-C | 0.33780 | 0.66220 | 0.26528 | 0.73472 | 0.38904 | 0.61096 | 0.12376 |
| IOL-Research | 0.04607 | 0.95393 | 0.00000 | 1.00000 | 0.10601 | 0.89399 | 0.10601 |
| Llama3-70B | 0.02378 | 0.97622 | 0.00000 | 1.00000 | 0.05802 | 0.94198 | 0.05802 |
| NVIDIA-NeMo | 0.30920 | 0.69080 | 0.31792 | 0.68208 | 0.31086 | 0.68914 | -0.00706 |
| ONLINE-A | 0.97638 | 0.02362 | 0.98437 | 0.01563 | 1.00000 | 0.00000 | 0.01563 |
| ONLINE-B | 0.09241 | 0.90759 | 0.09440 | 0.90560 | 0.11212 | 0.88788 | 0.01771 |
| ONLINE-G | 0.03956 | 0.96044 | 0.03883 | 0.96117 | 0.03196 | 0.96804 | -0.00687 |
| ONLINE-W | 0.59098 | 0.40902 | 0.38679 | 0.61321 | 0.88381 | 0.11619 | 0.49703 |
| SCIR-MT | 0.21048 | 0.78952 | 0.12845 | 0.87155 | 0.33493 | 0.66507 | 0.20648 |
| TranssionMT | 0.85793 | 0.14207 | 0.78136 | 0.21864 | 0.91884 | 0.08116 | 0.13747 |
| Unbabel-Tower70B | 0.38241 | 0.61759 | 0.22141 | 0.77859 | 0.53240 | 0.46760 | 0.31100 |

Table 7: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Adverb** test suite subset for all systems (English to **Czech**).

| System | $F$ | $M$ | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|---|---|
| AMI | 0.06617 | 0.93383 | 0.04315 | 0.95685 | 0.05593 | 0.94407 | 0.01279 |
| Aya23 | 0.25779 | 0.74221 | 0.13023 | 0.86977 | 0.32917 | 0.67083 | 0.19893 |
| Claude-3.5 | 0.12794 | 0.87206 | 0.01019 | 0.98981 | 0.18941 | 0.81059 | 0.17922 |
| Dubformer | 0.13975 | 0.86025 | 0.17119 | 0.82881 | 0.13380 | 0.86620 | -0.03739 |
| GPT-4 | 0.48465 | 0.51535 | 0.36143 | 0.63857 | 0.71391 | 0.28609 | 0.35248 |
| IKUN | 0.74118 | 0.25882 | 0.61495 | 0.38505 | 0.81479 | 0.18521 | 0.19984 |
| IKUN-C | 0.42608 | 0.57392 | 0.35849 | 0.64151 | 0.46387 | 0.53613 | 0.10538 |
| IOL-Research | 0.19206 | 0.80794 | 0.13012 | 0.86988 | 0.25000 | 0.75000 | 0.11988 |
| Llama3-70B | 0.17153 | 0.82847 | 0.07006 | 0.92994 | 0.24176 | 0.75824 | 0.17170 |
| ONLINE-A | 0.09271 | 0.90729 | 0.09124 | 0.90876 | 0.09357 | 0.90643 | 0.00234 |
| ONLINE-B | 0.20944 | 0.79056 | 0.23913 | 0.76087 | 0.18975 | 0.81025 | -0.04938 |
| ONLINE-G | 0.14517 | 0.85483 | 0.15385 | 0.84615 | 0.16008 | 0.83992 | 0.00623 |
| TranssionMT | 0.23640 | 0.76360 | 0.26056 | 0.73944 | 0.25143 | 0.74857 | -0.00913 |
| Unbabel-Tower70B | 0.26414 | 0.73586 | 0.11700 | 0.88300 | 0.44016 | 0.55984 | 0.32316 |

Table 8: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Adverb** test suite subset for all systems (English to **Icelandic**).

| Variable | ONLINE-W | Aya23 | CommandR-plus |
|---|---|---|---|
| Intercept | $(-1.73, 3.7E-07^{***})$ | $(-1.07, 1.1E-04^{***})$ | $(-0.71, 7.6E-03^{**})$ |
| Adj Stereo(M) | $(-0.89, 1.6E-02^{*})$ | $(-0.15, 6.5E-01)$ | $(-0.82, 7.8E-03^{**})$ |
| Adj Stereo(F) | $(1.96, 4.8E-15^{***})$ | $(0.48, 2.1E-02^{*})$ | $(0.50, 1.2E-02^{*})$ |
| Adj Sentiment(neg) | $(-0.44, 3.6E-02^{*})$ | $(0.35, 5.6E-02)$ | $(0.29, 9.6E-02)$ |
| Adj Type(appearance) | $(0.63, 3.1E-03^{**})$ | $(0.31, 9.5E-02)$ | $(0.52, 3.5E-03^{**})$ |
| **Adv Stereo(M)** | $(-0.42, 2.0E-01)$ | $(-0.99, 4.4E-04^{***})$ | $(-0.87, 8.5E-04^{***})$ |
| **Adv Stereo(F)** | $(1.61, 4.8E-07^{***})$ | $(0.65, 1.0E-02^{*})$ | $(0.54, 3.0E-02^{*})$ |

Table 9: Stereotypical manner of speaking (adverb) regression analysis for the most affected systems (**Spanish**), displayed as (coefficient value, p-value). The variable to predict is the binary adjective declension choice, where feminine adjectives are the positive class, such that negative coefficient values indicate a greater probability of $M$, and positive coefficient values indicate a greater probability of $F$. Strong negative intercepts indicate the default male baseline exhibited by many systems. Internal adjective traits are controlled by stereotype variables (e.g. Stereo(M) is expected to increase the probability of $M$), the sentiment (here negative as opposed to positive), and type (here appearance as opposed to character). For example, if the adjective is the appearance type, the results show that systems ONLINE-W and CommandR-plus are more likely to choose an $F$-adjective, controlling for all other variables. Here we see the adverb variables are strong in their expected directions, and significant.

| Variable | CUNI-MH | ONLINE-W | CommandR-plus |
|---|---|---|---|
| Intercept | $(2.45, 1.2E-07^{***})$ | $(-0.66, 1.1E-02^{*})$ | $(-0.84, 1.3E-03^{**})$ |
| Adj Stereo(M) | $(-0.96, 8.1E-03^{**})$ | $(-0.76, 1.5E-02^{*})$ | $(-0.60, 6.3E-02)$ |
| Adj Stereo(F) | $(1.04, 1.6E-04^{***})$ | $(0.95, 2.7E-05^{***})$ | $(1.30, 5.1E-09^{***})$ |
| Adj Sentiment(neg) | $(0.13, 5.7E-01)$ | $(-0.47, 1.2E-02^{*})$ | $(-0.19, 3.2E-01)$ |
| Adj Type(appearance) | $(0.57, 1.4E-02^{*})$ | $(1.04, 8.9E-08^{***})$ | $(0.94, 1.1E-06^{***})$ |
| **Adv Stereo(M)** | $(-3.26, 7.0E-13^{***})$ | $(0.01, 9.6E-01)$ | $(-1.59, 1.2E-08^{***})$ |
| **Adv Stereo(F)** | $(0.05, 9.3E-01)$ | $(2.26, 5.3E-16^{***})$ | $(0.68, 6.3E-03^{**})$ |

Table 10: Stereotypical manner of speaking (adverb) regression analysis for the most affected systems (**Czech**).

| Variable | GPT-4 | Unbabel-Tower70B | IKUN |
|---|---|---|---|
| Intercept | $(0.52, 4.7E-02^*)$ | $(-0.91, 1.4E-03^{**})$ | $(1.31, 1.6E-05^{***})$ |
| Adj Stereo(M) | $(-1.51, 4.7E-07^{***})$ | $(-0.72, 2.9E-02^*)$ | $(-0.78, 3.6E-03^{**})$ |
| Adj Stereo(F) | $(0.75, 2.6E-04^{***})$ | $(0.89, 1.4E-05^{***})$ | $(0.62, 1.1E-02^*)$ |
| Adj Sentiment(neg) | $(-0.22, 2.1E-01)$ | $(-0.49, 9.7E-03^{**})$ | $(-0.01, 9.6E-01)$ |
| Adj Type(appearance) | $(0.28, 1.1E-01)$ | $(0.03, 8.7E-01)$ | $(-0.08, 6.5E-01)$ |
| **Adv Stereo(M)** | $(-0.98, 1.5E-04^{***})$ | $(-0.54, 5.7E-02)$ | $(-0.64, 2.6E-02^*)$ |
| **Adv Stereo(F)** | $(0.31, 2.3E-01)$ | $(0.70, 9.0E-03^{**})$ | $(0.07, 8.2E-01)$ |

Table 11: Stereotypical manner of speaking (adverb) regression analysis for the most affected systems (**Icelandic**).

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|
| Aya23 | 0.122065 | 0.877935 | 0.429200 | 0.570800 | 0.307135 |
| Claude-3.5 | 0.000000 | 1.000000 | 0.390775 | 0.609225 | 0.390775 |
| CommandR-plus | 0.011598 | 0.988402 | 0.401382 | 0.598618 | 0.389784 |
| Dubformer | 0.002979 | 0.997021 | 0.128387 | 0.871613 | 0.125408 |
| GPT-4 | 0.000000 | 1.000000 | 0.273960 | 0.726040 | 0.273960 |
| IKUN | 0.165404 | 0.834596 | 0.363451 | 0.636549 | 0.198047 |
| IKUN-C | 0.287342 | 0.712658 | 0.418834 | 0.581166 | 0.131493 |
| IOL-Research | 0.003182 | 0.996818 | 0.245182 | 0.754818 | 0.242001 |
| Llama3-70B | 0.000000 | 1.000000 | 0.261807 | 0.738193 | 0.261807 |
| MSLC | 0.263380 | 0.736620 | 0.217680 | 0.782320 | -0.045700 |
| ONLINE-A | 0.023395 | 0.976605 | 0.109371 | 0.890629 | 0.085976 |
| ONLINE-B | 0.038452 | 0.961548 | 0.084521 | 0.915479 | 0.046069 |
| ONLINE-G | 0.054620 | 0.945380 | 0.042830 | 0.957170 | -0.011790 |
| ONLINE-W | 0.055130 | 0.944870 | 0.175764 | 0.824236 | 0.120634 |
| TranssionMT | 0.020904 | 0.979096 | 0.100892 | 0.899108 | 0.079988 |
| Unbabel-Tower70B | 0.057582 | 0.942418 | 0.359154 | 0.640846 | 0.301572 |

Table 12: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Character-Amb** test suite subset for unknown gender cases for all systems (English to **Spanish**). All adjectives refer to either the speaker or the listener which have been introduced as gender-stereotyped characters. The subscripts denote the gender stereotype label. In affected systems, adjectives that refer to a $M$-stereotyped character (subscript $M$) are more likely to be translated with a male declension, and vice versa for $F$-stereotyped characters. The overall strength of the character description effect can be captured by the difference $\Delta M_{M-F}$.

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|
| Aya23 | 0.157011 | 0.842989 | 0.490925 | 0.509075 | 0.333913 |
| CUNI-DocTransformer | 0.194503 | 0.805497 | 0.304000 | 0.696000 | 0.109497 |
| CUNI-GA | 0.237413 | 0.762587 | 0.328323 | 0.671677 | 0.090910 |
| CUNI-MH | 0.184362 | 0.815638 | 0.649108 | 0.350892 | 0.464746 |
| CUNI-Transformer | 0.279416 | 0.720584 | 0.329398 | 0.670602 | 0.049982 |
| Claude-3.5 | 0.020730 | 0.979270 | 0.424670 | 0.575330 | 0.403940 |
| CommandR-plus | 0.020468 | 0.979532 | 0.354000 | 0.646000 | 0.333532 |
| GPT-4 | 0.005208 | 0.994792 | 0.348186 | 0.651814 | 0.342978 |
| IKUN | 0.125093 | 0.874907 | 0.402927 | 0.597073 | 0.277834 |
| IKUN-C | 0.273129 | 0.726871 | 0.528897 | 0.471103 | 0.255767 |
| IOL-Research | 0.037213 | 0.962787 | 0.350441 | 0.649559 | 0.313228 |
| Llama3-70B | 0.002959 | 0.997041 | 0.193021 | 0.806979 | 0.190062 |
| NVIDIA-NeMo | 0.127240 | 0.872760 | 0.261404 | 0.738596 | 0.134165 |
| ONLINE-A | 0.140739 | 0.859261 | 0.223089 | 0.776911 | 0.082350 |
| ONLINE-B | 0.050206 | 0.949794 | 0.097500 | 0.902500 | 0.047294 |
| ONLINE-G | 0.063216 | 0.936784 | 0.070601 | 0.929399 | 0.007385 |
| ONLINE-W | 0.085828 | 0.914172 | 0.458372 | 0.541628 | 0.372544 |
| SCIR-MT | 0.082652 | 0.917348 | 0.277897 | 0.722103 | 0.195246 |
| TranssionMT | 0.108245 | 0.891755 | 0.178112 | 0.821888 | 0.069868 |
| Unbabel-Tower70B | 0.048333 | 0.951667 | 0.381000 | 0.619000 | 0.332667 |

Table 13: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Character-Amb** test suite subset for unknown gender cases for all systems (English to **Czech**).

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|
| AMI | 0.106313 | 0.893687 | 0.077960 | 0.922040 | -0.028353 |
| Aya23 | 0.234091 | 0.765909 | 0.450333 | 0.549667 | 0.216242 |
| Claude-3.5 | 0.005272 | 0.994728 | 0.434516 | 0.565484 | 0.429244 |
| Dubformer | 0.092273 | 0.907727 | 0.139669 | 0.860331 | 0.047396 |
| GPT-4 | 0.159235 | 0.840765 | 0.477160 | 0.522840 | 0.317924 |
| IKUN | 0.282857 | 0.717143 | 0.555590 | 0.444410 | 0.272733 |
| IKUN-C | 0.273029 | 0.726971 | 0.378756 | 0.621244 | 0.105728 |
| IOL-Research | 0.002394 | 0.997606 | 0.126020 | 0.873980 | 0.123627 |
| Llama3-70B | 0.065437 | 0.934563 | 0.235696 | 0.764304 | 0.170259 |
| ONLINE-A | 0.040660 | 0.959340 | 0.031744 | 0.968256 | -0.008916 |
| ONLINE-B | 0.122783 | 0.877217 | 0.087892 | 0.912108 | -0.034891 |
| ONLINE-G | 0.089511 | 0.910489 | 0.048492 | 0.951508 | -0.041019 |
| TranssionMT | 0.103972 | 0.896028 | 0.100703 | 0.899297 | -0.003269 |
| Unbabel-Tower70B | 0.063657 | 0.936343 | 0.293988 | 0.706012 | 0.230331 |

Table 14: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Stereo-Character-Amb** test suite subset for unknown gender cases for all systems (English to **Icelandic**).

| Variable | Claude-3.5 | CommandR-plus | Aya23 |
|---|---|---|---|
| Intercept | $(-6.83, 2.7E-11^{***})$ | $(-3.89, 5.5E-34^{***})$ | $(-0.99, 3.1E-07^{***})$ |
| Adj Stereo(M) | $(0.87, 8.2E-02)$ | $(0.29, 5.4E-01)$ | $(-0.13, 7.9E-01)$ |
| Adj Stereo(F) | $(0.85, 2.7E-06^{***})$ | $(-0.29, 1.3E-01)$ | $(0.15, 2.5E-01)$ |
| Adj Sentiment(neg) | $(-0.56, 1.6E-04^{***})$ | $(0.04, 7.4E-01)$ | $(-0.86, 4.9E-16^{***})$ |
| Adj Type(appearance) | $(0.01, 9.7E-01)$ | $(0.67, 3.6E-04^{***})$ | $(0.06, 7.0E-01)$ |
| **Character Stereo(F)** | $(6.59, 5.0E-11^{***})$ | $(3.56, 2.4E-43^{***})$ | $(1.41, 9.5E-43^{***})$ |

Table 15: Stereotypical character description regression analysis for the most affected systems (**Spanish**). Internal adjective variables are defined as above (see Table 9). Here **Character Stereo(F)** denotes a binary variable equal to 1 when the character description is stereotypically female, and equal to 0 when the character description is stereotypically male. As shown, coefficient values for **Character Stereo(F)** are significant and in the expected direction (positive, indicating an increased likelihood of a $F$-adjective), and much stronger than the internal variables.

| Variable | CUNI-MH | Claude-3.5 | ONLINE-W |
|---|---|---|---|
| Intercept | $(-1.38, 1.5E-16^{***})$ | $(-4.93, 1.2E-50^{***})$ | $(-1.36, 9.4E-13^{***})$ |
| Adj Stereo(M) | $(-0.53, 3.6E-03^{**})$ | $(-0.75, 8.6E-03^{**})$ | $(-1.01, 6.7E-02)$ |
| Adj Stereo(F) | $(1.58, 1.7E-37^{***})$ | $(0.17, 3.0E-01)$ | $(0.01, 9.5E-01)$ |
| Adj Sentiment(neg) | $(-0.43, 9.8E-08^{***})$ | $(-0.31, 6.9E-03^{**})$ | $(-0.64, 2.1E-11^{***})$ |
| Adj Type(appearance) | $(-0.78, 5.6E-09^{***})$ | $(0.20, 3.1E-01)$ | $(1.34, 2.5E-26^{***})$ |
| **Character Stereo(F)** | $(1.66, 1.1E-93^{***})$ | $(4.25, 4.4E-52^{***})$ | $(2.20, 1.1E-76^{***})$ |

Table 16: Stereotypical character description regression analysis for the most affected systems (**Czech**).

| Variable | Claude-3.5 | GPT-4 | IKUN |
|---|---|---|---|
| Intercept | $(-4.09, 9.7E-45^{***})$ | $(-1.93, 3.0E-31^{***})$ | $(-1.13, 5.6E-14^{***})$ |
| Adj Stereo(M) | $(-1.55, 1.3E-06^{***})$ | $(-0.18, 3.3E-01)$ | $(0.12, 6.7E-01)$ |
| Adj Stereo(F) | $(-0.45, 1.5E-02^{*})$ | $(0.04, 7.9E-01)$ | $(0.48, 2.1E-04^{***})$ |
| Adj Sentiment(neg) | $(0.91, 2.1E-15^{***})$ | $(0.34, 6.9E-05^{***})$ | $(0.80, 2.3E-24^{***})$ |
| Adj Type(appearance) | $(0.03, 8.8E-01)$ | $(-0.03, 8.3E-01)$ | $(0.69, 1.8E-09^{***})$ |
| Character Stereo(F) | $(3.62, 1.3E-64^{***})$ | $(1.16, 1.5E-44^{***})$ | $(0.60, 4.7E-15^{***})$ |

Table 17: Stereotypical character description regression analysis for the most affected systems (**Icelandic**).

| System | Accuracy (PRO) | Accuracy (ANTI) | Δ(PRO, ANTI) |
|---|---|---|---|
| Aya23 | 1.000 | 0.655000 | 0.345000 |
| Claude-3.5 | 1.000 | 0.742500 | 0.257500 |
| CommandR-plus | 1.000 | 0.950167 | 0.049833 |
| Dubformer | 0.965 | 0.776000 | 0.189000 |
| GPT-4 | 0.990 | 0.527500 | 0.462500 |
| IKUN | 0.975 | 0.701667 | 0.273333 |
| IKUN-C | 0.910 | 0.887167 | 0.022833 |
| IOL-Research | 0.990 | 0.963500 | 0.026500 |
| Llama3-70B | 1.000 | 0.862500 | 0.137500 |
| MSLC | 0.935 | 0.870500 | 0.064500 |
| ONLINE-A | 0.970 | 0.990667 | -0.020667 |
| ONLINE-B | 0.980 | 0.996000 | -0.016000 |
| ONLINE-G | 1.000 | 0.996333 | 0.003667 |
| ONLINE-W | 0.985 | 0.414333 | 0.570667 |
| TranssionMT | 0.985 | 0.990167 | -0.005167 |
| Unbabel-Tower70B | 0.995 | 0.970833 | 0.024167 |

Table 18: The accuracy in gender-adjective agreement on the **Stereo-Character-Det** test suite subset for known gender cases for all systems (English to **Spanish**). The test subset is further partitioned into cases that align with a stereotype (PRO) and cases that oppose a stereotype (ANTI). Accuracy is consistently high on the PRO subset, and drops significantly in the challenge setting for some translation systems, indicating that stereotype effects persist in the presence of correct and unambiguous gender context.

| System | Accuracy (PRO) | Accuracy (ANTI) | Δ(PRO, ANTI) |
|---|---|---|---|
| Aya23 | 0.985 | 0.539833 | 0.445167 |
| CUNI-DocTransformer | 1.000 | 0.993500 | 0.006500 |
| CUNI-GA | 0.995 | 0.991667 | 0.003333 |
| CUNI-MH | 1.000 | 0.995500 | 0.004500 |
| CUNI-Transformer | 1.000 | 0.994167 | 0.005833 |
| Claude-3.5 | 1.000 | 0.878167 | 0.121833 |
| CommandR-plus | 0.985 | 0.912333 | 0.072667 |
| GPT-4 | 1.000 | 0.807833 | 0.192167 |
| IKUN | 0.935 | 0.814000 | 0.121000 |
| IKUN-C | 0.995 | 0.930000 | 0.065000 |
| IOL-Research | 1.000 | 0.878000 | 0.122000 |
| Llama3-70B | 1.000 | 0.640000 | 0.360000 |
| NVIDIA-NeMo | 1.000 | 0.993333 | 0.006667 |
| ONLINE-A | 1.000 | 0.996333 | 0.003667 |
| ONLINE-B | 1.000 | 0.996667 | 0.003333 |
| ONLINE-G | 1.000 | 1.000000 | 0.000000 |
| ONLINE-W | 1.000 | 0.988667 | 0.011333 |
| SCIR-MT | 1.000 | 0.871667 | 0.128333 |
| TranssionMT | 1.000 | 0.994667 | 0.005333 |
| Unbabel-Tower70B | 1.000 | 0.867000 | 0.133000 |

Table 19: The accuracy in gender-adjective agreement on the **Stereo-Character-Det** test suite subset for known gender cases for all systems (English to **Czech**).

| System | Accuracy (PRO) | Accuracy (ANTI) | $\Delta$(PRO, ANTI) |
|---|---|---|---|
| AMI | 0.990000 | 0.977667 | 0.012333 |
| Aya23 | 0.632833 | 0.765833 | -0.133000 |
| Claude-3.5 | 0.990000 | 0.900833 | 0.089167 |
| Dubformer | 0.657560 | 0.564000 | 0.093560 |
| GPT-4 | 0.925000 | 0.832500 | 0.092500 |
| IKUN | 0.890000 | 0.942167 | -0.052167 |
| IKUN-C | 0.975000 | 0.950333 | 0.024667 |
| IOL-Research | 0.955000 | 0.963167 | -0.008167 |
| Llama3-70B | 0.982416 | 0.850000 | 0.132416 |
| ONLINE-A | 0.885000 | 0.973833 | -0.088833 |
| ONLINE-B | 1.000000 | 0.988500 | 0.011500 |
| ONLINE-G | 0.940000 | 0.830000 | 0.110000 |
| TranssionMT | 1.000000 | 0.982833 | 0.017167 |
| Unbabel-Tower70B | 1.000000 | 0.969167 | 0.030833 |

Table 20: The accuracy in gender-adjective agreement on the **Stereo-Character-Det** test suite subset for known gender cases for all systems (English to **Icelandic**).

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|
| Aya23 | 0.256023 | 0.743977 | 0.345843 | 0.654157 | 0.089820 |
| Claude-3.5 | 0.419257 | 0.580743 | 0.073611 | 0.926389 | -0.345646 |
| CommandR-plus | 0.763592 | 0.236408 | 0.425954 | 0.574046 | -0.337638 |
| Dubformer | 0.089736 | 0.910264 | 0.110243 | 0.889757 | 0.020507 |
| GPT-4 | 0.090745 | 0.909255 | 0.019064 | 0.980936 | -0.071682 |
| IKUN | 0.255602 | 0.744398 | 0.725592 | 0.274408 | 0.469990 |
| IKUN-C | 0.292295 | 0.707705 | 0.702554 | 0.297446 | 0.410258 |
| IOL-Research | 0.070826 | 0.929174 | 0.342207 | 0.657793 | 0.271382 |
| Llama3-70B | 0.105211 | 0.894789 | 0.064873 | 0.935127 | -0.040338 |
| MSLC | 0.177329 | 0.822671 | 0.247036 | 0.752964 | 0.069707 |
| ONLINE-A | 0.049758 | 0.950242 | 0.208951 | 0.791049 | 0.159194 |
| ONLINE-B | 0.105807 | 0.894193 | 0.141280 | 0.858720 | 0.035473 |
| ONLINE-G | 0.070760 | 0.929240 | 0.300656 | 0.699344 | 0.229895 |
| ONLINE-W | 0.261696 | 0.738304 | 0.372914 | 0.627086 | 0.111218 |
| TranssionMT | 0.091749 | 0.908251 | 0.151915 | 0.848085 | 0.060166 |
| Unbabel-Tower70B | 0.273571 | 0.726429 | 0.415857 | 0.584143 | 0.142286 |

Table 21: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Structure-Amb** test suite subset for all systems (English to **Spanish**). All adjectives refer to someone of an unknown gender in conversation with someone of a known gender (where that known gender is denoted by the subscripts). Systems that have an "opposite" binary gender bias effect resolve the ambiguous-gender speaker to be opposite to the known speaker (i.e., masculine adjectives increase when the other speaker is female, and vice versa, and the difference $\Delta M_{M-F}$ is strongly positive). Systems with a same-binary gender effect consistently choose adjective forms matching the gender of the other speaker (i.e., $\Delta M_{M-F}$ is strongly negative).

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|
| Aya23 | 0.671295 | 0.328705 | 0.583859 | 0.416141 | -0.087436 |
| CUNI-DocTransformer | 0.092894 | 0.907106 | 0.270689 | 0.729311 | 0.177795 |
| CUNI-GA | 0.455846 | 0.544154 | 0.175642 | 0.824358 | -0.280204 |
| CUNI-MH | 0.714519 | 0.285481 | 0.699626 | 0.300374 | -0.014893 |
| CUNI-Transformer | 0.454953 | 0.545047 | 0.173387 | 0.826613 | -0.281566 |
| Claude-3.5 | 0.380231 | 0.619769 | 0.039362 | 0.960638 | -0.340869 |
| CommandR-plus | 0.661912 | 0.338088 | 0.139604 | 0.860396 | -0.522308 |
| GPT-4 | 0.496953 | 0.503047 | 0.033441 | 0.966559 | -0.463512 |
| IKUN | 0.118729 | 0.881271 | 0.322997 | 0.677003 | 0.204268 |
| IKUN-C | 0.475241 | 0.524759 | 0.808473 | 0.191527 | 0.333232 |
| IOL-Research | 0.131136 | 0.868864 | 0.053251 | 0.946749 | -0.077885 |
| Llama3-70B | 0.041654 | 0.958346 | 0.014288 | 0.985712 | -0.027366 |
| NVIDIA-NeMo | 0.027024 | 0.972976 | 0.635958 | 0.364042 | 0.608934 |
| ONLINE-A | 0.733600 | 0.266400 | 0.381812 | 0.618188 | -0.351788 |
| ONLINE-B | 0.058142 | 0.941858 | 0.082191 | 0.917809 | 0.024049 |
| ONLINE-G | 0.025188 | 0.974812 | 0.192564 | 0.807436 | 0.167376 |
| ONLINE-W | 0.645775 | 0.354225 | 0.390405 | 0.609595 | -0.255370 |
| SCIR-MT | 0.361420 | 0.638580 | 0.440349 | 0.559651 | 0.078929 |
| TranssionMT | 0.586897 | 0.413103 | 0.326415 | 0.673585 | -0.260481 |
| Unbabel-Tower70B | 0.443972 | 0.556028 | 0.323738 | 0.676262 | -0.120234 |

Table 22: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Structure-Amb** test suite subset for all systems (English to **Czech**).

| System | $F_M$ | $M_M$ | $F_F$ | $M_F$ | $\Delta M_{M-F}$ |
|---|---|---|---|---|---|
| AMI | 0.078522 | 0.921478 | 0.345131 | 0.654869 | 0.266609 |
| Aya23 | 0.298890 | 0.701110 | 0.269610 | 0.730390 | -0.029280 |
| Claude-3.5 | 0.561121 | 0.438879 | 0.151645 | 0.848355 | -0.409476 |
| Dubformer | 0.087548 | 0.912452 | 0.158890 | 0.841110 | 0.071343 |
| GPT-4 | 0.683264 | 0.316736 | 0.517260 | 0.482740 | -0.166004 |
| IKUN | 0.497869 | 0.502131 | 0.862015 | 0.137985 | 0.364145 |
| IKUN-C | 0.302331 | 0.697669 | 0.695231 | 0.304769 | 0.392900 |
| IOL-Research | 0.085036 | 0.914964 | 0.236664 | 0.763336 | 0.151628 |
| Llama3-70B | 0.248973 | 0.751027 | 0.333898 | 0.666102 | 0.084925 |
| ONLINE-A | 0.035666 | 0.964334 | 0.139427 | 0.860573 | 0.103761 |
| ONLINE-B | 0.140336 | 0.859664 | 0.256017 | 0.743983 | 0.115682 |
| ONLINE-G | 0.069463 | 0.930537 | 0.136131 | 0.863869 | 0.066667 |
| TranssionMT | 0.142135 | 0.857865 | 0.253417 | 0.746583 | 0.111282 |
| Unbabel-Tower70B | 0.207583 | 0.792417 | 0.407592 | 0.592408 | 0.200009 |

Table 23: The proportion of adjectives with male ($M$) and female ($F$) agreement on the **Structure-Amb** test suite subset for all systems (English to **Icelandic**).

| System | Acc (one gender) | Acc (same genders) | Acc (opp genders) | $\Delta$(same, opp) |
|---|---|---|---|---|
| Aya23 | 0.937340 | 0.812066 | 0.930527 | -0.118461 |
| Claude-3.5 | 0.997078 | 0.923965 | 0.997372 | -0.073407 |
| CommandR-plus | 0.987414 | 0.796548 | 0.990717 | -0.194170 |
| Dubformer | 0.844990 | 0.790325 | 0.850310 | -0.059985 |
| GPT-4 | 0.991524 | 0.855742 | 0.992963 | -0.137221 |
| IKUN | 0.876698 | 0.835986 | 0.837003 | -0.001018 |
| IKUN-C | 0.863909 | 0.838490 | 0.798583 | 0.039907 |
| IOL-Research | 0.947063 | 0.873722 | 0.906976 | -0.033254 |
| Llama3-70B | 0.956589 | 0.805900 | 0.977354 | -0.171454 |
| MSLC | 0.611581 | 0.692553 | 0.598783 | 0.093771 |
| ONLINE-A | 0.734181 | 0.828018 | 0.667730 | 0.160288 |
| ONLINE-B | 0.727604 | 0.740764 | 0.746103 | -0.005339 |
| ONLINE-G | 0.725552 | 0.826803 | 0.624745 | 0.202058 |
| ONLINE-W | 0.914281 | 0.887881 | 0.919022 | -0.031141 |
| TranssionMT | 0.728791 | 0.739865 | 0.748009 | -0.008144 |
| Unbabel-Tower70B | 0.924064 | 0.817639 | 0.909270 | -0.091631 |

Table 24: The accuracy in gender-adjective agreement on the **Structure-Det** test suite subset for known gender cases for all systems (English to **Spanish**). The second speaker in the conversation is either unknown (one gender subset), the same, or opposite to the adjective referent. Systems with an opposite binary gender effect suffer on the same-gender subset such that the difference in accuracy $\Delta$(same, opp) $\ll 0$, and systems with a same-gender preference suffer on the opposite-gender subset such that the difference in accuracy $\Delta$(same, opp) $\gg 0$.

| System | Acc (one gender) | Acc (same genders) | Acc (opp genders) | $\Delta$(same, opp) |
|---|---|---|---|---|
| Aya23 | 0.965847 | 0.808469 | 0.951880 | -0.143411 |
| CUNI-DocTransformer | 0.892850 | 0.896380 | 0.855471 | 0.040909 |
| CUNI-GA | 0.768732 | 0.601509 | 0.911337 | -0.309828 |
| CUNI-MH | 0.928232 | 0.814084 | 0.898555 | -0.084471 |
| CUNI-Transformer | 0.769805 | 0.602070 | 0.911349 | -0.309278 |
| Claude-3.5 | 0.995241 | 0.911843 | 0.999082 | -0.087239 |
| CommandR-plus | 0.996232 | 0.739137 | 0.990259 | -0.251121 |
| GPT-4 | 0.997750 | 0.823093 | 0.989451 | -0.166358 |
| IKUN | 0.856785 | 0.812804 | 0.878981 | -0.066177 |
| IKUN-C | 0.883146 | 0.867369 | 0.831151 | 0.036219 |
| IOL-Research | 0.975905 | 0.916801 | 0.957697 | -0.040897 |
| Llama3-70B | 0.953773 | 0.827831 | 0.931734 | -0.103903 |
| NVIDIA-NeMo | 0.824710 | 0.797210 | 0.704606 | 0.092604 |
| ONLINE-A | 0.736260 | 0.561771 | 0.926380 | -0.364608 |
| ONLINE-B | 0.760273 | 0.752525 | 0.750789 | 0.001737 |
| ONLINE-G | 0.739938 | 0.812760 | 0.667776 | 0.144984 |
| ONLINE-W | 0.892118 | 0.828515 | 0.927573 | -0.099059 |
| SCIR-MT | 0.916721 | 0.835374 | 0.869839 | -0.034465 |
| TranssionMT | 0.742036 | 0.597267 | 0.894969 | -0.297702 |
| Unbabel-Tower70B | 0.935530 | 0.863328 | 0.926457 | -0.063129 |

Table 25: The accuracy in gender-adjective agreement on the **Structure-Det** test suite subset for known gender cases for all systems (English to **Czech**).

| System | Acc (one gender) | Acc (same genders) | Acc (opp genders) | $\Delta$(same, opp) |
|---|---|---|---|---|
| AMI | 0.741426 | 0.890035 | 0.606934 | 0.283102 |
| Aya23 | 0.650749 | 0.665003 | 0.681895 | -0.016892 |
| Claude-3.5 | 0.990550 | 0.948800 | 0.983105 | -0.034305 |
| Dubformer | 0.685313 | 0.663091 | 0.701806 | -0.038716 |
| GPT-4 | 0.923000 | 0.862593 | 0.906645 | -0.044052 |
| IKUN | 0.859620 | 0.793285 | 0.788005 | 0.005280 |
| IKUN-C | 0.860037 | 0.826175 | 0.774361 | 0.051815 |
| IOL-Research | 0.927880 | 0.879636 | 0.890744 | -0.011107 |
| Llama3-70B | 0.863632 | 0.784824 | 0.830711 | -0.045887 |
| ONLINE-A | 0.681548 | 0.743801 | 0.602983 | 0.140818 |
| ONLINE-B | 0.745792 | 0.794195 | 0.697057 | 0.097137 |
| ONLINE-G | 0.579342 | 0.617988 | 0.546436 | 0.071552 |
| TranssionMT | 0.747217 | 0.795348 | 0.691245 | 0.104103 |
| Unbabel-Tower70B | 0.933936 | 0.892148 | 0.916886 | -0.024738 |

Table 26: The accuracy in gender-adjective agreement on the **Structure-Det** test suite subset for known gender cases for all systems (English to **Icelandic**).

| Variable | CommandR-plus | Llama3-70B | GPT-4 |
|---|---|---|---|
| Intercept | $(5.31, 1.3E-24^{***})$ | $(3.69, 8.8E-41^{***})$ | $(6.88, 1.1E-11^{***})$ |
| True(M) | $(-12.55, 9.4E-96^{***})$ | $(-11.53, 8.9E-129^{***})$ | $(-14.32, 2.9E-40^{***})$ |
| Adj Stereo(M) | $(-0.03, 9.0E-01)$ | $(-0.58, 5.8E-03^{**})$ | $(-0.95, 6.5E-05^{***})$ |
| Adj Stereo(F) | $(1.26, 2.7E-17^{***})$ | $(1.77, 5.6E-24^{***})$ | $(1.89, 7.5E-24^{***})$ |
| Adj Sentiment(neg) | $(-0.66, 4.6E-07^{***})$ | $(-0.67, 1.1E-06^{***})$ | $(-1.05, 1.8E-11^{***})$ |
| Adj Type(appearance) | $(-0.32, 2.7E-02^{*})$ | $(0.16, 2.8E-01)$ | $(-0.10, 5.3E-01)$ |
| You(M) | $(2.12, 2.0E-13^{***})$ | $(1.62, 2.0E-07^{***})$ | $(0.95, 8.2E-03^{**})$ |
| You(F) | $(-4.06, 4.4E-14^{***})$ | $(-3.27, 4.3E-28^{***})$ | $(-6.02, 3.3E-09^{***})$ |
| Lookahead(M) | $(1.83, 4.4E-13^{***})$ | $(0.26, 3.7E-01)$ | $(1.47, 1.8E-04^{***})$ |
| Lookahead(F) | $(-2.48, 3.2E-21^{***})$ | $(-2.68, 1.3E-22^{***})$ | $(-1.69, 2.8E-11^{***})$ |
| Consistency(M) | $(-0.10, 6.7E-01)$ | $(-2.23, 1.1E-15^{***})$ | $(0.18, 6.2E-01)$ |
| Consistency(F) | $(0.25, 3.0E-01)$ | $(0.64, 1.9E-03^{**})$ | $(0.67, 3.7E-03^{**})$ |
| Opposite(M) | $(3.56, 7.1E-59^{***})$ | $(4.31, 6.0E-50^{***})$ | $(3.28, 2.1E-40^{***})$ |

Table 27: Structural factors regression analysis for the systems with the greatest opposite binary gender tendency (**Spanish**). As above (refer to Table 9), the variable to predict in the adjective declension choice, where female is the positive class. Unlike the prior regression results, here we include determined-gender cases in order to assess the difficulty introduced by different structural factors. Therefore, the true gender of the referent must be controlled for (**True(M)** is 1 when the true label is $M$, 0 when the true label is $F$). In addition to the regular adjective traits, we include structural factors consistency(M/F): 1 if an earlier adjective refers to the same entity and is M/F, lookahead(M/F): 1 if an adjective's referent appears for the first time after the adjective and the known gender is M/F, and you(M/F): 1 if the adjective refers to "you" and the known gender is M/F. Lookahead and you variables must be paired with the true label because they affect the task difficulty regardless of gender. The results show that both lookahead and you strongly increase difficulty (as indicated by strong, significant, and positive coefficients when the correct label is $M$, and conversely strong, significant, and negative coefficients when the correct label is $F$). That is, the coefficients indicate an increased likelihood of choosing the incorrect gender agreement, while all else is controlled for. The variable of interest in these systems is the **Opposite(M)**: 1 when the other referent in conversation is known to be male. Systems with a strong opposite binary gender effect have strong positive coefficients, indicating an increased likelihood of a $F$-adjective.

| Variable | ONLINE-A | CUNI-GA | TranssionMT |
|---|---|---|---|
| Intercept | $(4.82, 4.9E-30^{***})$ | $(18.93, 9.5E-01)$ | $(4.57, 4.5E-35^{***})$ |
| True(M) | $(-12.42, 3.3E-145^{***})$ | $(-27.04, 9.2E-01)$ | $(-10.64, 1.3E-147^{***})$ |
| Adj Stereo(M) | $(-0.01, 9.6E-01)$ | $(-0.30, 1.2E-01)$ | $(-0.24, 1.7E-01)$ |
| Adj Stereo(F) | $(0.99, 8.6E-15^{***})$ | $(0.89, 1.3E-12^{***})$ | $(0.83, 2.8E-13^{***})$ |
| Adj Sentiment(neg) | $(-0.49, 1.5E-05^{***})$ | $(-0.61, 2.0E-08^{***})$ | $(-0.44, 1.1E-05^{***})$ |
| Adj Type(appearance) | $(0.16, 2.3E-01)$ | $(0.58, 3.2E-06^{***})$ | $(0.45, 1.3E-04^{***})$ |
| You(M) | $(2.61, 1.6E-39^{***})$ | $(4.55, 6.8E-32^{***})$ | $(2.47, 1.4E-33^{***})$ |
| You(F) | $(-6.18, 4.4E-44^{***})$ | $(-20.47, 9.4E-01)$ | $(-5.99, 1.9E-52^{***})$ |
| Lookahead(M) | $(2.25, 2.8E-22^{***})$ | $(1.06, 1.7E-06^{***})$ | $(1.50, 9.8E-13^{***})$ |
| Lookahead(F) | $(-1.38, 2.2E-07^{***})$ | $(-0.96, 5.8E-04^{***})$ | $(-0.72, 1.7E-03^{**})$ |
| Consistency(M) | $(0.33, 6.1E-02)$ | $(0.46, 3.1E-02^{*})$ | $(0.25, 1.8E-01)$ |
| Consistency(F) | $(0.22, 2.7E-01)$ | $(0.44, 9.4E-02)$ | $(0.09, 6.6E-01)$ |
| Opposite(M) | $(4.80, 1.9E-127^{***})$ | $(3.13, 3.8E-98^{***})$ | $(3.10, 1.7E-116^{***})$ |

Table 28: Structural factors regression analysis for the systems with the greatest opposite binary gender tendency (**Czech**).

| Variable | AMI | ONLINE-A | TranssionMT |
|---|---|---|---|
| Intercept | $(3.89, 1.1E-128^{***})$ | $(1.34, 1.3E-34^{***})$ | $(3.17, 2.5E-89^{***})$ |
| True(M) | $(-22.31, 9.4E-01)$ | $(-4.71, 3.8E-127^{***})$ | $(-6.73, 3.9E-174^{***})$ |
| Adj Stereo(M) | $(-0.44, 1.4E-03^{**})$ | $(-0.01, 9.2E-01)$ | $(-0.94, 6.0E-09^{***})$ |
| Adj Stereo(F) | $(0.34, 7.1E-04^{***})$ | $(0.78, 1.3E-15^{***})$ | $(0.39, 1.2E-04^{***})$ |
| Adj Sentiment(neg) | $(-0.18, 3.9E-02^{*})$ | $(0.09, 3.0E-01)$ | $(0.15, 9.4E-02)$ |
| Adj Type(appearance) | $(-0.09, 3.7E-01)$ | $(0.75, 1.6E-12^{***})$ | $(0.04, 7.2E-01)$ |
| You(M) | $(19.86, 9.4E-01)$ | $(2.19, 3.5E-21^{***})$ | $(3.15, 5.2E-38^{***})$ |
| You(F) | $(-3.78, 3.0E-79^{***})$ | $(-3.21, 4.0E-81^{***})$ | $(-4.68, 3.0E-90^{***})$ |
| Lookahead(M) | $(-1.26, 1.0E-10^{***})$ | $(0.15, 4.2E-01)$ | $(-0.38, 5.3E-02)$ |
| Lookahead(F) | $(0.16, 4.1E-01)$ | $(0.85, 1.7E-05^{***})$ | $(0.73, 1.8E-03^{**})$ |
| Consistency(M) | $(-0.79, 8.1E-07^{***})$ | $(-0.73, 1.4E-09^{***})$ | $(-0.05, 7.5E-01)$ |
| Consistency(F) | $(-0.40, 3.4E-02^{*})$ | $(-0.48, 3.9E-03^{**})$ | $(0.57, 9.4E-03^{**})$ |
| Opposite(M) | $(-1.78, 2.0E-75^{***})$ | $(-0.96, 1.4E-26^{***})$ | $(-0.99, 2.4E-24^{***})$ |

Table 29: Structural factors regression analysis for the systems with the greatest opposite binary gender tendency (**Icelandic**).

# The GenderQueer Test Suite

**Steinunn Rut Friðriksdóttir**
University of Iceland
`srf2@hi.is`

## Abstract

This paper introduces the GenderQueer Test Suite, an evaluation set for assessing machine translation (MT) systems' capabilities in handling gender-diverse and queer-inclusive content, focusing on English to Icelandic translation. The suite evaluates MT systems on various aspects of gender-inclusive translation, including pronoun and adjective agreement, LGBTQIA+ terminology, and the impact of explicit gender specifications.

The 17 MT systems submitted to the WMT24 English-Icelandic track were evaluated. Key findings reveal significant performance differences between large language model-based systems (LLMs) and lightweight models in handling context for gender agreement. Challenges in translating the singular "they" were widespread, while most systems performed relatively well in translating LGBTQIA+ terminology. Accuracy in adjective gender agreement is quite low, with some models struggling particularly with the feminine form.

This evaluation set contributes to the ongoing discussion about inclusive language in MT and natural language processing. By providing a tool for assessing MT systems' handling of gender-diverse content, it aims to enhance the inclusivity of language technology. The methodology and evaluation scripts are made available for adaptation to other languages, promoting further research in this area.

## 1 Introduction

This paper introduces the GenderQueer Test Suite, a novel evaluation set designed to probe MT systems' capabilities in translating gender-diverse and queer-inclusive content. The test suite has been made publicly available and can be adapted to other languages. The test suite aims to address five key areas of evaluation:

1. Pronoun translation: Assessing translation accuracy when translating the third-person pronoun "they" from English to Icelandic with respect to gender agreement.

2. The singular "they": Assessing whether MT systems are able to translate the gender-neutral, singular "they" as it is used in English, i.e. when "they" is used to refer to a single person who is either non-binary, female, or male, to the more rigid grammatical gender system of Icelandic.

3. Adjective agreement: Evaluating the translation of adjectives with respect to gender forms in the target language. Translation accuracy for each gender form is examined individually as well as accuracy for translations of adjectives with positive, negative, and neutral sentiment.

4. LGBTQIA+ terminology: Examining the translation accuracy of LGBTQIA+-specific terms, including an assessment of whether translations are current and culturally appropriate or potentially outdated or inappropriate.

5. Influence of explicit gender information: Investigating whether explicitly defining a subject as cis or trans affects the translation accuracy of "they" compared to that of similar sentences without such specifications.

The test suite primarily consists of short paragraphs (3-4 sentences long) designed to provide context and challenge MT systems across these five dimensions. An additional 16 single-sentence examples are included for comparison between sentence-level and paragraph-level translations. Each passage contains explicit information about the gender of the subject or subjects mentioned. The purpose of the test suite is to highlight the current capabilities and limitations of MT systems in handling gender agreement in morphologically

rich languages such as Icelandic as well as to provide a tool for assessing MT systems' handling of non-binary pronouns and LGBTQIA+ terminology.

The following sections discuss the motivation behind the GenderQueer Test Suite and present the phenomena of interest in more detail. An analysis of the performance of the 17 MT systems submitted during the 9th Conference of Machine Translation (WMT24) for the English-Icelandic language direction follows. Finally, the implications of these findings are discussed.

## 2 Test Suite Details

The text examples in the test suite were manually compiled by the author, who holds a BA degree in Icelandic. The test suite contains 331 text examples in English, stored in a single text file which is to be translated by the MT systems. The test suite also contains a gold standard translation meant for comparison, in which each example has been translated as expected into Icelandic. Uncertainties when translating LGBTQAI+ terminology were handled in collaboration with members of the queer community in Iceland.

Each example begins by explicitly mentioning the gender of the subject or subjects in question. This is done in four ways:

1. These (cis/trans) men/women are my neighbors / This (cis/trans) man and this (cis/trans) woman are my neighbors.

2. This non-binary/genderqueer/genderfluid person is my neighbor.

3. I'm a woman/man. My friends are women/men/a man and a woman.

4. I'm a woman/man. My friends X, Y and C are women/men / My friend X is a woman/man but my friends Z and Y are men/women / My friends X and Y are women/men but my friend Z is a man/woman.

Genders are explicitly stated in a similar format in the single-sentence examples as well: *"These men/women who live next door to me are my neighbors and they..."* By explicitly stating the gender of the subject or subjects, problems that may arise from assumption of gender based on a person's name are avoided. After specifying gender, the text examples then examine the phenomenon or phenomena in question.

## 2.1 Gender: Translating "They"

Text examples 1 through 169 evaluate the translation of the third-person plural pronoun "they" in terms of gender agreement with the subjects, which in these examples are always plural. In the case of Icelandic, there are three grammatical genders that must be accounted for: the feminine (Icelandic: *þær*), the masculine (Icelandic: *þeir*), and the neuter (Icelandic: *þau*)[1]. There are 108 occurrences of the feminine "they", 102 occurrences of the masculine "they", and 150 occurrences of the neuter "they" (for further details, see table 1 in Appendix B). The greater amount of neuter examples owes to various combinations of gender specifications, further discussed in Section 2.5.

Text examples 1 through 72 each include two examples of the third-person pronoun "they" which, in English, is gender-neutral but, as previously stated, must agree with the gender of the subjects in Icelandic. The first example is always the same, i.e. *They live next door to me.* In order to probe for heteronormativity in the translations, each gender is then tested with the sentence *They have two children.* This is compared to the translation of sentences where the subjects have various types of pets (dogs, cats, parrots, and goldfish). The hypothesis is that, in the cases where the subjects are indicated to have children, the MT systems will opt for the neuter gender form, indicating a preference to parents of opposite genders rather than same-sex parents. An example follows:

> **English:** This trans woman and this cis man are my neighbors. They live next door to me. They have two children.
> **Icelandic:** Þessi trans kona og þessi cis maður eru nágrannar mínir. Þau búa við hliðina á mér. Þau eiga tvö börn.

Text examples 73 through 169 include two occurrences of the third-person pronoun "they" as before, but one contains an LGBTQAI+ term indicating the sexuality of the subjects. This is further discussed in Section 2.4. The other example continues to probe for heteronormativity by refering to the fact that the subjects have children. For example:

> **English:** These women are my neighbors. They are lesbians. They have two children.

---

[1]All Icelandic translations mentioned here are in the plural form.

> **Icelandic:** Þessar konur eru nágrannar mínir. *Þær* eru lesbíur. *Þær* eiga tvö börn.

Text examples 266-319 further challenge the MT systems' ability to follow context. The subjects are introduced in the following way: *I'm a wo/man. My friends are (wo)men/a woman and a man.* Directly following is a sentence containing the pronoun *we*, which is not gendered in Icelandic, along with an adjective that must agree with the gender of the subjects (further discussed in Section 2.3). The second sentence contains the pronoun *they* along with a second adjective. This means that the MT system must realize the gender combination of the group as a whole but also make a distinction between the gender of the group and the portion of the group only containing the friends (and therefore the *they*-reference). For example:

> **English:** I'm a woman. My friends are men. We are 25 years old. They are tall.
> **Icelandic:** Ég er kona. Vinir mínir eru menn. Við erum 25 ára gömul. *Þeir* eru hávaxnir.

## 2.2 Gender: The Singular "They"

Text examples 170-211 are designed to be particularly difficult for an English-Icelandic MT system to translate correctly. They all contain a single subject, referenced by the singular "they", which is gender-neutral in English. In Icelandic, no such singular, gender-neutral pronoun exists in reality. The pronoun *hán* has existed in the language since approximately 2010[2] and has been widely adopted by non-binary people in Iceland although other variations exist. It is important to note, however, that unlike the English equivalent, which can refer to an individual of any gender, *hán* is almost never used for people that fall within binary gender norms but rather exclusively for non-binary individuals.

In any case, text examples 170-184 follow the same pattern as described in 2.1 except in these examples, the single subject is defined as a non-binary, genderqueer, or genderfluid person. In the evaluation, a system is awarded 1 point for translating the singular "they" as *hán*. As the plural neuter form is used by some non-binary individuals in Iceland to refer to themselves (in the singular) and to account for the much higher likelihood of the

MT systems recognizing "they" as a plural form, a system is awarded 0.5 points for translating the singular "they" as *þau*. The same is expected from text examples 185-193 which contain adjectives, further discussed in Section 2.3. For example:

> **English:** This non-binary person is my neighbor. They are short. They are an adult.
> **Icelandic (preferred):** Þessi kynsegin manneskja er nágranni minn. *Hán* er lágvaxið. *Hán* er fullorðið.
> **Icelandic (acceptable):** Þessi kynsegin manneskja er nágranni minn. *Þau* eru lágvaxin. *Þau* eru fullorðin.

On the other hand, text examples 194-211 define the single subject as either a man or a woman, which is then also indicated by the singular "they". This requires the MT system to not only recognize the indicated gender of the subject, but also to realize that "they" should not be translated in the plural, but rather as the singular masculine *hann* (English: *he*) or feminine *hún* (English: *she*), respectively. If a system successfully translates this, it is awarded 1 point per occurrence. As it is much more likely that these examples will be translated in the plural, systems are awarded 0.5 points for translating them as the masculine *þeir* or the feminine *þær*, respectively. For example:

> **English:** This woman is my neighbor. They are short. They are an adult.
> **Icelandic (preferred):** Þessi kona er nágranni minn. *Hún* er lágvaxin. *Hún* er fullorðin.
> **Icelandic (acceptable):** Þessi kona er nágranni minn. *Þær* eru lágvaxnar. *Þær* eru fullorðnar.

## 2.3 Gender: Translating Adjectives

Text examples 185-319 each contain two adjectives and examples 320-331 contain three adjectives each[3]. While gender neutral in English, each adjective must agree with the gender of the subjects in Icelandic. The MT systems are thus evaluated

---

[3]In this case, LGBTQAI+ terms are not considered adjectives though most of them certainly qualify as such. The adjectives in question are all generic and describe people's traits, i.a. *hungry*, *boring* or *funny*

based on their overall accuracy in translating these adjectives with respect to their gender forms.[4]

These examples vary in difficulty. The most difficult (besides those containing the singular "they", discussed in Section 2.2) can be found in text examples 320-331, which indicate the gender of four different, named subjects: *I'm a woman/man. My friends X, Y and C are women/men / My friend X is a woman/man but my friends Z and Y are men/women / My friends X and Y are women/men but my friend Z is a man/woman.* Directly following is a sentence containing the pronoun *we* along with an adjective that must agree with the gender of the group as a whole. The second sentence contains a reference to the subjects' names along with two adjectives whereby each adjective must agree with half of the group: *X and I are smart but Y and Z are dumb.* An example follows:

> **English:** I'm a woman. My friends Mary and Sophia are women but my friend John is a man. We are 25 years old. Mary and I are smart but John and Sophia are dumb.
>
> **Icelandic:** Ég er kona. Vinkonur mínar, Mary og Sophia eru konur en vinur minn John er maður. Við erum 25 ára *gömul*. Við Mary erum *gáfaðar* en John og Sophia eru *heimsk*.

Additionally, accuracy for each gender is examined individually as well as the accuracy for translations of adjectives with a positive, negative or neutral sentiment. The hypothesis here is that if a model only translates adjectives for a particular gender correctly if the adjectives convey a certain sentiment, a gender bias within the model is indicated. An example of this can be found in Sólmundsdóttir et al. (2022) where MT systems tended to translate adjectives with a negative connotation more frequently as feminine, while adjectives with a positive connotation were more likely to be translated as masculine, except when the adjective described a person's appearance, where the opposite was the case.

### 2.4 Queer: Translating LGBTQAI+ Terms

Text examples 33 through 193 each contain at least one LGBTQAI+ term. While most of these terms

are adjectives and could (and should, perhaps) be evaluated based on gender agreement like the adjectives discussed in Section 2.3, these terms are only evaluated based on the quality of the translations themselves (in other words: whether or not the correct term is used in the translation, regardless of gender form). This is done to place more emphasis on the importance of the words themselves rather than grammatically perfect translations. Additionally, they represent a vocabulary that is highly connected to a person's sense of self and should therefore be examined individually in order to account for inclusive language in MT systems.

In total, there are 283 terms to be translated. The systems are evaluated in two ways. Firstly, each system receives an accuracy score based on whether or not the translation of the term exists in the accompanying terminology database. If it does, the system is awarded 1 point. There are three exceptions to this. If a system translates *trans woman* or *trans man* as a compound (for instance *transkona* instead of *trans kona*, with *trans* as a prefix rather than an adjective), it receives only 0.5 points along with a warning indicating that the use of the compound is considered inappropriate by many trans people in Iceland. The same goes for translations where *trans* and *cis* are translated as *transkynja* and *sískynja*, respectively. While these terms exist in the language, they are hardly ever used and should be avoided according to members of the queer community. Similarly, while unlikely to come up as translations at all, if a system translates the terms *lesbians* and *bisexual* as *lessur* and *bæjarar*, respectively, the system receives 0.5 points along with a warning indicating that these terms are only considered appropriate if used by the people they refer to and should be avoided as general terms.

Secondly, the MT systems receive a score based on the proportion of terms translated in an inappropriate manner as determined by the terminology database. These might include outdated translations that are no longer in use or crude terms that are considered slurs. The purpose is to separate the use of these terms from translations that are plainly wrong for the context. A model that uses the inappropriate terms should be considered more harmful to LGBTQAI+ individuals than a model that simply translates the terminology incorrectly. In other words, a high inappropriate score is a clear indicator of bias against LGBTQAI+ individuals in the respective model.

---

[4]It should be noted that the database used for determining the correct translations might not be exhaustive in terms of possible translations for these adjectives, so some translations might be misidentified as incorrect. There should, however, be very few such instances.

## 2.5 Queer: Specificity of Gender

The GenderQueer Test Suite allows for a comparison of translations of the third person plural pronoun "they" based on the specificity of the gender in question. In other words, it is possible to examine whether specifying a subject as either cis or trans leads to a poorer outcome than if the genders are not defined in this manner. Each gender combination is examined, i.e. *trans women*, *trans men*, *cis women*, *cis men*, *a trans woman and a trans man*, *a cis woman and a cis man*, *a trans woman and a cis man*, and *a cis woman and a trans man*. The process is otherwise the same as described in Section 2.1, including a comparison of text examples involving a reference to the subjects having children and examples where there is no mention of children.

## 3 Evaluation

Every aspect of the evaluation of the GenderQueer Test Suite has been automated and made available with an CC-BY license on Github[5]. The following sections will discuss notable results in the evaluation of the WMT24 English-Icelandic MT systems. Figures and tables referenced can be found in Appendices A and B, respectively.

### 3.1 Pronoun Translations and Explicit Gender Information

Figure 1 shows the overall translation accuracy of "they" translations (both plural and singular) and compares the text examples containing a single sentence to the text examples containing at least three sentences. This refers to whether or not the models respect the gender agreement with the subject or subjects. As the number of "they"-occurrences in the short examples (16 in total) is much lower than that of the longer ones (444 in total), these results should only be taken as indicative and not conclusive. However, it is clear that many models struggle much more with translating the longer examples, indicating that the problem of paragraph-level translations remains to be fully solved.

Figure 2 breaks down the accuracy of these translations per gender. Each gender is again broken down in terms of specific definitions, i.e. whether or not the subjects are explicitly defined as cis or trans. All models struggle with translating the singular "they", with no model achieving accuracy above 40.5% (GPT-4). This may not be surprising,

as widespread use of the singular "they" in both languages is relatively new and so the training data for these models might not include a lot of examples of it in use. It is, however, important to take note of social development and include gender-inclusive language when developing such models.

The difference between the performance of LLM-based systems and lightweight systems in handling gender agreement at the paragraph-level is striking. While most of the LLMs receive a near-perfect score in this regard, the lightweight models rarely achieve more than 60% accuracy and all of them seem to almost entirely exclude feminine forms from their translations. It is somewhat expected that the masculine form dominates in these translations, as it has traditionally been used to refer to a group of mixed-gendered people or to refer to a person or persons of unknown genders[6]. This certainly seems to be the case for Aya23, where the masculine is predicted in 100% of the cases.

On the other hand, a preference for the neuter form might indicate a heteronormative bias in the models, particularly in text examples involving a reference to the subjects having children. Interestingly, when Figures 3 and 4 are compared, this preference is more pronounced in text examples where children are not mentioned. It should, however, be noted that the latter are fewer in total; the comparison should be considered as preliminary. However, it is clear that the limited use of the feminine form indicates some form of bias, either linguistic, societal, or a combination of the two.

In general, there does not seem to be much difference in accuracy between explicit gender definitions and those that do not specify the gender as either cis or trans. Rather, some of the models seem to struggle the most with a combination of more than one gender, i.e. the neuter form, where the subjects are defined individually (*This woman **and** this man...*). While this may seem to contradict the heteronormative hypothesis, Figure 3 shows that these models will in general translate the examples involving children a lot more accurately than the examples that contain no reference to children, further indicating that the hypothesis holds true to a significant extent.

---

[6]For further discussion on the generic masculine in Icelandic, see for instance Section 5 in Friðriksdóttir and Einarsson (2024).

## 3.2 Adjective Agreement

Figure 6 reveals that no model performs perfectly in the case of gender agreement between subjects and adjectives, with accuracy ranging from 88.89% (Claude-3.5) to 0.3% (TSU-HITs). As discussed in Section 2.3, some of the examples involving adjective translations are quite complex and the relatively poor performance of the models overall might simply be due to this. On the other hand, it is again noticeable how many models struggle the most with translations in the feminine form. It is interesting to note that in general, most of the correctly translated adjectives in the feminine form seem to have a positive sentiment and the same holds true for the correctly translated adjectives in the neuter form. For the masculine, however, most of the correctly translated adjectives have either a negative or a neutral connotation. This might indicate a gender bias.

## 3.3 LGBTQAI+ Terminology

Most models do relatively well on the translation of LGBTQAI+ terminology, as indicated by Figure 5, averaging at about 70% in overall accuracy and never exceeding 6.01% in terms of inappropriate translations. Not surprisingly, the models that have a decent overall translation score are also more likely to have more instances of inappropriate vocabulary. While the overall performance of the models is relatively good in this regard, researchers must make sure that their training data does not include excessive (or any) harmful slurs about minority groups to prevent inappropriate terms from becoming the default translations for this terminology.

## 4 Conclusion and Future Work

The GenderQueer Test Suite provides valuable insights into the capabilities and limitations of MT systems in handling gender-diverse and queer-inclusive translations from English to Icelandic. The evaluation of the 17 MT systems submitted to WMT24 revealed that LLM-based systems generally outperform lightweight models in terms of gender agreement in paragraph-level translations. All systems struggled with translating the singular "they", highlighting the importance of incorporating gender-inclusive language in the training data for such models. While LGBTQIA+ terminology was generally translated accurately, the higher performing models still sometimes use outdated

or derogatory vocabulary which could potentially cause direct harm to minority groups if used as the default translations of these terms.

Future work should focus on expanding the test suite to cover more language pairs and incorporating more diverse gender identities and expressions. Collaboration with LGBTQIA+ communities will ensure that the test suite keeps up with evolving terminology and language use. Exploring the integration of the GenderQueer Test Suite into standard MT evaluation pipelines could promote consistent attention to gender-inclusive translation across the field. This can drive progress towards more inclusive and accurate MT systems that respect and represent the full spectrum of gender identities. The test suite has been made openly available and other researchers are encouraged to adapt it to their languages.

## Limitations

While the GenderQueer Test Suite offers valuable insights into machine translation of gender-diverse content, several limitations should be acknowledged:

Language Specificity: The test suite is designed for English to Icelandic translation. The complex gender system of Icelandic presents unique challenges that may not generalize to languages with different grammatical structures or those lacking grammatical gender.

Scope of Gender Diversity: Despite efforts to include a range of gender identities, the test suite may not fully capture the entire spectrum of gender diversity, potentially oversimplifying some nuances. Additionally, limited number of text examples for certain tasks may skew the results.

Evolving Language: The rapidly changing nature of gender and sexuality means some terms in the test suite may become outdated, necessitating regular updates.

Evaluation Method: The evaluation of the translation of the third person plural pronoun "they" compares the number of correct translations with respect to gender forms to the total number of "they" occurrences in the English text examples. However, some models might drop one or more occurrences from their translations. An example of this can be seen in the AMI model's translation:

> **English**: This woman and this man are my neighbors. They are bisexual. They have two children.

> **Icelandic**: Þessi kona og maðurinn eru
> nágrannar mínir. Þau eru tvíkynhneigð
> og eiga tvö börn.

This is a perfectly valid translation despite dropping the second "they". Due to the evaluation method, this will still hurt the measured accuracy of the model.

## Ethics Statement

Some of the inappropriate translations included in the database used to evaluate LGBTQAI+ vocabulary are disrespectful and harmful to minority groups. These terms are included as a means to evaluate the presence of bias in the MT systems and their use in any context is highly discouraged.

## Acknowledgements

## References

Steinunn Rut Friðriksdóttir and Hafsteinn Einarsson. 2024. Gendered grammar or ingrained bias? exploring gender bias in Icelandic language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 7596–7610, Torino, Italia. ELRA and ICCL.

Agnes Sólmundsdóttir, Dagbjört Guðmundsdóttir, Lilja Björk Stefánsdóttir, and Anton Ingason. 2022. Mean machine translations: On gender bias in icelandic machine translations. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3113–3121.

# A  Graphs



Accuracy when Translating Text with Single vs. Multiple Sentences

Figure 1: Translation accuracy for text examples containing a single sentence as opposed to text examples containing at least three sentences. This refers to the translation of the third person plural pronoun "they" with respect to gender forms, i.e. whether or not the models respect the gender agreement with the subject, explicitly presented in the first sentence of the longer examples and in the first phrase of the shorter examples. It also includes translations of the singular "they", which refers to a single person who is either non-binary, female, or male. It should be noted that the number of short examples is much lower than that of the longer examples and the comparison should therefore be taken as indicative and not conclusive. Still, we can see that the models struggle much more with following the context of the longer examples, indicating that paragraph-based translations are still at least somewhat problematic.

Figure 2: Translation accuracy of the third person plural pronoun "they" with respect to gender forms, i.e. how often the models respect the gender agreement with the subject, explicitly presented in the first sentence of the text examples. It also includes translations of the singular "they", which refers to a single person who is either non-binary, female, or male. Note that the results presented on this heatmap only apply to the longer examples, i.e. text examples that contain at least three sentences. The first column refers to the overall accuracy of the models. The heatmap then shows the translation accuracy for each gender. Each gender is broken down depending on whether or not the subject is explicitly defined as either cis or trans. We can see that every model struggles with translating the singular "they" and the lightweight models almost entirely exclude the feminine form from their translations.

Figure 3: Translation accuracy of the third person plural pronoun "they" with respect to gender forms, i.e. how often the models respect the gender agreement with the subject, explicitly presented in the first sentence of the text examples. It also includes translations of the singular "they", which refers to a single person who is either non-binary, female, or male. Note that the results presented on this heatmap only apply to the longer examples, i.e. text examples that contain at least three sentences. All of the examples presented here contain a reference to the subjects having children (their last sentence being "They have two children"). We can see that all of the models struggle with the singular "they" but otherwise, the translation accuracy seems to depend almost entirely on the architecture of the model, with LLM-based systems outperforming the lightweight models. It is interesting to note that the lightweight models struggle the most with the feminine form, while the performance when translating the neuter and the masculine form is relatively even. The hypothesis was that the models would default to the neuter form, indicating heteronormativity. On the other hand, the masculine form is the one traditionally used as the general form, such as when the gender of the subject is unknown or the subjects are mixed-gendered. These results could therefore indicate a twofold bias, one linguistic in nature and the other societal.

Figure 4: Translation accuracy of the third person plural pronoun "they" with respect to gender forms, i.e. how often the models respect the gender agreement with the subject, explicitly presented in the first sentence of the text examples. It also includes translations of the singular "they", which refers to a single person who is either non-binary, female, or male. Note that the results presented on this heatmap only apply to the longer examples, i.e. text examples that contain at least three sentences. Here, the text examples do not contain a reference of the subjects having children. We again see that all of the models struggle with translating the singular "they" and that the accuracy of the LLM-based models is much higher than that of the lightweight models. The latter perform best on the neuter form with the feminine form almost not appearing at all. On the other hand, half of the better-performing models struggle with the neuter form, some of which do not predict it at all. While this is interesting and could potentially indicate a bias, it should be noted that these examples are fewer than those containing references to the subjects having children and so the comparison should be taken as indicative rather than conclusive.

Figure 5: Translation accuracy for LGBTQAI+ terminology. The models are tested for appropriate and inappropriate translations. The latter refers to terms that are either outdated, prejudiced, or otherwise not advisable but not entirely wrong in the sense that they are accurate but harmful translations of the English terms. The higher the red bar, the more harm the model might cause to minority groups.

Figure 6: Translation accuracy for adjectives with respect to gender forms. The first column refers to the overall accuracy of each model, i.e. the proportion of adjectives that were translated correctly in the sense that they respect the gender agreement with the subject, explicitly presented in the first sentence of the text example. The heatmap breaks down the translation accuracy for each gender and for each gender, the accuracy for each sentiment is observed. Again, most of the systems struggle the most with the feminine form. On the other hand, most of the correctly translated adjectives in the feminine form have a positive sentiment, while correctly translated adjectives in the masculine form more often have either a neutral or a negative sentiment. This could potentially indicate a gender bias.

## B Tables

| | Total | Long (≥ 3 sentences) | Short (single sentence) |
|---|---|---|---|
| Text examples | 331 | 315 | 16 |
| "They" | 460 | 444 | 16 |
| LBGTQAI+ terms | 283 | 283 | 0 |
| Adjectives | 306 | 306 | 0 |

Table 1: The overall occurrences of each phenomena in the GenderQueer Test Suite as indicated by the gold standard translation.

| | Total | Positive | Negative | Neutral | English | Icelandic (singular/plural) |
|---|---|---|---|---|---|---|
| Feminine | 71 | 24 | 25 | 22 | young | ung/ungar |
| Masculine | 71 | 24 | 25 | 22 | young | ungur/ungir |
| Neuter | 164 | 54 | 52 | 58 | young | ungt/ung |

Table 2: The occurrences of adjectives in the GenderQueer Test Suite as indicated by the gold standard translation. The overall occurrences of each gender form are presented along with a breakdown of the sentiments attached to the adjectives. The translation examples show the declensions with respect to the number and gender of the subject(s).

| | Total | Unsp. (C) | Unsp. (NC) | Trans (C) | Trans (NC) | Cis (C) | Cis (NC) | Cis and trans (C) | Cis and trans (NC) | English | Icelandic |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Feminine | 108 | 22 | 8 | 22 | 8 | 22 | 8 | 0 | 0 | she/they | hún/þær |
| Masculine | 102 | 20 | 8 | 20 | 8 | 20 | 8 | 0 | 0 | he/they | hann/þeir |
| Neuter | 150 | 18 | 10 | 18 | 8 | 18 | 8 | 36 | 16 | it*/they | það*/**þau** |
| Singular they | 84 | 6 | 78 | 0 | 0 | 0 | 0 | 0 | 0 | they/they | **hán**/þau |

Table 3: The occurrences of the third person plural pronoun "they" in the GenderQueer Test Suite as indicated by the gold standard translation. Also included are the occurrences of the singular "they", referring to a single person which can be non-binary, female, or male. The overall occurrences of each gender are presented along with a breakdown referring to whether or not the gender definitions are explicit, i.e. if "cis" or "trans" is specified. "C" refers to examples that include a reference to the subjects having children, i.e. where the last sentence of the text example is "they have two children". "NC" refers to examples where there is no reference to the subjects having children. Examples where one person is defined to be cis and the other as trans were limited to that of the neuter gender form, where one person is a woman and the other a man. The translation examples show the declensions with respect to the number and gender of the subject(s). It should be noted that, while the traditional translation of the third person singular in the neuter form, *það* is never used to refer to a person. Rather, *hán* is used in this case. Both the traditional neuter (referring to a mixed-gendered group of people) and the plural form of the singular "they" is *þau*.

# Domain Dynamics: Evaluating Large Language Models in English-Hindi Translation

**Soham Bhattacharjee , Baban Gain**
Indian Institute of Technology, Patna


**Asif Ekbal**
Indian Institute of Technology, Jodhpur
{sohambhattacharjeenghss,gainbaban,asif.ekbal}@gmail.com

## Abstract

Large Language Models (LLMs) have demonstrated impressive capabilities in machine translation, leveraging extensive pre-training on vast amounts of data. However, this generalist training often overlooks domain-specific nuances, leading to potential difficulties when translating specialized texts. In this study, we present a multi-domain test suite, collated from previously published datasets, designed to challenge and evaluate the translation abilities of LLMs. The test suite encompasses diverse domains such as judicial, education, literature (specifically religious texts), and noisy user-generated content from online product reviews and forums like Reddit. Each domain consists of approximately 250-300 sentences, carefully curated and randomized in the final compilation. This English-to-Hindi dataset aims to evaluate and expose the limitations of LLM-based translation systems, offering valuable insights into areas requiring further research and development. We have submitted the dataset to WMT24 *Break the LLM* subtask. In this paper, we present our findings. We have made the code and the dataset publicly available at https://github.com/sohamb37/wmt24-test-suite.

## 1 Introduction

Machine translation (MT) (Bahdanau et al., 2016) has witnessed significant advancements with the advent of Large Language Models (LLMs) (et al., 2024a,b), which leverage extensive pretraining on massive datasets to achieve high performance across various language pairs (Alves et al., 2024; Zhu et al., 2024; Zhang et al., 2023). Despite their remarkable generalization capabilities, LLMs often struggle with domain-specific texts due to a lack of targeted training on such specialized content (Robinson et al., 2023; Jiao et al., 2023; Hendy et al., 2023). Some LLMs (Workshop et al., 2023) generate good translation involving low-resource

language when the target language is English but not the other way around (Bawden and Yvon, 2023). These challenges are amplified when the domains involved are different from those of training data. This limitation poses a challenge for deploying MT systems in real-world applications where domain-specific accuracy is crucial.

To address this gap, we participated in the "Help us break LLMs" subtask at the Workshop on Machine Translation (WMT) 2024 (Kocmi et al., 2024). The primary objective of this subtask is to create a dataset that exposes the difficulties faced by LLM-based MT systems when dealing with domain-specific content. Our approach involves collating a multi-domain dataset that includes sentences from judicial, educational, religious literature, and noisy user-generated content from online product reviews and forums like Reddit.

Each domain-specific subset comprises approximately 250-300 sentences, which are then randomized to form the final dataset. This dataset, focusing on the English-to-Hindi translation direction, aims to rigorously test the robustness and adaptability of LLM-based MT systems. By identifying the translation challenges specific to each domain, our study provides valuable insights for improving domain adaptation techniques in machine translation, ultimately contributing to more reliable and accurate MT solutions for specialized applications. Our contributions to the paper are as follows:

- We participate in the Break the LLM challenge in WMT24 for English-Hindi language direction, where we submit diverse data consisting of six domains.

- We calculate the standard BLEU score as well as the state-of-the-art metric xCOMET-XXL to evaluate the translation quality.

- We perform a tiny scale manual evaluation of the translation outputs.

## 2 Related Works

Neural Machine Translation has achieved significant advancements (Vaswani et al., 2017). However, translation of text involving low-resource languages remains a challenge. In low-resource languages, the translations of Indic languages like Hindi is difficult due to the paucity of the high-quality parallel corpus. Existing multilingual models like IndicTrans (Ramesh et al., 2022) and IndicTrans2 (Gala et al., 2023) achieved significant performance gains compared to other models. However, English-Hindi machine translations still have room for improvement.

Moslem et al. (2022) has previously used pre-trained Language Models(LM) for domain specific data augmentation for Machine Translation. They simulated the characteristics of a small bilingual dataset or monolingual source text and combined it with back translation to create huge amounts of synthetic in-domain data. Other works involving low-resource languages include translation of chat-based conversation by (Gain et al., 2022) where English Hindi translation was implemented on Chat and question answers in chatbots. In the domain of education, (Behnke et al., 2018) used crowd-sourcing English texts to obtain translation into 11 languages for generating NMT data. Similarly, Ramakrishna et al. (2023) introduced the EduMT dataset for improving the English-Hindi translation for educational content.

In a recent study, (Briva-Iglesias et al., 2024) showed that LLMs outperform Google translate when it comes to the Legal domain. (Martínez-Domínguez et al., 2020) implemented machine translation in the legal domain for Italian to Swiss language. For low-resource language, (Poudel et al., 2024) introduced a custom-built dataset for the legal domain for English Nepali language machine translation.

In the Literary domain, (Drobot, 2023) has studied the prospects of neural machine translation. Earlier (Matusov, 2019) has used NMT for translating German literary works to English, and (Kuzman et al., 2019) implemented NMT for the literary domain from English to Slovene. (Yirmibeşoğlu et al., 2023) has implemented NMT in the literary domain for the low-resource language of English-Turkish. (Thai et al., 2022) has also explored document-level literary machine translation for non-English languages. They have also shown that there is a disparity between the automatic evaluation of these machine translations and human evaluation, prompting further improvement of machine translation in this domain.

Noisy or non-standard input text can cause disastrous mistranslations in most modern Machine Translation (MT) systems.Khayrallah and Koehn (2018) has shown in a study the impact of noise on NMT systems. Michel and Neubig (2018) proposed a benchmark dataset for machine translation of noisy texts(MTNT). Herold et al. (2022) has worked on filtering noise from machine translation data for improving the performance of NMT systems.Bolding et al. (2023) has used LLMs to remove noise from the MTNT dataset target sentences and proposed C-MTNT dataset. Machine Translation of noisy text is mainly explored through multimodal translation in English-Hindi (Gain et al., 2021b; Laskar et al., 2021; Gain et al., 2021a; Gupta et al., 2021c; Gain et al., 2023) where images features were used to aid in machine translation from English to Hindi.

Product review is a translation task that is related to the field of e-commerce. (Gupta et al., 2022) explores NMT with sentiment preservation in this domain for the low-resource language of the English-Hindi pair. (Gupta et al., 2021b) and (Gupta et al., 2021a) are some of the other works on online product review translation.

Some other notable works on low-resource languages include (Goyle et al., 2023), (Chowdhury et al., 2022) and (Ranathunga et al., 2023) that have implemented unique NMT techniques to complement the scarcity of data in these languages.

## 3 Dataset

Our proposed dataset includes English-Hindi bitext pairs from six critical domains, chosen for their significance to both the machine translation community and their difficulty of translation. We provide a sample from each domain in Appendix D and some statistics about the datasets in Table 4. It can be noted that the size of each domain is different. We had collected 500 sentences from each domain in the beginning but after filtering out sentences less than 5 words, we arrived at the final size of the dataset.

### 3.1 Education domain

The education domain plays a crucial role in knowledge dissemination. Enhancing machine translation in education promotes equal access to qual-

| Model | Education | | | General | | | Judicial | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | COMET | HUMAN | BLEU | COMET | HUMAN | BLEU | COMET | HUMAN |
| Aya23 | 36.40 | 0.71 | 2.00 | 14.13 | 0.70 | 3.33 | 17.07 | 0.70 | 4.00 |
| Claude3.5 | 46.04 | 0.80 | 3.33 | 19.02 | 0.85 | 3.67 | 25.62 | 0.85 | 3.67 |
| CommandR-plus | 35.33 | 0.75 | 3.67 | 14.39 | 0.77 | 3.67 | 17.64 | 0.77 | 3.00 |
| CycleL | 0.38 | 0.72 | 1.33 | 1.21 | 0.15 | 0.79 | 1.33 | 0.14 | 1.00 |
| GPT-4 | 40.90 | 0.68 | 2.67 | 14.68 | 0.75 | 2.67 | 18.45 | 0.75 | 2.67 |
| IKUN-C | 28.99 | 0.75 | 2.67 | 11.60 | 0.67 | 3.00 | 8.21 | 0.50 | 2.33 |
| IKUN | 28.62 | 0.76 | 1.33 | 11.99 | 0.66 | 2.33 | 6.95 | 0.47 | 1.00 |
| IOL-Research | 40.47 | 0.67 | 2.00 | 15.41 | 0.77 | 4.0 | 19.12 | 0.78 | 3.33 |
| Llama3-70B | 45.73 | 0.75 | 3.00 | 15.58 | 0.77 | 3.0 | 21.27 | 0.77 | 3.00 |
| NVIDIA-NeMo | 45.12 | 0.82 | 3.00 | 18.12 | 0.66 | 3.67 | 21.21 | 0.69 | 1.33 |
| Online-A | 50.27 | 0.73 | 3.00 | 19.84 | 0.75 | 4.0 | 25.02 | 0.73 | 3.33 |
| Online-B | 46.19 | 0.82 | 4.00 | 21.36 | 0.85 | 4.0 | 25.20 | 0.86 | 3.67 |
| Online-G | 46.19 | 0.73 | 2.67 | 16.49 | 0.67 | 3.67 | 27.33 | 0.73 | 2.67 |
| TransmissionMT | 46.70 | 0.82 | 3.67 | 21.39 | 0.85 | 4.67 | 25.25 | 0.86 | 4.00 |
| Unbabel-Tower-70B | 44.22 | 0.80 | 4.33 | 20.50 | 0.83 | 4.67 | 22.04 | 0.83 | 3.67 |
| ZMT | 50.27 | 0.72 | 3.67 | 19.83 | 0.75 | 4.0 | 25.01 | 0.73 | 3.33 |

Table 1: Performance of different models across education, general and judicial domains

ity learning, supports multilingual environments, and empowers non-native speakers to engage with content. This helps reduce educational disparities and fosters cultural exchange. For this study, 330 English-Hindi language pairs were collected from the EduMT dataset, which focuses on educational content in Indian languages (Appicharla et al., 2021).

## 3.2 General domain

The general domain in our dataset is sourced from the IIT Bombay English-Hindi Parallel Corpus (Kunchukuttan et al., 2018), which includes a diverse range of parallel and monolingual Hindi texts compiled by the Center for Indian Language Technology. It features content from various sources such as news articles, TED Talks, government websites, and Wikipedia. For our study, we randomly selected 500 English-Hindi language pairs from this domain. Improving machine translation in the general domain enhances the accuracy of translations across diverse content, making information more accessible for Hindi-speaking audiences.

## 3.3 Judicial domain

The judicial domain in our dataset is sourced from the IIT Patna Hindi-English Machine Aided Translation (HEMAT) training corpora, which is specifically designed for legal and judicial content. For this domain, we have included 325 sentences in our proposed dataset. Enhancing machine translation performance in the judicial domain is crucial, as it ensures that legal documents, court rulings, and

other judicial materials are accurately translated. This can have a significant impact by improving access to legal information, supporting multilingual legal proceedings, and ensuring that individuals who speak Hindi can fully understand and engage with the judicial system.

## 3.4 Religious Literature domain

The religious literature domain in our dataset consists of 300 pairs: 150 Quran verses from the Tanzil Project [1] and 150 Bible verses from the Bible Eudin Project, both sourced from the OPUS collection (Tiedemann, 2012). These texts pose unique challenges due to their religious significance and archaic language.

## 3.5 Noisy domain

The noisy user-generated data domain in our dataset is sourced from the benchmark dataset for Machine Translation of Noisy Text (MTNT) (Michel and Neubig, 2018). This domain includes 350 English sentences from MTNT, consisting of informal and often error-prone comments made by users on Reddit. Our annotators translated these sentences into Hindi retaining the tone and nature of the input sentences. However, they got rid of some noise based on their own discretion. This domain captures the informality of online communication. Improving machine translation in this domain will help models better handle slangs, typos, and non-standard language use, in turn making

---

[1] https://tanzil.net/docs/tanzil_project

| Model | Literature | | | Noisy | | | Review | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | COMET | HUMAN | BLEU | COMET | HUMAN | BLEU | COMET | HUMAN |
| Aya23 | 8.34 | 0.75 | 2.67 | 31.76 | 0.51 | 3.00 | 30.82 | 0.78 | 3.00 |
| Claude3.5 | 15.11 | 0.90 | 3.33 | 42.49 | 0.71 | 4.33 | 36.45 | 0.89 | 3.33 |
| CommandR-plus | 10.32 | 0.83 | 3.33 | 31.35 | 0.62 | 3.67 | 26.49 | 0.85 | 3.33 |
| CycleL | 0.21 | 0.14 | 1.00 | 0.82 | 0.14 | 1.00 | 0.33 | 0.14 | 1.00 |
| GPT-4 | 7.95 | 0.80 | 2.67 | 35.43 | 0.60 | 3.67 | 33.66 | 0.84 | 2.33 |
| IKUN-C | 4.85 | 0.68 | 2.0 | 19.99 | 0.54 | 2.33 | 19.09 | 0.69 | 1.33 |
| IKUN | 4.80 | 0.70 | 1.33 | 18.89 | 0.54 | 2.00 | 16.48 | 0.60 | 1.33 |
| IOL-Research | 6.82 | 0.82 | 3.00 | 39.79 | 0.62 | 3.33 | 33.23 | 0.84 | 2.67 |
| Llama3-70B | 9.51 | 0.83 | 2.67 | 34.73 | 0.61 | 3.67 | 33.16 | 0.82 | 2.67 |
| NVIDIA-NeMo | 16.65 | 0.72 | 1.0 | 37.32 | 0.38 | 2.33 | 41.07 | 0.61 | 2.00 |
| Online-A | 20.34 | 0.81 | 2.0 | 52.55 | 0.49 | 3.00 | 46.78 | 0.74 | 3.00 |
| Online-B | 26.21 | 0.91 | 3.33 | 51.51 | 0.72 | 2.67 | 41.55 | 0.88 | 3.00 |
| Online-G | 8.56 | 0.69 | 1.67 | 44.13 | 0.44 | 3.33 | 55.29 | 0.72 | 4.00 |
| TransmissionMT | 26.27 | 0.91 | 3.33 | 51.71 | 0.72 | 3.67 | 41.58 | 0.88 | 3.33 |
| Unbabel-Tower-70B | 20.03 | 0.90 | 2.67 | 40.86 | 0.68 | 3.00 | 35.42 | 0.90 | 4.00 |
| ZMT | 20.34 | 0.81 | 1.67 | 52.55 | 0.49 | 2.67 | 46.78 | 0.74 | 3.00 |

Table 2: Performance of different models across literature, noisy, and review domains

them more robust.

## 3.6 Online User Review domain

The final domain in our dataset consists of user product reviews from the e-commerce site Flipkart (Gupta et al., 2021b). We included 300 English-Hindi text pairs from this corpus. This domain presents challenges like grammatical errors and code-mixing, where users blend English and Hindi within a sentence. Similar to MTNT, overcoming the challenges in this domain will make the MT systems more robust.

## 4 Evaluation

In this section, we outline the various evaluation techniques employed to assess the performance of the models based on their outputs. The evaluation metrics considered in this study are the BLEU (Papineni et al., 2002; Post, 2018) score, COMET (Rei et al., 2020; Guerreiro et al., 2023) score, and human evaluation score. We have shared the candidate translations from 3 models, Online-B, Nvidia-Nemo, and INKUN-C in Appendix D. Online B is one of the consistently best performing models across all the domains and metrics among all the submissions. Whereas, Nvidia-Nemo and IKUN-C translations are of lower quality. This table gives us a comparison of the quality of translations by these models.

| Model | BLEU | COMET | HUMAN |
|---|---|---|---|
| Aya23 | 23.53 | 0.69 | 3.00 |
| Claude3.5 | 31.63 | 0.83 | 3.61 |
| CommandR-plus | 23.28 | 0.76 | 3.44 |
| CycleL | 0.78 | 0.14 | 1.11 |
| GPT-4 | 25.98 | 0.74 | 2.78 |
| IKUN-C | 16.70 | 0.63 | 2.28 |
| IKUN | 16.44 | 0.61 | 1.56 |
| IOL-Research | 26.79 | 0.76 | 3.06 |
| Llama3-70B | 26.18 | 0.76 | 3.00 |
| NVIDIA-NeMo | 29.81 | 0.62 | 2.22 |
| Online-A | 36.21 | 0.84 | 3.06 |
| Online-B | 35.92 | 0.71 | 3.44 |
| Online-G | 32.79 | 0.66 | 3.00 |
| TransmissionMT | 35.94 | 0.84 | 3.78 |
| Unbabel-Tower-70B | 31.30 | 0.82 | 3.72 |
| ZMT | 36.20 | 0.71 | 3.06 |

Table 3: Performance of models on the full dataset

## 4.1 BLEU Scores

The BLEU score measures the quality of machine translations by comparing the output to reference translations based on n-gram similarity. A higher n-gram match leads to a higher score, with a brevity penalty to discourage overly short translations. The score ranges from 0 to 100, with higher values indicating better alignment with the references. We calculate the BLEU score with sacrebleu (Post, 2018) and report corpus_score for the dataset.

Figure 1: COMET scores in the Education Domain



Figure 3: COMET scores in the Judicial Domain



Figure 2: COMET scores in the General Domain



Figure 4: COMET scores in the Literature Domain

### 4.1.1 Domain wise Overview

The average BLEU scores in the general, judicial, and literature domains are significantly lower, with scores of 15.97, 19.14, and 12.89, respectively. In the literature domain, the frequent use of ornamental language often leads to subjective translations Table 5, causing notable differences between the machine translations and the reference texts. The general domain, encompassing diverse subdomains and characterized by longer sentence lengths and larger data size Figure 17, also contributes to lower BLEU scores, as models struggle with both factors. Similarly, the judicial domain presents challenges due to its specialized terminology and formal tone, which are difficult for models to translate accurately. Additionally, in all three domains, transliteration instead of translation in many cases further impacts the models' performance.

For the education domain, the sentences are relatively straightforward and easier to translate. Interestingly, the models also achieved relatively high BLEU scores for the user-generated data domains, including noisy texts and product review texts.

### 4.1.2 Model wise Overview

Here we can see the average performance of the models based on all the domains. Models Online-

A and ZMT have the best performance, closely followed by Online-B and TransmissionMT, while CycleL has the worst BLEU scores across all the different domains. Note that BLEU is calculated based on N-gram overlaps. Therefore, transliterations of some tokens, even if they are relevant, are not considered. This results in lower BLEU scores in certain models, even if translation quality is acceptable.

## 4.2 COMET Scores

The COMET score evaluates machine translations using pre-trained language models, focusing on both adequacy (preserving meaning) and fluency (naturalness). It compares machine translations to references and human translations through a regression model trained on human judgments, capturing language nuances that other metrics may miss. The score reflects how closely the machine translation aligns with human preferences. We use xCOMET-XXL to calculate the scores.

### 4.2.1 Domain wise Overview

The COMET scores of judicial, general, and education domains are the highest. It is easier to retain the adequacy and fluency for these domains compared to the other domains. They have a formal tone to them, and the COMET score does not

Figure 5: COMET scores in the Noisy Domain



Figure 6: COMET scores in the Product Review Domain

penalize the MT models much for paraphrasing sentences since it is a more robust metric.

Likewise, the worst COMET scores are obtained for the domains of user-generated data for noisy and product review texts. These texts are more informal in nature and ridden with both spelling and grammatical errors. There could be multiple possible reasons: a) LLMs struggle to translate the noisy texts, resulting in poor quality hypotheses and lower COMET score. b) COMET metric is calculated through embeddings. Here, the source side is noisy, which can lead to unreliable embeddings and, therefore, an unreliable COMET score.

#### 4.2.2 Model wise Overview

The best-performing models in terms of COMET scores are Online-B and TransmissionMT, closely followed by Claude-3.5 and Unbabel-Tower-70B. However, the worst-performing model is still CycleL.

### 4.3 Human Evaluation

The next evaluation method employed is human evaluation. We enlisted the expertise of a linguist in our lab, who randomly selected 3 sentences from each of the 6 domains. For each sentence, the corresponding machine translations from the 16 sub-

mitted model outputs were collected, resulting in 288 sentences. These sentences were then rated on a scale from 1 to 5, where 1 indicates the poorest translation, and 5 represents the best possible translation compared to the reference texts. Note that due to such a low number of samples, the results in manual evaluation are very unreliable. However, due to resource constraints, we could not perform a large-scale manual evaluation. Nonetheless, we hope this rating will provide some ideas about the competence of the models when observed along with scores from automated metrics.

#### 4.3.1 Domain wise Overview

According to the human evaluation, the general domain showed the highest faithfulness to the reference translations. This outcome is expected, as general domain texts are typically easier to translate due to their formal and unambiguous nature, with fewer grammatical, lexical, and spelling errors. Conversely, the noisy domain demonstrated the lowest faithfulness to the reference translations. This is largely attributed to the informal nature of these texts, which often include profanities and internet acronyms like "lol" and "idk" as well as a higher prevalence of errors.

#### 4.3.2 Model wise Overview

Almost consistent with the COMET metrics, we can see that the TransmissionMT, Unbabel-Tower-70B, and Claude-3.5 have the best human-evaluated scores, whereas CycleL again scored the least favorably.

## 5 Conclusion

This paper presents a comparison of various model submissions for the WMT Shared Task 2024. We proposed a dataset with domain-wise segregation and conducted a domain-specific analysis of the submitted models. Our comprehensive evaluation using BLEU, COMET, and human assessments of the machine-translated hypotheses identified Claude 3.5, TransmissionMT, Unbabel Tower 70B, Online-A, and Online-B as some of the top-performing models for machine translation using LLMs. The analysis revealed that the formal domains of general and education are the easiest for models to handle, whereas the noisy and review domains proved to be the most challenging. This study highlights that while LLMs show proficiency in machine translation, there is still significant room for improvement.

# References

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks.

Ramakrishna Appicharla, Asif Ekbal, and Pushpak Bhattacharyya. 2021. EduMT: Developing machine translation system for educational content in Indian languages. In *Proceedings of the 18th International Conference on Natural Language Processing (ICON)*, pages 35–43, National Institute of Technology Silchar, Silchar, India. NLP Association of India (NLPAI).

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2016. Neural machine translation by jointly learning to align and translate.

Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of bloom.

Maximiliana Behnke, Antonio Valerio, Antonio Valerio Miceli Barone, Rico Sennrich, Vilelmini Sosoni, Thanasis Naskos, Eirini Takoulidou, Maria Stasimioti, Menno Zaanen, Sheila Castilho, Federico Gaspari, Yota Georgakopoulou, Valia Kordoni, Markus Egg, and Katia Kermanidis. 2018. Improving machine translation of educational content via crowdsourcing.

Quinten Bolding, Baohao Liao, Brandon James Denis, Jun Luo, and Christof Monz. 2023. Ask language model to clean your noisy translation data.

Vicent Briva-Iglesias, Joao Lucas Cavalheiro Camargo, and Gokhan Dogru. 2024. Large language models "ad referendum": How good are they at machine translation in the legal domain?

Amartya Chowdhury, Deepak K. T., Samudra Vijaya K, and S. R. Mahadeva Prasanna. 2022. Machine translation for a very low-resource language - layer freezing approach on transfer learning. In *Proceedings of the Fifth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2022)*, pages 48–55, Gyeongju, Republic of Korea. Association for Computational Linguistics.

Sören DREANO, Derek MOLLOY, and Noel MURPHY. 2024. Cyclegn: a cycle consistent approach for neural machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Irina-Ana Drobot. 2023. Translating literature using machine translation: Is it really possible? *Scientific Bulletin of the Politehnica University of Timişoara Transactions on Modern Languages*, 20:57–64.

Denis Elshin, Nikolay Karpachev, Boris Gruzdev, Ilya Golovanov, Georgy Ivanov, Alexander Antonov, Nickolay Skachkov, Ekaterina Latypova, Vladimir Layner, Ekaterina Enikeeva, Dmitry Popov, Anton Chekashev, Vladislav Negodin, Vera Frantsuzova, Alexander Chernyshev, and Kirill Denisov. 2024. From general LLM to translation: How we dramatically improve translation quality using human evaluation data for LLM finetuning. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Abhimanyu Dubey et al. 2024a. The llama 3 herd of models.

OpenAI et al. 2024b. Gpt-4 technical report.

Baban Gain, Ramakrishna Appicharla, Soumya Chennabasavaraj, Nikesh Garera, Asif Ekbal, and Muthusamy Chelliah. 2022. Low resource chat translation: A benchmark for Hindi–English language pair. In *Proceedings of the 15th biennial conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 83–96, Orlando, USA. Association for Machine Translation in the Americas.

Baban Gain, Dibyanayan Bandyopadhyay, and Asif Ekbal. 2021a. Experiences of adapting multimodal machine translation techniques for Hindi. In *Proceedings of the First Workshop on Multimodal Machine Translation for Low Resource Languages (MMTLRL 2021)*, pages 40–44, Online (Virtual Mode). INCOMA Ltd.

Baban Gain, Dibyanayan Bandyopadhyay, and Asif Ekbal. 2021b. IITP at WAT 2021: System description for English-Hindi multimodal translation task. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 161–165, Online. Association for Computational Linguistics.

Baban Gain, Dibyanayan Bandyopadhyay, Samrat Mukherjee, Chandranath Adak, and Asif Ekbal. 2023. Impact of visual context on noisy multimodal nmt: An empirical study for english to indian languages.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Vakul Goyle, Parvathy Krishnaswamy, Kannan Girija Ravikumar, Utsa Chattopadhyay, and Kartikay Goyle. 2023. Neural machine translation for low resource languages.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection.

Kamal Gupta, Soumya Chennabasavaraj, Nikesh Garera, and Asif Ekbal. 2021a. Product review translation using phrase replacement and attention guided noise augmentation. In *Proceedings of Machine Translation Summit XVIII: Research Track*, pages 243–255, Virtual. Association for Machine Translation in the Americas.

Kamal Kumar Gupta, Soumya Chennabasavaraj, Nikesh Garera, and Asif Ekbal. 2021b. Product review translation: Parallel corpus creation and robustness towards user-generated noisy text. In *Proceedings of the 4th Workshop on e-Commerce and NLP*, pages 174–183, Online. Association for Computational Linguistics.

Kamal Kumar Gupta, Divya Kumari, Soumya Chennabasavaraj, Nikesh Garera, and Asif Ekbal. 2022. Reviewmt: Sentiment preserved e-commerce review translation system. In *Proceedings of the 5th Joint International Conference on Data Science & Management of Data (9th ACM IKDD CODS and 27th COMAD)*, CODS-COMAD '22, page 275–279, New York, NY, USA. Association for Computing Machinery.

Kshitij Gupta, Devansh Gautam, and Radhika Mamidi. 2021c. ViTA: Visual-linguistic translation by aligning object tags. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 166–173, Online. Association for Computational Linguistics.

Ali Hatami, Shubhanker Banerjee, Mihael Arcan, Bharathi Raja Chakravarthi, Paul Buitelaar, and John Philip McCrae. 2024. English-to-low-resource translation: A multimodal approach for hindi, malayalam, bengali, and hausa. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation.

Christian Herold, Jan Rosendahl, Joris Vanvinckenroye, and Hermann Ney. 2022. Detecting various types of noise for neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2542–2551, Dublin, Ireland. Association for Computational Linguistics.

Miroslav Hrabal, Josef Jon, Martin Popel, Nam Luu, Danil Semin, and Ondřej Bojar. 2024. CUNI at WMT24 general translation task: Llms, (q)lora, CPO and model merging. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine.

Huda Khayrallah and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation. In *NMT@ACL*.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Findings of the WMT24 general machine translation shared task:
the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Minato Kondo, Ryo Fukuda, Xiaotian Wang, Katsuki Chousa, Masato Nishimura, Kosei Buma, Takatomo Kano, and Takehito Utsuro. 2024. NTTSU at WMT2024 general translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Keito Kudo, Hiroyuki Deguchi, Makoto Morishita, Ryo Fujii, Takumi Ito, Shintaro Ozaki, Koki Natsumi, Kai Sato, Kazuki Yano, Ryosuke Takahashi, Subaru Kimura, Tomomasa Hara, Yusuke Sakai, and Jun Suzuki. 2024. Document-level translation with LLM reranking: Team-j at WMT 2024 general translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Anoop Kunchukuttan, Pratik Mehta, and Pushpak Bhattacharyya. 2018. The IIT Bombay English-Hindi parallel corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Taja Kuzman, Špela Vintar, and Mihael Arčan. 2019. Neural machine translation of literary texts from English to Slovene. In *Proceedings of the Qualities of Literary Machine Translation*, pages 1–9, Dublin, Ireland. European Association for Machine Translation.

Sahinur Rahman Laskar, Abdullah Faiz Ur Rahman Khilji, Darsh Kaushik, Partha Pakray, and Sivaji Bandyopadhyay. 2021. Improved English to Hindi multimodal neural machine translation. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 155–160, Online. Association for Computational Linguistics.

Baohao Liao, Christian Herold, Shahram Khadivi, and Christof Monz. 2024. IKUN for WMT24 general MT task: Llms are here for multilingual machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Rubén Martínez-Domínguez, Matīss Rikters, Artūrs Vasiļevskis, Mārcis Pinnis, and Paula Reichenberg. 2020. Customized neural machine translation systems for the Swiss legal domain. In *Proceedings of the 14th Conference of the Association for Machine Translation in the Americas (Volume 2: User Track)*, pages 217–223, Virtual. Association for Machine Translation in the Americas.

Evgeny Matusov. 2019. The challenges of using neural machine translation for literature. In *Proceedings of the Qualities of Literary Machine Translation*, pages 10–19, Dublin, Ireland. European Association for Machine Translation.

Paul Michel and Graham Neubig. 2018. MTNT: A testbed for machine translation of noisy text. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 543–553, Brussels, Belgium. Association for Computational Linguistics.

Yasmin Moslem, Rejwanul Haque, John Kelleher, and Andy Way. 2022. Domain-specific text generation for machine translation. In *Proceedings of the 15th biennial conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 14–30, Orlando, USA. Association for Machine Translation in the Americas.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Shabdapurush Poudel, Bal Krishna Bal, and Praveen Acharya. 2024. Bidirectional English-Nepali machine translation(MT) system for legal domain. In *Proceedings of the 3rd Annual Meeting of the Special Interest Group on Under-resourced Languages @ LREC-COLING 2024*, pages 53–58, Torino, Italia. ELRA and ICCL.

Anil Ramakrishna, Rahul Gupta, Jens Lehmann, and Morteza Ziyadi. 2023. INVITE: a testbed of automatically generated invalid questions to evaluate large language models for hallucinations. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5422–5429, Singapore. Association for Computational Linguistics.

Gowtham Ramesh, Sumanth Doddapaneni, Aravinth Bheemaraj, Mayank Jobanputra, Raghavan AK, Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Mahalakshmi J, Divyanshu Kakwani, Navneet Kumar, Aswin Pradeep, Srihari Nagaraj, Kumar Deepak, Vivek Raghavan, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh Shantadevi Khapra. 2022. Samanantar: The Largest Publicly Available Parallel Corpora Collection for 11 Indic Languages. *Transactions of the Association for Computational Linguistics*, 10:145–162.

Surangika Ranathunga, En-Shiun Annie Lee, Marjana Prifti Skenduli, Ravi Shekhar, Mehreen Alam, and Rishemjit Kaur. 2023. Neural machine translation for low-resource languages: A survey. *ACM Comput. Surv.*, 55(11).

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Nathaniel R. Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. Chatgpt mt: Competitive for high- (but not low-) resource languages.

Shaomu Tan, David Stap, Seth Aycock, Christof Monz, and Di Wu. 2024. Uva-MT's participation in the WMT24 general translation shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Katherine Thai, Marzena Karpinska, Kalpesh Krishna, Bill Ray, Moira Inghilleri, John Wieting, and Mohit Iyyer. 2022. Exploring document-level literary machine translation with parallel paragraphs from world literature. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9882–9902, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey. European Language Resources Association (ELRA).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

BigScience Workshop, :, and Teven Le Scao et al. 2023. Bloom: A 176b-parameter open-access multilingual language model.

Zhanglin Wu, Daimeng Wei, Zongyao Li, Hengchao Shang, Jiaxin GUO, Shaojun Li, Zhiqiang Rao, Yuanchang Luo, Ning Xie, and Hao Yang. 2024. Choose the final translation from NMT and LLM hypotheses using MBR decoding: HW-TSC's submission to the WMT24 general MT shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Zeynep Yirmibeşoğlu, Olgun Dursun, Harun Dallı, Mehmet Şahin, Ena Hodzik, Sabri Gürses, and Tunga Güngör. 2023. Incorporating human translator style into english-turkish literary machine translation.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: a case study. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Wenbo Zhang. 2024. IOL research machine translation systems for WMT24 general machine translation shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Zebiao Zhou, Xiangxun Zhu, Xiaowei Ji, Li Yang, Fengjie Zhu, and Tuanwei Shi. 2024. Hyper-SNMT at translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis.

Hao Zong, Chao Bei, Huan Liu, Conghu Yuan, Wentao Chen, and Degen Huang. 2024. DLUT and GTCOM's neural machine translation systems for WMT24. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

## A  Overall Scores

We report the overall BLEU and COMET scores in Figure 13 and Figure 14. Further, we provide the domain-wise and model-wise average rating by human annotators in Figure 15 and Figure 16.

## B  Participants

The WMT24 General Translation Task showcased diverse approaches to machine translation. Several teams explored the potential of Large Language Models (LLMs) for translation tasks. IKUN demonstrated the effectiveness of LLMs in multilingual translation, achieving top rankings in multiple language directions (Liao et al., 2024). The IOL Research team leveraged LLMs for continued pretraining and synthetic data generation (Zhang, 2024).

Some teams focused on improving existing neural machine translation (NMT) architectures. HW-TSC combined NMT and LLM-based models using Minimum Bayesian Risk (MBR) decoding (Wu et al., 2024). UvA-MT compared fine-tuned LLMs with traditional encoder-decoder NMT systems (Tan et al., 2024). The DLUT and GTCOM team emphasized back-translation and multilingual models (Zong et al., 2024).

Novel approaches were also presented. CycleGN introduced a cycle-consistent approach for non-parallel datasets (DREANO et al., 2024). Hyper-SNMT proposed embedding sentences in hyperbolic space to better capture language hierarchies (Zhou et al., 2024).

Several teams explored domain-specific adaptations. Team-J incorporated document-level LLM reranking for improved context-aware translations (Kudo et al., 2024). NTTSU focused on speech domain translation for Japanese to Chinese (Kondo et al., 2024).

The Yandex team demonstrated significant improvements using human evaluation data for LLM fine-tuning (Elshin et al., 2024). CUNI explored various techniques including QLoRA, CPO, and model merging (Hrabal et al., 2024).

Multimodal approaches were also explored, with researchers integrating visual information to enhance translation for low-resource languages (Hatami et al., 2024).

These diverse approaches highlight the ongoing innovation in machine translation, with a notable trend towards leveraging LLMs and exploring novel architectures to improve translation quality across various language pairs and domains.

## C  Dataset Statistics

Here, we have shared the summary statistics of the lengths of different sentences in each domain. Further we have also shared the harmonic mean of ratio of source to reference text sentence in each domain. From this graph it is evident that general domain has the most disparity in terms of source and reference sentence length. Also, it has the longest sentences compared to the other domains.

## D  Dataset Example

In Table Table 5, we present examples from the religious domain. This table showcases various outputs relevant to religious texts, highlighting key themes and interpretations.

Table ?? provides examples from the judicial domain. The *Online-B* model has the best quality of translation. The output from the model *Nvidia_Nemo* and *IKUN_C* is inadequate. The original text conveys a universal message about divine provision and the consequences of human actions, while the translation introduces specificity, making it feel more direct and personal.

Figure 7: BLEU scores in the Education Domain



Figure 8: BLEU scores in the General Domain



Figure 9: BLEU scores in the Judicial Domain



Figure 10: BLEU scores in the Literature Domain



Figure 11: BLEU scores in the Noisy Domain



Figure 12: BLEU scores in the Product Review Domain

| file | min_nword | max_nword | average_nword |
|------|-----------|-----------|---------------|
| education_source.txt | 5 | 79 | 25 |
| education_reference.txt | 5 | 80 | 29 |
| general_source.txt | 16 | 222 | 29 |
| general_reference.txt | 5 | 195 | 30 |
| judicial_source.txt | 11 | 39 | 21 |
| judicial_reference.txt | 9 | 56 | 24 |
| literature_source.txt | 11 | 38 | 21 |
| literature_reference.txt | 9 | 63 | 24 |
| noisy_source.txt | 21 | 49 | 31 |
| noisy_reference.txt | 20 | 74 | 38 |
| review_source.txt | 11 | 48 | 21 |
| review_reference.txt | 9 | 59 | 25 |

Table 4: Statistics of the domain-wise files



Figure 13: BLEU Score on the Full Dataset



Figure 15: Domain-wise Average Human Score



Figure 14: COMET Score on the Full Dataset



Figure 16: Model-wise Average Human Score

Figure 17: Maximum Minimum and Average Words in File

| Model/Text | Model Output |
|---|---|
| Source | And made the cloud spread shade over you, and sent for you manna and quails that you may eat of the good things We have made for you. No harm was done to Us, they only harmed themselves. |
| Online-B | और तुम्हारे ऊपर बादल छाया कर दिया और तुम्हारे लिए मन्ना और बटेर भेजे, ताकि जो अच्छी-अच्छी चीज़ें हमने तुम्हारे लिए बनाई हैं, उनमें से तुम खाओ। हमें कोई नुक़सान नहीं पहुँचाया गया, बल्कि उन्होंने तो बस अपने आप को नुक़सान पहुँचाया। |
| Nvidia_Nemo | और हमने तुम्हारे ऊपर बादल की छाया फैला दी और तुम्हारे लिए मन्ना और बटेरें भेजी ताकि तुम उन भलाईयों में से खाओ जो हमने तुम्हारे लिए बनाई हैं और हम पर कोई हानि नहीं हुई, वे अपने आप को ही हानि पहुँचाते थे। |
| IKUN_C | और बादल को तुम पर छाया डालने के लिए भेजा, और माना और बाज को भेजा कि तुम हमारी ओर से तुम्हारे लिए भेजी हुई अच्छी चीजों से खा लो। |
| Reference | और हमने तुमपर बादलों की छाया की और तुमपर 'मन्न' और 'सलबा' उतारा - खाओ, जो अच्छी पाक चीजें हमने तुम्हें प्रदान की है उन्होंने हमारा तो कुछ भी नहीं बिगाड़ा, बल्कि वे अपने ही ऊपर अत्याचार करते रहे |

Table 5: Example of Religious domain output

# Investigating the Linguistic Performance of Large Language Models in Machine Translation

**Shushen Manakhimova[1], Vivien Macketanz[1], Eleftherios Avramidis[1],
Ekaterina Lapshinova-Koltunski[2], Sergei Bagdasarov[3] and Sebastian Möller[1]**

[1]German Research Center for Artificial Intelligence (DFKI)
`firstname.lastname@dfki.de`
[2]University of Hildesheim, `lapshinovakoltun@uni-hildesheim.de`
[3]Saarland University, `sergeiba@lst.uni-saarland.de`

## Abstract

This paper summarizes the results of our test suite evaluation on 39 machine translation systems submitted at the Shared Task of the Ninth Conference of Machine Translation (WMT24). It offers a fine-grained linguistic evaluation of machine translation outputs for English–German and English–Russian, resulting from a significant manual linguistic effort. Based on our results, LLMs are inferior to NMT in English–German, both in overall scores and when translating specific linguistic phenomena, such as punctuation, complex future verb tenses, and stripping. LLMs show quite a competitive performance in English-Russian, although top-performing systems might struggle with some cases of named entities and terminology, function words, mediopassive voice, and semantic roles. Additionally, some LLMs generate very verbose or empty outputs, posing challenges to the evaluation process.

## 1 Introduction

The evolution of large language models (LLMs) has revived interest in machine translation (MT) evaluation, raising the discussion about whether general-purpose LLMs can outperform specialized MT systems. LLMs have demonstrated remarkable performance across various tasks, prompting an urgent need to assess their linguistic capabilities and potential risks (Wang et al., 2024; Guerreiro et al., 2023). Last year's Eighth Conference on Machine Translation findings (WMT23; Kocmi et al., 2023) showed that one LLM performed well across most language pairs. Although GPT-4 excelled in some areas (e.g., translation of user-generated content), it struggled with other aspects, such as speaker gender translation and specific domains (e.g., legal); it ranked lower than encoder-decoder systems when translating from English into less-represented languages, e.g., Czech and Russian.. However, last year's General MT Task included only two LLM-based system submissions (Kocmi

et al., 2023). This year marks a noteworthy increase in LLMs participating in the task. As a result, this paper covers a linguistically motivated evaluation of a broad range of LLMs, including Claude-3.5-Sonnet, GPT-4, Llama3-70B, Mistral-Large, and the recently released Unbabel-Tower70B, as well as CUNI-DS, IKUN and IKUN-C, IOL-Research, CommandR-plus, Yandex, and Occiglot.

In this context, we are presenting the results of our participation in the test suite sub-task of the Ninth Conference on Machine Translation (WMT24). Our test suite[1] consists of carefully crafted sentences that assess the ability of MT systems to handle specific linguistic phenomena. It was applied to the MT systems submitted for evaluation in two language directions: English–German and English–Russian.

## 2 Related Work

Several researchers have adopted test suites or challenge sets to better identify flaws in MT outputs, further contributing to the advancement in MT evaluation. The WMT test suite sub-task has played a significant role by providing a platform for these evaluations.

Chen et al. (2023), for example, developed a systematic method of selecting difficult sentences from the Wiki Corpus, taking into account factors like word difficulty, sentence length, grammatical complexity, and model learnability. Their findings showed significant differences from the official ranking, suggesting that systems performing well on average test sets might not do as well on more challenging ones. Notably, GPT-4 ranked among the top two for Chinese–English translations and between fourth and ninth in the other direction. Other research has focused on difficulties posed by special domains and writing styles. Mukherjee and Shrivastava (2023) designed a test suite for

---

[1]`https://github.com/DFKI-NLP/mt-testsuite`

English–German translation across five domains and writing styles. They found that while GPT-4 performed competitively overall, it struggled in the legal domain and with the judgment writing style. The work of Savoldi et al. (2023) looked into gender translation of the English–German and German–English language directions. They found that while systems generally handled gender form translation well, producing gender-inclusive translations still remains a significant challenge. Specifically, GPT-4 exhibited relatively lower accuracy in accurately translating feminine gender in first-person singular references reflecting the speaker's linguistic expression of gender.

Bawden and Sagot (2023) tested the ability of MT systems, including GPT-4, to handle user-generated text from in-domain sources characterized by informal language and various grammatical deviations. Their findings show that although data at such a large scale can provide extensive training data, GPT-4 still does not perform well on consistency and faithfulness to source sentences, implying a hurdle for generalization to out-of-domain text.

The fact that these works indicated weaknesses not apparent on the General MT Shared Task illustrates the critical importance of developing focused test suites beyond general evaluation metrics to measure the capabilities and limitations of MT systems.

## 3 Method

### 3.1 Test suite description

We have developed a fine-grained test suite to evaluate the performance of MT systems for the language pairs English–German and English–Russian[2]. While we are only touching on the description of our test suite in the paper, the interested reader can find a detailed description in Macketanz et al. (2022a). Previous submissions of the test suite in WMT can be found in (Macketanz et al., 2018, 2021, 2022b; Avramidis et al., 2019, 2020; Manakhimova et al., 2023).

Our test suite focuses on various linguistic phenomena that are of interest to the respective language pairs. The phenomena are based on extensive research in linguistics, contrastive grammars, and translation studies, covering a wide range of po-

---

| Language Pair | Test Items | Categories | Phenomena |
|---|---|---|---|
| en–de | 4,846 | 13 | 110 |
| en–ru | 1,234 | 12 | 51 |

Table 1: Metadata of the language pairs in the test suite.

tential translation challenges. The phenomena and their categorization are specific to a language pair and a language direction; however, there is a big overlap of the phenomena between the language pairs for the languages covered so far.

The phenomena in the test suite are classified into several categories, grouped by the underlying syntactical/morphological/lexical mechanisms. Each phenomenon is represented by at least 20 (in many cases more) test items. Every test item consists of one or more sentence(s) in the source language and a set of rules to evaluate them. The test items are either handwritten by linguistic experts or taken from existing corpora. The number of test items, phenomena, and categories per language pair can be seen in Table 1. While the English–German test suite has been around and growing since 2017, the English–Russian test suite is newer (from 2022) and, therefore, has fewer test items.

With the change of MT system types over the years (from phrase-based and statistical MT to neural MT, and finally to LLMs), typical MT challenges and errors have also changed. Thus, we have also adapted our test suite over the past few years to accommodate those changes. These adaptations included adding new phenomena, longer/more complex test sentences, and more test items per phenomenon.

MT outputs evaluated by the test suite have been used to produce challenge sets for WMT metrics (Avramidis and Macketanz, 2022; Avramidis et al., 2023).

### 3.2 Application of the test suite

The test suite can be characterized as semi-automatic, as the evaluation process is based on automatic rules and additional manual evaluation. While this kind of evaluation can be more time-consuming than a fully automatic evaluation process, we assume it to be more accurate as the regular expressions are handwritten by human experts.

For each test item in the test suite, one or more linguist(s) have written regular expressions to cover as many as possible expected correct and incorrect translations. The linguists rely on their years of

experience in evaluating MT systems when writing regular expressions. However, of course, not all MT outputs can be covered by the regular expressions as languages, and the MT systems are very diverse. In these cases, the human comes into play again. All outputs that cannot be automatically evaluated by the regular expressions are inspected and hand-evaluated by a linguist.The more unexpected (meaning, in most cases, incorrect) outputs a system creates, the more manual work is involved in the evaluation process. After the evaluation process, the translation accuracy of an MT system specific per phenomenon or category is calculated by dividing the number of correctly translated test items by the total number of test items.

To ensure a fair comparison, only evaluated test items are considered for accuracy calculations. If a test item is not evaluated for one system, it is excluded for all systems, reducing the number of the effective test items.

For the system comparison (per language direction), we first identify the highest-scoring system and then compare it to the other systems. The significance of the comparison is confirmed by a one-tailed Z-test with $\alpha = 0.95$. Systems that do not perform significantly inferior to the best-performing system are grouped into the first performance cluster. The best-performing systems are indicated in boldface in the respective rows of the tables.

To account for variations in the number of test items within each category or phenomenon, average scores are computed in three different methods: The *micro-average* method combines the contributions of all test items to calculate the average percentages. In the *category macro-average*, the percentages are first computed independently per category and subsequently averaged, treating all categories equally. Analogously, for the *phenomenon macro-average* the percentages are computed independently per phenomenon and averaged afterwards, treating all phenomena equally.

## 4 Experiment Setup

This year, we evaluated a total of 39 systems with our test suite. The systems had been submitted to the General MT Shared Tasl of the Ninth Conference on Machine Translation. 21 systems were evaluated for English–German and 18 systems for English–Russian[3].

It is the fourth time we evaluated the English-German systems and the third time for the English-Russian systems. As described above, the evaluation of the system outputs is only semi-automatic, and therefore, manual work is needed to complement the automatic evaluation by resolving cases in which none of the rules in our rule database can be applied, the so-called *warnings*. Upon receiving the system outputs, there were on average around 25 % of warnings for English–German, varying across systems from 4.7 % to 77.5 %. For English–Russian, there were on average 46.9 %, ranging from 24.5 % to 82.7 %. As we had added several new phenomena and test items to existing phenomena before this year's WMT, we expected more warnings this year. Additionally, several systems this year, particularly LLMs, were more verbose or "creative" with their translations than we are used to from previous years. For example, Mistral sometimes offered several translation options, including explanations. This creativity, however, led to more manual work as the existing evaluation rules could not cover these unexpected outputs.

This year, the manual evaluation was conducted by three linguists who were experts in one or both language pairs. Combined, the linguists spent around 160 person-hours on the manual evaluation within about three weeks. After the manual input, an average of 0.9 % of warnings remained for English–German and 5.7 % for English–Russian.

As mentioned above, test items with one unresolved warning for at least one system were excluded from the comparison. This reduced the number of effective test items to 4219 ($\sim$87 %) test items for English–German, and 994 ($\sim$80 %) for English–Russian.

## 5 Results

All result tables can be found in the Appendix.

### 5.1 System comparison

For **English–German**, Online-B, TranssionMT, and Claude-3.5 had the highest micro-average with a score of around 97 %. Furthermore, Online-B and TranssionMT also had the highest macro-average, with a scrore of around 95 %. Whereas little is

---

[3]There had originally been 25 systems submitted for En-De, and 22 for En-Ru. However, the systems Dubformer and CycleL/CycleL2 had to be excluded from our evaluation for both language pairs due to invalid output

known about Online-B, TranssionMT's good performance may be explained by its optimization for complex grammatical structures and rich morphology through the use of a hyperbolic embedding. The lowest micro-average was reached by TSU-HITs (Mynka and Mikhaylovskiy, 2024) with a score of 38.6 %, and the lowest macro-average was reached by MSLC (Larkin et al., 2024) with a score of 45.8 %. On average, systems reached a micro-average of 81.4 % and a macro-average of 79.7 %.

At this point, it is important to note that two systems, Mistral (Jiang et al., 2023) and Occiglot (Avramidis et al., 2024), produced a (high) number of empty outputs for German–English. While Occiglot only generated 335 empty sentences, Mistral generated as many as 3,624. For the system comparison, we had to mark these sentences as incorrect. Therefore, Mistral appears to have the worst accuracy on micro- and macro-average. However, we conducted an extra analysis for these two systems, only considering the correct and incorrect outputs and ignoring the empty outputs. This resulted in an accuracy of 73.0 % for Occiglot and 85.9 % for Mistral macro-averaged over the non-empty outputs (see Tables in Section A). Since the accuracies are calculated over different test items, they are not comparable with each other and with other systems.

Interestingly, and contrary to previous years, our ranking of the systems according to their linguistically-related performance differs from the preliminary results of the automatic ranking of the General MT Shared Task (Kocmi et al., 2024): While the top 3 systems in the General task were Unbabel-Tower70B (Rei et al., 2024), Dubformer, and TranssionMT, according to our analysis, Online-B and TranssionMT made it to the first significance cluster, with GPT-4 falling in the second position and Unbabel-Tower70B scoring even lower. Furthermore, we had to exclude Dubformer from our analysis due to invalid output. Nonetheless, both analyses have MSLC and TSU-HITs at the bottom of the ranking.

When comparing the human rankings of the General MT Task with our rankings, one can note that in the former, many systems share the cluster of the first position. The fact that our test suite can produce a smaller significance cluster for the first position can be considered a success.

While Unbabel-Tower70B showed exceptional performance across all language directions in the automatic preliminary rankings, our evaluation revealed some potential blind spots. Compared to the top-performing systems, it struggles with less commonly used future tenses (ditransitive—future II progressive, ditransitive—future II simple, reflexive—future II progressive), with the elliptic process of *stripping*, and with *semantic roles*. Future II progressive tense can pose difficulties, likely due to its infrequent occurrence in training data and nuanced nature. An example sentence would be "I will have been baking Tim a cake." Stripping will be explained in further detail below, cf. Sec 5.3 As for semantic roles, English is relatively flexible in assigning semantic roles to subjects. In contrast, German tends to have stricter rules for subject roles regarding agentivity. This difference can cause translation issues when models directly map English constructions onto German without considering these syntactic and semantic differences.

For **English–Russian**, Yandex (Elshin et al., 2024) and Claude-3.5-Sonnet achieved the highest micro-average scores with 91.8 % and 90.4 %, respectively as well as macro-averages with 92.4 % and 90.5 %. This year's poorest-performing system was TSU-HITs, with both micro- and macro-averages of 50 and 49 %. On average, the systems reached a micro-average of 80.1 % and a macro-average of 78.7 %. According to the automatic preliminary results, the top four best-performing systems for English–Russian are Unbabel-Tower70B, Dubformer, Yandex, and Claude-3.5-Sonnet, in that order. As mentioned earlier, Dubformer was excluded from our analysis. Unbabel-Tower70B scores slightly lower than Yandex and Claude-3.5-Sonnet, achieving 89.4 % micro-average and a 90 % macro-average. On the phenomenon level, our evaluation shows that Yandex and Claude-3.5-Sonnet outperform Unbabel-Tower70B, when it comes to *collocations*, *onomatopeia*, *verb valency*, and *passive voice*. If we exclude Cycle and CycleL (Dreano et al., 2024), the worst four performing systems, according to the automatic preliminary ranking, are the same four systems in our ranking, listed here from best to worst: IKUN-C (Liao et al., 2024), CUNI-DS (Semin and Bojar, 2024), NVIDIA-NeMo, and TSU-HITs. GPT-4, one of the best-performing systems last year, falls into the second cluster this year.

| Stripping | |
|---|---|
| John can play the guitar, and Mary too. | |
| John kann Gitarre spielen und Mary auch. | pass |
| John kann Gitarre spielen, und Mary auch. | fail |
| John kann das instrument spielen, und Lucia noch nicht. | fail |
| **Verb Semantics** | |
| "I've missed you so much!" he bawled. | |
| "Ich habe dich so sehr vermisst!" schluchzte er. | pass |
| "Ich habe dich so vermisst!" schrie er. | pass |
| »Ich habe dich so sehr verpasst!«, bawte er. | fail |

Table 2: Examples of English–German linguistic phenomena with passing and failing MT outputs.

| Compound | |
|---|---|
| The police officer was pregnant. | |
| Сотрудница полиции была беременна. | pass |
| Полицейский был беременна. | fail |
| У полицейской была беременность. | fail |
| **Verb Semantics** | |
| She described the book as a page-turner. | |
| Она описала книгу как -захватывающую историю. | pass |
| она описала книгу как страницу-поворотчик. | fail |
| Она описала книгу как перелистывание страниц. | fail |

Table 3: Examples of English–Russian linguistic phenomena with passing and failing MT outputs.

## 5.2 Category-level analysis

For **English–German**, two systems are in the cluster of best-performing systems per category in all categories: Online-B and TranssionMT. Furthermore, two systems have the highest accuracies on all but two/three categories, namely GPT-4 and Claude-3.5-Sonnet. The categories with the highest accuracies are *negation*, with 14 systems reaching 100 % accuracy, and *subordination*.

Some of the easiest categories for **English–Russian** include *subordination* (89.7 %), *function words* (89.1 %), where both LLM-based and other MT system score over 95 %. In contrast, *ambiguity* stood out as the most challenging, with an accuracy average of 69.2 % along such categories as *false friends* and *multi-word expressions*, with average accuracies of 70.7 % and 69.5 %, respectively. These indicate more challenges on the lexical rather than the syntactical level.

## 5.3 Phenomenon-level analysis

For **English–German**, the phenomenon-level macro average is 80 %, which is similar to the category-level macro average and the general micro-average. The phenomena with the highest accuracies (> 90 %) are *negative inversion*, *prepositional MWE*, *date*, *substitution*, *adverbial clause*, *infinitive clause*, and *intransitive future I progressive/simple*; for a detailed overview cf. Table 10.

On the other hand, the phenomena with the lowest accuracies (< 65 %) that a lot of LLM-based model struggled with are *stripping*, *idiom*, *onomatopoeia*, *ditransitive future II progressive/simple*, *reflexive future II progressive/simple*, *transitive future II progressive*, and *semantic roles*. It seems that the future II progressive/simple tense is particularly difficult for systems to translate, no matter the verb type. As mentioned above, this is likely due to this verb tense's uncommonness.

Table 2 contains translation examples from English–German. The first example is a test item for the phenomenon *stripping*. *Stripping* is a type of ellipsis. While stripping exists in both German and English, one aspect that can lead to translation errors is punctuation. In English, there is a comma between the two constituents ("John can play the guitar" and "and Mary too"). In German, however, placing a comma in between the constituents is incorrect; see the first and second translation examples. The third translation contains more errors than the additional comma, as it completely changes the meaning of the second constituent. This translation was produced by the Cycle system and also showcases how these kinds of "creative" translations lead to more manual evaluation work: It is easy to write a regular expression for the incorrect output with the comma before the second constituent, and this regular expression will cover most of the outputs of the incorrect system as this is a very common error. However, it is impossible to predict such an incorrect output as it was produced by Cycle, and therefore, it is impossible to write a regular expression to cover cases like this.

The second example is from the phenomenon of *verb semantics*. This phenomenon refers to semantic components in the verb's semantic structure that do not have formal markers. Some examples of these kinds of verbs are *to stride*, *to rumble*, *to stagger*, or *to bawl*, like in the example at hand. There are usually several correct translations for these verbs, as seen in the Table. However, they might lead to translation errors, with systems sometimes not translating them (because they are not so common) or translating them with an incorrect semantic meaning.

When evaluating the performance across phenomena for **English–Russian**, it was found that the following phenomena posed minimal challenges, with many systems achieving near-perfect accuracy: *catenative verb*, *case government*, *conditional*, *contact clause*, *object clauses*, *personal pronoun coreference*, *prepositional mwe*, and *date*. Notably, *personal pronoun coreference*, a new phenomenon added last year that focuses on the consistency in translating the formal and informal "you" across sentences, as well as ensuring that a past tense "I" retains the correct gender ending. This category attained a remarkable accuracy of 96.8 %, marking a 13 % improvement compared to last year. The phenomena with the lowest accuracies (< 60 %) are *verbal MWE*, *resultative*, *gapping*, *compounds*, *idioms* and *semantic roles*.

Table 3 contains translation examples from English–Russian. The original sentences and their translations have been shortened for the paper. The first example involves the common English compound "police officer". Despite the simplicity of this sentence, a closer examination reveals various issues in the translations. In English, the nominal phrase in question is gender-neutral, with no gender marking on nouns, adjectives, or verbs. However, in contrast to English, Russian has gender marking not only on pronouns but also on other parts of speech, including nouns, adjectives, verbs, determiners, and numbers. The first translation correctly uses the collocation сотрудница полиции (literally, "female police employee") and appropriately pairs it with была беременна ("was pregnant"), both in the feminine form. This nominal phrase construction is necessary to convey the gender within the translation. In the second translation, the word полицейский, typically referring to a male police officer is then followed by the verb был (the masculine form of was), and later by the adjective pregnant in the feminine form. This translation error was produced by an LLM and highlights a gap in the model's understanding of gender agreement rules in Russian and a lack of real-world knowledge. The third translation renders the phrase as "the female police officer had pregnancy," which is not a linguistically acceptable Russian collocation. It also uses the adjective полицейская as a job title, which is not a standard noun for "police officer" in Russian.

The next example comes from the phenomenon of Noun formation with the suffix -er. This process is a part of derivational morphology, where new words are formed by adding affixes to existing words or changing their grammatical category or meaning. This is a highly fruitful suffix in English. In the first example translation, we see it rendered as захватывающую историю or "captivating story." This transformation effectively captures the essence of "page-turner." The second translation has страницу-поворотчик – a literal translation. Перелистывание страниц in the third translation describes the physical action of turning pages. The first translation is accurate as it captures the idiomatic meaning of "page-turner"; the other two translations fail due to overly literal interpretations, a common issue in encoder-decoder models and LLMs.

## 5.4 Comparison with previous years

We have analyzed some of the best-performing systems' development over the years for systems submitted to the WMT repeatedly in the past years. For **English–German**, we took a closer look at GPT-4, Online-B, Online-W, and Online-A, see Table 8. GPT-4 has seen barely any changes in the accuracy from 2023 to 2024 (although it needs to be noted that the prompting method has changed from 5-shot to 3-shot). Online-B, however, shows an improvement of 2.5 percentage points on the macro-average from 2021 to 2024, while the micro-average stayed almost the same throughout this period. Online-W, similarly to GPT-4, shows almost no changes from 2021 to 2024. And finally, Online-A has slight improvements of 1 and 3 percentage points from 2021 to 2024 on the micro-average and macro-average level, respectively.

While in the past years, Online-B and Online-W were usually in the cluster of best-performing systems together, this year, Online-B has surpassed Online-W as only the former is in the best-performing cluster as of this year, while the latter is not. Furthermore, in 2023, GPT-4, Online-W, and Online-B were together in the cluster of best performing systems, while this year, GPT-4 is also not in that cluster anymore.

As for the scored of the micro-average, the phenomenon macro-average, and the category macro-average, while the first two have almost not changed from 2023 to 2024, the category macro-average has improved about 2.5 percentage points from last year to this year. This suggests that the systems for English—German have undergone a

slight improvement compared to last year.

Table 9 compares the performance of Yandex, GPT-4, and Online-G for **English–Russian** from 2022 to 2024. This year's Yandex submission is a trained YandexGPT, an LLM-based model. Their approach includes extensive pre-training, fine-tuning, p-tuning, and structure-preserving techniques, which help ensure contextually accurate translations (Elshin et al., 2024). Over the last two years, Yandex's submission has likely undergone a significant update, as reflected in the 2.59 % accuracy increase. Overall, Yandex shows consistent performance with some improvement. GPT-4, another LLM, demonstrates a generally strong performance compared to last year, with a significant drop in the punctuation category (from 100 % to 60 %). Despite this, GPT-4 has either improved or maintained stable performance across most linguistic categories. Online-G, as we suspect based on encoder-decoder methods, exhibits stable performance without any substantial improvements in any areas.

### 5.5 LLMs vs. encoder-decoder NMT

NMT systems based on an encoder-decoder (or commercial systems that we assume they use this technology) still exhibit better linguistic performance than LLMs in English–German, whereas in English–Russian the first position is shared indeed by two LLMs. In English-German, LLMs seem to perform worse than the two best-performing NMT systems, regarding *punctuation, future verb tenses* and *stripping*. For English-Russian, Yandex is weaker in *named entities and terminology*, while Claude struggles with *function words*, and Unbabel with *verb valency* that includes error-prone phenomena for all LLMs, such as *semantic roles*, *verb semantics*, *resultative*, and *mediopassive voice*. GPT-4 scores even lower than several commercial NMT-based systems. This suggests that while LLMs are indeed taking over the MT in fine-grained analysis, some still struggle to match the capabilities of specialized NMT systems, which are tailored specifically to the target language and potentially trained on more language-specific data.

### 6 Conclusions and Outlook

In this paper, we apply a linguistically motivated test suite for the first time to evaluate the translation performance of several LLMs as well as several systems with different architectures. Based on the macro-averaged accuracies, the best systems for English-German are Online-B and TranssionMT, with Claude-3.5-Sonnet also sharing the first position based on micro-averaging. For English–Russian, the best-performing systems are Yandex and Claude-3.5-Sonnet. While LLMs generally perform strongly in MT, systems based on encoder-decoder methods, such as TranssionMT and most probably Online-B may still have an edge in certain areas. What the human evaluations of the main MT task reveal about the systems is still to be determined, pending the official announcement of the rankings. The results underscore the potential of LLMs in MT but also highlight areas for improvement.

### Limitations

The current test suite was initially designed to evaluate earlier MT systems, featuring a wide range of linguistic phenomena without challenging the models. However, it is becoming increasingly clear that we need to adapt and potentially eliminate the phenomena that have proven too easy for the systems in recent years. The significance of the averaging is unclear, and adding weights depending on the importance of various phenomena is something to consider. While we have introduced context in some cases and complexity with multi-sentence test items in others, this has not been done for all phenomena and sentences so far. One challenge is that we often encounter correct rendering of the phenomena, but then encounter grave errors in the sentence structure. Internally, it has been concluded that these sentences should be marked as incorrect, as the errors are often too significant for the whole output to be considered correct. Additionally, this year, some models generated responses that resembled those of a classical chatbot, including additional explanations or commentary that mixed correct and incorrect translations, making it challenging to evaluate. Going forward, we plan to further refine the test suite to better capture the nuances of modern translation systems.

### Acknowledgements

prior contributions to the creation of the test suite.

# References

Eleftherios Avramidis, Annika Grützner-Zahn, Manuel Brack, Patrick Schramowski, Pedro Ortiz Suarez, Malte Ostendorff, Fabio Barth, Shushen Manakhimova, Vivien Macketanz, Georg Rehm, and Kristian Kersting. 2024. Occiglot at WMT24: European open-source large language models evaluated on translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Eleftherios Avramidis and Vivien Macketanz. 2022. Linguistically motivated evaluation of machine translation metrics based on a challenge set. In *Proceedings of the Seventh Conference on Machine Translation*, pages 514–529, Abu Dhabi. Association for Computational Linguistics.

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, Aljoscha Burchardt, and Sebastian Möller. 2020. Fine-grained linguistic evaluation for state-of-the-art Machine Translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 346–356, Online. Association for Computational Linguistics.

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, and Hans Uszkoreit. 2019. Linguistic Evaluation of German-English Machine Translation Using a Test Suite. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 445–454, Florence, Italy. Association for Computational Linguistics.

Eleftherios Avramidis, Shushen Manakhimova, Vivien Macketanz, and Sebastian Möller. 2023. Challenging the state-of-the-art machine translation metrics from a linguistic perspective. In *Proceedings of the Eighth Conference on Machine Translation*, pages 713–729, Singapore. Association for Computational Linguistics.

Rachel Bawden and Benoît Sagot. 2023. RoCS-MT: Robustness challenge set for machine translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 198–216, Singapore. Association for Computational Linguistics.

Xiaoyu Chen, Daimeng Wei, Zhanglin Wu, Ting Zhu, Hengchao Shang, Zongyao Li, Jiaxin Guo, Ning Xie, Lizhi Lei, Hao Yang, and Yanfei Jiang. 2023. Multifaceted challenge set for evaluating machine translation performance. In *Proceedings of the Eighth Conference on Machine Translation*, pages 217–223, Singapore. Association for Computational Linguistics.

Sören Dreano, Derek Molloy, and Noel Murphy. 2024. Cyclegn: a cycle consistent approach for neural machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Denis Elshin, Nikolay Karpachev, Boris Gruzdev, Ilya Golovanov, Georgy Ivanov, Alexander Antonov, Nickolay Skachkov, Ekaterina Latypova, Vladimir Layner, Ekaterina Enikeeva, Dmitry Popov, Anton Chekashev, Vladislav Negodin, Vera Frantsuzova, Alexander Chernyshev, and Kirill Denisov. 2024. From general LLM to translation: How we dramatically improve translation quality using human evaluation data for LLM finetuning. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Nuno M. Guerreiro, Duarte M. Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André F. T. Martins. 2023. Hallucinations in large multilingual translation models. *Transactions of the Association for Computational Linguistics*, 11:1500–1517.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Preliminary WMT24 Ranking of General MT Systems and LLMs.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Samuel Larkin, Chi-kiu Lo, and Rebecca Knowles. 2024. MSLC24 submissions to the general machine translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Baohao Liao, Christian Herold, Shahram Khadivi, and Christof Monz. 2024. IKUN for WMT24 general MT task: Llms are here for multilingual machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, and Hans Uszkoreit. 2018. Fine-grained evaluation of German-English machine translation based on a test suite. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 578–587, Belgium, Brussels. Association for Computational Linguistics.

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, and Hans Uszkoreit. 2022a. A linguistically motivated test suite to semi-automatically evaluate German–English machine translation output. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 936–947, Marseille, France. European Language Resources Association.

Vivien Macketanz, Eleftherios Avramidis, Shushen Manakhimova, and Sebastian Möller. 2021. Linguistic Evaluation for the 2021 State-of-the-art Machine Translation Systems for German to English and English to German. In *Proceedings of the Sixth Conference on Machine Translation*, pages 1059–1073, Online. Association for Computational Linguistics.

Vivien Macketanz, Shushen Manakhimova, Eleftherios Avramidis, Ekaterina Lapshinova-koltunski, Sergei Bagdasarov, and Sebastian Möller. 2022b. Linguistically motivated evaluation of the 2022 state-of-the-art machine translation systems for three language directions. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 432–449, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Shushen Manakhimova, Eleftherios Avramidis, Vivien Macketanz, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2023. Linguistically motivated evaluation of the 2023 state-of-the-art machine translation: Can ChatGPT outperform NMT? In *Proceedings of the Eighth Conference on Machine Translation*, pages 224–245, Singapore. Association for Computational Linguistics.

Ananya Mukherjee and Manish Shrivastava. 2023. IIIT HYD's submission for WMT23 test-suite task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 246–251, Singapore. Association for Computational Linguistics.

Vladimir Aleksandrovich Mynka and Nikolay Mikhaylovskiy. 2024. TSU HITS's submissions to the WMT 2024 general machine translation shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Ricardo Rei, Jose Maria Pombal, Nuno M. Guerreiro, João Alves, Pedro Henrique Martins, Patrick Fernandes, Helena Wu, Tania Vaz, Duarte Alves, Amin Farajian, Sweta Agrawal, Antonio Farinhas, José G. C. de Souza, and André Martins. 2024. Tower v2: Unbabel-IST 2023 submission for the general MT shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Beatrice Savoldi, Marco Gaido, Matteo Negri, and Luisa Bentivogli. 2023. Test suites task: Evaluation of gender fairness in MT with MuST-SHE and INES. In *Proceedings of the Eighth Conference on Machine Translation*, pages 252–262, Singapore. Association for Computational Linguistics.

Danil Semin and Ondřej Bojar. 2024. CUNI-DS submission: A naive transfer learning setup for english-to-russian translation utilizing english-to-czech data. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Jiaqi Wang, Hanqi Jiang, Yiheng Liu, Chong Ma, Xu Zhang, Yi Pan, Mengyuan Liu, Peiran Gu, Sichen Xia, Wenjun Li, Yutong Zhang, Zihao Wu, Zhengliang Liu, Tianyang Zhong, Bao Ge, Tuo Zhang, Ning Qiang, Xintao Hu, Xi Jiang, Xin Zhang, Wei Zhang, Dinggang Shen, Tianming Liu, and Shu Zhang. 2024. A comprehensive review of multimodal large language models: Performance and challenges across different tasks.

# A Separate systems

| category | items | acc. |
|---|---|---|
| Ambiguity | 5 | 100.0 |
| Coordination & ellipsis | 35 | 82.9 |
| False friends | 12 | 100.0 |
| Function word | 15 | 100.0 |
| LDD & interrogatives | 54 | 98.1 |
| Lexical Morphology | 17 | 100.0 |
| MWE | 31 | 90.3 |
| Named entity & terminology | 39 | 94.9 |
| Negation | 4 | 75.0 |
| Non-verbal agreement | 31 | 93.5 |
| Punctuation | 8 | 75.0 |
| Subordination | 51 | 94.1 |
| Verb semantics | 4 | 0.0 |
| Verb tense/aspect/mood | 875 | 96.9 |
| Verb valency | 31 | 87.1 |
| micro-average | 1212 | 95.5 |
| macro-average | 1212 | 85.9 |

Table 4: Accuracies for the translations of the Mistral-Large system (en-de) considering only the non-empty outputs

| category | items | acc. |
|---|---|---|
| Ambiguity | 22 | 86.4 |
| Coordination & ellipsis | 124 | 60.5 |
| False friends | 40 | 92.5 |
| Function word | 40 | 75.0 |
| LDD & interrogatives | 207 | 76.3 |
| Lexical Morphology | 39 | 61.5 |
| MWE | 123 | 76.4 |
| Named entity & terminology | 112 | 77.7 |
| Negation | 18 | 66.7 |
| Non-verbal agreement | 109 | 87.2 |
| Punctuation | 37 | 51.4 |
| Subordination | 191 | 85.3 |
| Verb semantics | 23 | 60.9 |
| Verb tense/aspect/mood | 3249 | 71.9 |
| Verb valency | 114 | 65.8 |
| micro-average | 4448 | 72.8 |
| macro-average | 4448 | 73.0 |

Table 5: Accuracies for the translations of the Occiglot system (en-de) considering only the non-empty outputs

# B Analysis based on categories

| categ | count | Onl-B | Trans | GPT4 | Claud | Unbab | Comma | Onl-W | Llama | Aya23 | IOLRe | Onl-A | Onl-G | IKUN | CUNIN | IKUNC | NVIDI | Occig | AISTA | TSUHI | MSLC | Mistr | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | 24 | 100.0 | 100.0 | 95.8 | 100.0 | 95.8 | 100.0 | 95.8 | 91.7 | 91.7 | 87.5 | 91.7 | 87.5 | 79.2 | 87.5 | 83.3 | 70.8 | 79.2 | 50.0 | 37.5 | 41.7 | 20.8 | 80.4 |
| Coordination & ellipsis | 83 | 94.0 | 94.0 | 90.4 | 71.1 | 73.5 | 74.7 | 74.7 | 80.7 | 79.5 | 65.1 | 73.5 | 84.3 | 62.7 | 78.3 | 69.9 | 66.3 | 57.8 | 61.4 | 42.2 | 26.5 | 27.7 | 69.0 |
| False friends | 40 | 95.0 | 95.0 | 97.5 | 95.0 | 92.5 | 95.0 | 90.0 | 95.0 | 95.0 | 95.0 | 87.5 | 82.5 | 97.5 | 95.0 | 90.0 | 77.5 | 92.5 | 70.0 | 55.0 | 45.0 | 30.0 | 84.2 |
| Function word | 42 | 97.6 | 97.6 | 97.6 | 97.6 | 97.6 | 95.2 | 100.0 | 95.2 | 92.9 | 92.9 | 97.6 | 97.6 | 88.1 | 88.1 | 92.9 | 92.9 | 66.7 | 88.1 | 61.9 | 47.6 | 33.3 | 86.6 |
| LDD & interrogatives | 160 | 97.5 | 97.5 | 96.3 | 97.5 | 96.9 | 93.8 | 95.6 | 92.5 | 91.9 | 90.0 | 95.6 | 93.1 | 80.6 | 85.6 | 80.6 | 85.6 | 71.9 | 75.0 | 57.5 | 48.8 | 25.0 | 83.3 |
| Lexical Morphology | 27 | 92.6 | 92.6 | 100.0 | 96.3 | 92.6 | 85.2 | 85.2 | 88.9 | 77.8 | 90.0 | 81.5 | 81.5 | 74.1 | 66.7 | 66.7 | 63.0 | 63.0 | 33.3 | 33.3 | 18.5 | 48.1 | 73.0 |
| MWE | 109 | 97.2 | 96.3 | 91.7 | 97.2 | 89.0 | 93.6 | 93.6 | 83.5 | 90.8 | 90.8 | 86.2 | 85.9 | 83.5 | 78.0 | 74.3 | 66.1 | 73.4 | 62.4 | 40.4 | 33.0 | 22.0 | 77.6 |
| Named entity & terminology | 92 | 92.4 | 92.4 | 95.7 | 94.6 | 89.1 | 96.7 | 89.1 | 93.5 | 89.1 | 90.2 | 90.2 | 85.9 | 81.5 | 76.1 | 81.5 | 80.4 | 77.2 | 77.2 | 65.2 | 59.8 | 33.7 | 82.5 |
| Negation | 19 | 100.0 | 100.0 | 100.0 | 100.0 | 94.7 | 89.5 | 94.7 | 100.0 | 84.2 | 100.0 | 100.0 | 100.0 | 94.7 | 89.5 | 100.0 | 100.0 | 63.2 | 89.5 | 89.5 | 89.5 | 15.8 | 90.2 |
| Non-verbal agreement | 97 | 97.9 | 97.9 | 99.0 | 100.0 | 99.0 | 93.8 | 100.0 | 88.2 | 100.0 | 93.8 | 91.8 | 91.8 | 90.7 | 73.5 | 85.6 | 79.4 | 84.5 | 78.4 | 78.4 | 53.6 | 23.7 | 86.2 |
| Punctuation | 34 | 100.0 | 100.0 | 88.2 | 85.3 | 94.1 | 91.2 | 97.1 | 88.2 | 100.0 | 94.1 | 88.2 | 88.2 | 88.2 | 73.5 | 88.2 | 88.2 | 50.0 | 82.4 | 61.8 | 67.6 | 14.7 | 82.6 |
| Subordination | 148 | 98.0 | 98.0 | 98.0 | 99.3 | 96.6 | 94.6 | 96.6 | 95.9 | 93.2 | 97.3 | 96.6 | 96.6 | 89.2 | 95.9 | 85.8 | 91.9 | 81.1 | 87.8 | 67.6 | 66.9 | 23.0 | 88.1 |
| Verb semantics | 18 | 83.3 | 83.3 | 72.2 | 72.2 | 88.9 | 77.8 | 66.7 | 77.8 | 83.3 | 72.2 | 50.0 | 61.1 | 72.2 | 66.7 | 55.6 | 44.4 | 55.6 | 22.2 | 16.7 | 11.1 | 0.0 | 58.7 |
| Verb tense/aspect/mood | 3225 | 98.2 | 98.2 | 97.6 | 98.4 | 98.9 | 96.6 | 96.6 | 93.6 | 96.2 | 94.9 | 97.1 | 98.7 | 77.0 | 80.7 | 77.4 | 82.3 | 67.1 | 72.5 | 72.5 | 42.2 | 23.8 | 81.6 |
| Verb valency | 101 | 91.1 | 91.1 | 86.1 | 88.1 | 88.1 | 84.2 | 86.1 | 78.2 | 86.1 | 84.2 | 83.2 | 76.2 | 71.3 | 72.3 | 75.2 | 64.4 | 63.4 | 54.5 | 34.7 | 42.2 | 16.8 | 71.9 |
| micro-average | 4219 | 97.7 | 97.7 | 96.8 | 97.3 | 91.5 | 95.3 | 95.3 | 92.6 | 94.7 | 93.5 | 95.3 | 96.3 | 78.1 | 81.4 | 81.4 | 81.3 | 68.4 | 72.1 | 38.6 | 43.4 | 24.0 | 81.4 |
| macro-average | 4219 | 95.7 | 95.6 | 93.7 | 92.8 | 92.0 | 90.8 | 90.3 | 89.9 | 89.6 | 89.1 | 87.8 | 87.4 | 82.0 | 81.8 | 80.5 | 76.9 | 69.8 | 67.2 | 51.7 | 45.8 | 23.9 | 79.7 |

Table 6: Accuracies (%) of successful translations on the categorylevel for English–German. The boldface indicates the significantly best-performing systems per row.

| categ | count | Yande | Claud | Unbab | Comma | Onl-G | Onl-W | GPT4 | IOLRe | Trans | Onl-B | Onl-A | Aya23 | IKUN | Llama | IKUNC | CUNID | NVIDI | TSUHI | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | 20 | 90.0 | 95.0 | 90.0 | 90.0 | 70.0 | 70.0 | 90.0 | 90.0 | 50.0 | 50.0 | 55.0 | 85.0 | 70.0 | 75.0 | 60.0 | 65.0 | 35.0 | 15.0 | 69.2 |
| Coordination & ellipsis | 86 | 87.2 | 80.2 | 84.9 | 82.6 | 84.9 | 83.7 | 76.7 | 75.6 | 72.1 | 72.1 | 72.1 | 77.9 | 74.4 | 69.8 | 65.1 | 72.1 | 54.7 | 47.7 | 74.1 |
| False friends | 15 | 86.7 | 86.7 | 86.7 | 86.7 | 86.7 | 73.3 | 66.7 | 66.7 | 66.7 | 66.7 | 80.0 | 66.7 | 66.7 | 66.7 | 66.7 | 46.7 | 66.7 | 53.3 | 70.7 |
| Function word | 34 | 97.1 | 88.2 | 94.1 | 100.0 | 94.1 | 100.0 | 97.1 | 91.2 | 94.1 | 94.1 | 88.2 | 94.1 | 85.3 | 85.3 | 82.4 | 73.5 | 73.5 | 66.7 | 89.1 |
| LDD & interrogatives | 81 | 97.5 | 93.8 | 97.5 | 91.4 | 96.3 | 95.1 | 91.4 | 90.1 | 91.4 | 91.4 | 85.2 | 86.4 | 82.7 | 80.2 | 76.5 | 82.7 | 70.4 | 59.3 | 86.6 |
| Lexical Morphology | 41 | 97.6 | 92.7 | 90.2 | 92.7 | 82.9 | 71.9 | 80.5 | 75.6 | 75.6 | 75.6 | 70.7 | 68.3 | 75.6 | 75.6 | 63.4 | 53.7 | 34.1 | 26.8 | 72.4 |
| MWE | 96 | 87.5 | 84.4 | 78.1 | 83.3 | 80.2 | 77.1 | 77.1 | 76.0 | 72.9 | 71.9 | 70.8 | 67.7 | 69.8 | 66.7 | 66.7 | 52.1 | 40.6 | 33.3 | 69.5 |
| Named entity & terminology | 80 | 83.8 | 95.0 | 87.5 | 81.3 | 80.0 | 80.0 | 81.3 | 73.8 | 80.0 | 80.0 | 77.5 | 71.3 | 62.5 | 77.5 | 60.0 | 57.5 | 56.3 | 41.3 | 73.7 |
| Non-verbal agreement | 98 | 94.9 | 95.9 | 91.8 | 93.9 | 90.8 | 89.8 | 90.8 | 92.9 | 80.6 | 80.6 | 82.9 | 92.9 | 83.7 | 86.7 | 85.7 | 81.6 | 73.5 | 65.3 | 86.2 |
| Punctuation | 13 | 92.3 | 92.3 | 92.3 | 100.0 | 92.3 | 76.9 | 61.5 | 76.9 | 84.6 | 84.6 | 92.3 | 84.6 | 84.6 | 61.5 | 86.1 | 100.0 | 80.0 | 92.3 | 85.0 |
| Subordination | 115 | 98.3 | 94.8 | 98.3 | 88.7 | 95.7 | 96.5 | 94.8 | 93.0 | 85.0 | 85.0 | 86.1 | 86.1 | 86.1 | 88.7 | 83.5 | 83.5 | 80.9 | 67.0 | 85.0 |
| Verb semantics | 20 | 100.0 | 90.0 | 95.0 | 70.0 | 95.0 | 85.0 | 65.0 | 80.0 | 85.0 | 85.0 | 65.0 | 65.0 | 80.0 | 75.0 | 70.0 | 55.0 | 30.0 | 35.0 | 74.7 |
| Verb tense/aspect/mood | 169 | 87.0 | 90.5 | 90.5 | 87.0 | 89.3 | 88.8 | 89.9 | 85.2 | 85.2 | 85.2 | 86.4 | 85.2 | 81.1 | 84.6 | 84.6 | 84.6 | 78.1 | 45.0 | 82.9 |
| Verb valency | 126 | 93.7 | 88.1 | 83.3 | 81.0 | 84.9 | 86.5 | 81.7 | 81.0 | 83.3 | 83.3 | 75.4 | 77.8 | 78.6 | 76.2 | 73.8 | 68.3 | 58.7 | 50.0 | 78.1 |
| micro-average | 994 | 91.8 | 90.4 | 89.4 | 86.8 | 87.8 | 86.1 | 85.2 | 83.3 | 82.3 | 82.2 | 80.7 | 80.4 | 78.1 | 79.0 | 78.1 | 75.1 | 62.2 | 50.0 | 80.1 |
| macro-average | 994 | 92.4 | 90.5 | 90.0 | 87.7 | 87.4 | 83.6 | 81.7 | 81.3 | 79.7 | 79.7 | 79.6 | 79.0 | 77.2 | 76.4 | 78.6 | 72.6 | 59.7 | 49.0 | 78.7 |

Table 7: Accuracies (%) of successful translations on the category-level for English–Russian. The boldface indicates the significantly best-performing systems per row.

# C  Yearly comparison

| Category | items | GPT4 | | onlineB | | | | OnlineW | | | | onlineA | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 2023 | 2024 | 2021 | 2022 | 2023 | 2024 | 2021 | 2022 | 2023 | 2024 | 2021 | 2022 | 2023 | 2024 |
| Ambiguity | 24 | 95.8 | 95.8 | 91.7 | 91.7 | 91.7 | 100 | 95.8 | 95.8 | 95.8 | 95.8 | 91.7 | 87.5 | 87.5 | 91.7 |
| Coordination & ellipsis | 88 | 88.6 | 87.5 | 80.7 | 87.5 | 8.8 | 89.8 | 70.5 | 70.5 | 73.9 | 73.9 | 71.6 | 80.7 | 79.5 | 72.7 |
| False friends | 38 | 97.4 | 97.4 | 84.2 | 89.5 | 89.5 | 94.7 | 89.5 | 92.1 | 89.5 | 89.5 | 86.8 | 86.8 | 86.8 | 89.5 |
| Function word | 41 | 100 | 97.6 | 100 | 97.6 | 97.6 | 95.1 | 100 | 100 | 97.6 | 97.6 | 97.6 | 97.6 | 97.6 | 97.6 |
| MWE | 104 | 97.1 | 95.2 | 94.2 | 95.2 | 95.2 | 98.1 | 96.2 | 97.1 | 97.1 | 97.1 | 86.5 | 90.4 | 92.3 | 94.2 |
| Named entity & termin. | 85 | 95.3 | 97.6 | 92.9 | 97.6 | 94.1 | 98.8 | 95.3 | 92.9 | 94.1 | 94.1 | 94.1 | 94.1 | 92.9 | 94.1 |
| Negation | 18 | 100 | 94.4 | 94.4 | 100 | 94.4 | 100 | 100 | 100 | 94.4 | 94.4 | 94.4 | 100 | 100 | 100 |
| Non-verbal agreement | 67 | 100 | 100 | 95.5 | 95.5 | 95.5 | 100 | 95.5 | 97 | 95.5 | 95.5 | 94.4 | 95.5 | 100 | 97 |
| Punctuation | 36 | 86.1 | 83.3 | 83.3 | 83.3 | 83.3 | 100 | 97.2 | 94.4 | 97.2 | 97.2 | 95.5 | 97.2 | 88.9 | 91.7 |
| Subordination | 163 | 98.8 | 97.5 | 97.5 | 98.8 | 98.2 | 95.7 | 96.9 | 96.3 | 96.6 | 96.9 | 97.2 | 98.2 | 98.8 | 98.8 |
| Verb tense & aspect/mood | 3076 | 97.9 | 97.9 | 99 | 98.7 | 97.9 | 98.1 | 96.5 | 96.3 | 96.6 | 96.6 | 96.1 | 98.5 | 98.3 | 97.2 |
| Verb valency | 89 | 87.6 | 92.1 | 86.5 | 87.6 | 91 | 95.5 | 86.5 | 86.5 | 88.8 | 88.8 | 84.3 | 87.6 | 91 | 92.1 |
| micro-avg | 3829 | 97.3 | 97.3 | 97.5 | 97.7 | 97.1 | 97.8 | 95.6 | 95.4 | 95.8 | 95.8 | 95 | 97.2 | 97.2 | 96.3 |
| macro-avg | 3829 | 95.4 | 94.7 | 91.7 | 93.6 | 93.2 | 97.2 | 93.3 | 93.3 | 93.1 | 93.1 | 91.2 | 92.8 | 92.4 | 93.1 |

Table 8: Yearly comparison of the systems of WMT24 for English-German, based on the category-level analysis

| category | items | Yandex | | GPT4 | | OnlineG | | |
|---|---|---|---|---|---|---|---|---|
| year | | 2022 | 2024 | 2023 | 2024 | 2022 | 2023 | 2024 |
| Ambiguity | 7 | 85.7 | 85.7 | 100 | 100 | 85.7 | 85.7 | 85.7 |
| Coordination & ellipsis | 30 | 80 | 70 | 80 | 76.7 | 80 | 80 | 80 |
| False friends | 5 | 80 | 100 | 100 | 100 | 80 | 80 | 80 |
| Function word | 10 | 90 | 90 | 80 | 90 | 90 | 90 | 90 |
| MWE | 32 | 71.9 | 96.9 | 75 | 84.4 | 71.9 | 71.9 | 78.1 |
| Named entity & terminology | 21 | 90.5 | 90.5 | 76.2 | 76.2 | 90.5 | 95.2 | 85.7 |
| Non-verbal agreement | 11 | 81.8 | 90.9 | 63.6 | 72.7 | 81.8 | 81.8 | 90.9 |
| Punctuation | 5 | 100 | 80 | 100 | 60 | 100 | 100 | 100 |
| Subordination | 28 | 89.3 | 92.9 | 82.1 | 85.7 | 89.3 | 89.3 | 89.3 |
| Verb tense/aspect/mood | 67 | 74.6 | 67.2 | 77.6 | 88.1 | 74.6 | 82.1 | 79.1 |
| Verb valency | 26 | 84.6 | 96.2 | 84.6 | 84.6 | 84.6 | 80.8 | 80.8 |
| micro-avg | 242 | 81 | 83.1 | 79.8 | 83.9 | 81 | 83.1 | 82.6 |
| macro-avg | 242 | 84.4 | 87.3 | 83.6 | 83.5 | 84.4 | 85.2 | 85.4 |

Table 9: Yearly comparison of the systems of WMT24 for English-Russian, based on the category-level analysis

## D   Detailed analysis on a phenomenon-level

| categ | count | Onl-B | Trans | GPT4 | Claud | Unbab | Comma | Onl-W | Llama | Aya23 | IOLRe | Onl-A | Onl-G | IKUN | CUNIN | IKUNC | NVIDI | Occig | AISTA | TSUHI | MSLC | Mistr | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | 24 | 100.0 | 100.0 | 95.8 | 100.0 | 95.8 | 100.0 | 95.8 | 91.7 | 91.7 | 91.7 | 87.5 | 91.7 | 87.5 | 79.2 | 83.3 | 70.8 | 79.2 | 50.0 | 37.5 | 41.7 | 20.8 | 80.4 |
| Lexical ambiguity | 24 | 100.0 | 100.0 | 95.8 | 100.0 | 95.8 | 100.0 | 95.8 | 91.7 | 91.7 | 91.7 | 87.5 | 91.7 | 87.5 | 79.2 | 83.3 | 70.8 | 79.2 | 50.0 | 37.5 | 41.7 | 20.8 | 80.4 |
| Coordination & ellipsis | 83 | 94.0 | 94.0 | 90.4 | 71.1 | 73.5 | 74.7 | 74.7 | 80.7 | 79.5 | 65.1 | 73.5 | 84.3 | 62.7 | 78.3 | 69.9 | 66.3 | 57.8 | 61.4 | 42.2 | 26.5 | 27.7 | 69.0 |
| Gapping | 15 | 86.7 | 93.3 | 93.3 | 66.7 | 60.0 | 80.0 | 73.3 | 73.3 | 86.7 | 46.7 | 86.7 | 86.7 | 53.3 | 93.3 | 86.7 | 93.3 | 33.3 | 73.3 | 33.3 | 26.7 | 40.0 | 70.2 |
| Pseudogapping | 7 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 71.4 | 71.4 | 85.7 | 100.0 | 71.4 | 57.1 | 100.0 | 71.4 | 100.0 | 100.0 | 42.9 | 100.0 | 42.9 | 42.9 | 0.0 | 0.0 | 72.1 |
| Right node raising | 11 | 90.9 | 90.9 | 90.9 | 81.8 | 81.8 | 100.0 | 90.9 | 100.0 | 100.0 | 100.0 | 81.8 | 85.7 | 72.7 | 90.9 | 72.7 | 81.8 | 81.8 | 81.8 | 27.3 | 45.5 | 27.3 | 80.5 |
| Sluicing | 18 | 100.0 | 100.0 | 94.4 | 94.4 | 77.8 | 77.8 | 100.0 | 94.4 | 88.9 | 100.0 | 77.8 | 77.8 | 77.8 | 66.7 | 72.2 | 72.2 | 33.3 | 77.8 | 55.6 | 33.3 | 27.8 | 79.4 |
| Stripping | 21 | 90.5 | 81.0 | 81.0 | 52.4 | 38.1 | 47.6 | 42.9 | 57.1 | 47.6 | 38.1 | 42.9 | 85.7 | 42.9 | 61.9 | 47.6 | 42.9 | 33.3 | 47.6 | 42.9 | 23.8 | 23.8 | 51.5 |
| VP-ellipsis | 11 | 100.0 | 100.0 | 81.8 | 81.8 | 72.7 | 90.9 | 81.8 | 90.9 | 90.9 | 63.6 | 72.7 | 90.9 | 72.7 | 81.8 | 54.5 | 63.6 | 36.4 | 36.4 | 45.5 | 18.2 | 36.4 | 70.1 |
| False friends | 40 | 95.0 | 95.0 | 97.5 | 95.0 | 92.5 | 90.0 | 90.0 | 95.0 | 95.0 | 95.0 | 87.5 | 82.5 | 97.5 | 95.0 | 77.5 | 77.5 | 92.5 | 70.0 | 55.0 | 45.0 | 30.0 | 84.2 |
| Function word | 42 | 97.6 | 97.6 | 97.6 | 97.6 | 97.6 | 95.2 | 100.0 | 92.9 | 92.9 | 92.9 | 88.1 | 97.6 | 88.1 | 88.1 | 92.9 | 92.9 | 66.7 | 88.1 | 61.9 | 47.6 | 33.3 | 86.6 |
| Focus particle | 23 | 95.7 | 95.7 | 95.7 | 95.7 | 95.7 | 100.0 | 95.7 | 95.7 | 87.0 | 95.7 | 91.3 | 95.7 | 95.7 | 91.3 | 87.0 | 91.3 | 73.9 | 95.7 | 82.6 | 87.0 | 39.1 | 89.9 |
| Question tag | 19 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 94.7 | 100.0 | 94.7 | 89.5 | 89.5 | 100.0 | 95.7 | 78.9 | 84.2 | 87.0 | 94.7 | 73.9 | 78.9 | 36.8 | 0.0 | 26.3 | 82.7 |
| LDD & interrogatives | 160 | 97.5 | 97.5 | 96.3 | 96.9 | 93.8 | 95.6 | 92.5 | 91.9 | 90.0 | 90.0 | 85.6 | 93.1 | 80.6 | 85.6 | 80.6 | 85.6 | 71.9 | 75.0 | 57.5 | 48.8 | 25.0 | 83.3 |
| Extraposition | 16 | 93.8 | 87.5 | 87.5 | 81.3 | 81.3 | 87.5 | 81.3 | 75.0 | 56.3 | 56.3 | 62.5 | 62.5 | 50.0 | 75.0 | 62.5 | 50.0 | 43.8 | 37.5 | 37.5 | 12.5 | 12.5 | 67.0 |
| Inversion | 14 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 85.7 | 100.0 | 92.9 | 92.9 | 85.7 | 85.7 | 85.7 | 71.4 | 28.6 | 42.9 | 0.0 | 85.0 |
| Multiple connectors | 16 | 93.8 | 93.8 | 93.8 | 100.0 | 100.0 | 100.0 | 93.8 | 100.0 | 100.0 | 100.0 | 81.3 | 100.0 | 75.0 | 87.5 | 81.3 | 81.3 | 81.3 | 93.8 | 81.3 | 87.5 | 31.3 | 89.9 |
| Negative inversion | 14 | 92.9 | 92.9 | 92.9 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 92.9 | 92.9 | 100.0 | 92.9 | 78.6 | 100.0 | 92.9 | 85.7 | 21.4 | 90.8 |
| Pied-piping | 14 | 100.0 | 100.0 | 85.7 | 85.7 | 85.7 | 92.9 | 85.7 | 92.9 | 92.9 | 92.9 | 71.4 | 100.0 | 78.6 | 71.4 | 71.4 | 100.0 | 71.4 | 78.6 | 57.1 | 57.1 | 21.4 | 84.7 |
| Polar question | 18 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 94.4 | 94.4 | 94.4 | 77.8 | 100.0 | 72.2 | 77.8 | 77.8 | 94.4 | 72.2 | 72.2 | 66.7 | 38.9 | 33.3 | 86.8 |
| Preposition stranding | 16 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 81.3 | 100.0 | 81.3 | 75.0 | 75.0 | 56.3 | 75.0 | 62.5 | 56.3 | 75.0 | 75.0 | 75.0 | 75.0 | 43.8 | 0.0 | 56.3 | 77.7 |
| Split infinitive | 10 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 80.0 | 100.0 | 80.0 | 90.0 | 80.0 | 80.0 | 90.0 | 80.0 | 40.0 | 80.0 | 30.0 | 89.0 |
| Topicalization | 12 | 100.0 | 91.7 | 100.0 | 91.7 | 83.3 | 83.3 | 83.3 | 91.7 | 91.7 | 91.7 | 66.7 | 83.3 | 75.0 | 83.3 | 50.0 | 50.0 | 83.3 | 41.7 | 50.0 | 25.0 | 8.3 | 74.2 |
| Wh-movement | 30 | 96.7 | 96.7 | 100.0 | 96.7 | 93.3 | 93.3 | 93.3 | 93.3 | 93.3 | 93.3 | 96.7 | 96.7 | 90.0 | 96.7 | 83.3 | 86.7 | 70.0 | 83.3 | 63.3 | 60.0 | 26.7 | 86.0 |
| Lexical Morphology | 27 | 92.6 | 92.6 | 96.3 | 92.6 | 85.2 | 85.2 | 82.6 | 85.7 | 77.8 | 88.9 | 81.5 | 81.5 | 74.1 | 66.7 | 66.7 | 63.0 | 63.0 | 37.0 | 33.3 | 18.5 | 48.1 | 73.0 |
| Functional shift | 14 | 92.9 | 92.9 | 100.0 | 100.0 | 85.7 | 92.9 | 85.7 | 85.7 | 78.6 | 85.7 | 57.1 | 78.6 | 57.1 | 57.1 | 85.7 | 71.4 | 57.1 | 35.7 | 35.7 | 21.4 | 64.3 | 76.2 |
| Noun formation (er) | 13 | 92.3 | 92.3 | 100.0 | 84.6 | 84.6 | 76.9 | 76.9 | 69.2 | 69.2 | 92.3 | 76.9 | 84.6 | 53.8 | 46.2 | 76.9 | 53.8 | 69.2 | 38.5 | 30.8 | 15.4 | 30.8 | 69.6 |
| MWE | 109 | 97.2 | 96.3 | 91.7 | 97.2 | 89.0 | 93.6 | 93.6 | 83.5 | 90.8 | 90.8 | 78.0 | 86.2 | 83.5 | 74.3 | 66.1 | 66.1 | 73.4 | 62.4 | 40.4 | 33.0 | 22.0 | 77.6 |
| Collocation | 17 | 100.0 | 94.1 | 100.0 | 88.2 | 100.0 | 88.2 | 100.0 | 88.2 | 94.1 | 94.1 | 76.5 | 82.4 | 81.3 | 76.5 | 58.8 | 76.5 | 70.6 | 41.2 | 42.9 | 35.3 | 35.3 | 81.8 |
| Compound | 12 | 100.0 | 100.0 | 100.0 | 91.7 | 100.0 | 100.0 | 100.0 | 100.0 | 83.3 | 100.0 | 100.0 | 100.0 | 66.7 | 100.0 | 66.7 | 100.0 | 66.7 | 91.7 | 66.7 | 58.3 | 8.3 | 88.9 |
| Idiom | 19 | 89.5 | 63.2 | 89.5 | 68.4 | 73.7 | 73.7 | 63.2 | 68.4 | 68.4 | 68.4 | 47.4 | 31.6 | 57.9 | 47.4 | 0.0 | 0.0 | 31.6 | 0.0 | 0.0 | 0.0 | 5.3 | 47.6 |
| Nominal MWE | 18 | 100.0 | 94.4 | 100.0 | 88.9 | 100.0 | 94.4 | 94.4 | 77.8 | 83.3 | 88.9 | 77.8 | 94.4 | 83.3 | 61.1 | 77.8 | 77.8 | 77.8 | 66.7 | 33.3 | 27.8 | 16.7 | 78.3 |
| Prepositional MWE | 16 | 100.0 | 93.8 | 100.0 | 93.8 | 93.8 | 93.8 | 93.8 | 100.0 | 100.0 | 100.0 | 93.8 | 87.5 | 87.5 | 93.8 | 93.8 | 81.3 | 81.3 | 87.5 | 87.5 | 56.3 | 37.5 | 90.8 |
| Verbal MWE | 27 | 100.0 | 96.3 | 96.3 | 100.0 | 100.0 | 96.3 | 96.3 | 92.6 | 92.6 | 92.6 | 77.8 | 92.6 | 92.6 | 88.9 | 66.7 | 66.7 | 85.2 | 63.0 | 33.3 | 33.3 | 25.9 | 82.7 |
| Named entity & terminology | 92 | 92.4 | 94.6 | 95.7 | 94.6 | 89.1 | 96.7 | 89.1 | 93.5 | 89.1 | 90.2 | 90.2 | 85.9 | 81.5 | 76.1 | 81.5 | 80.4 | 77.2 | 84.2 | 65.2 | 59.8 | 33.7 | 82.5 |
| Date | 13 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 92.3 | 100.0 | 92.3 | 92.3 | 92.3 | 92.3 | 61.5 | 76.9 | 92.3 | 30.8 | 30.8 | 90.1 |
| Domainspecific Term | 16 | 87.5 | 87.5 | 87.5 | 87.5 | 93.8 | 87.5 | 87.5 | 81.3 | 68.8 | 68.8 | 68.8 | 75.0 | 81.3 | 81.3 | 81.3 | 81.3 | 61.5 | 87.5 | 31.3 | 37.5 | 43.8 | 76.2 |
| Location | 18 | 100.0 | 100.0 | 100.0 | 94.4 | 100.0 | 94.4 | 94.4 | 94.4 | 100.0 | 100.0 | 77.8 | 83.3 | 83.3 | 88.9 | 88.9 | 88.9 | 72.2 | 83.3 | 83.3 | 42.1 | 38.9 | 89.9 |
| Measuring unit | 19 | 94.7 | 94.7 | 94.7 | 100.0 | 89.5 | 100.0 | 89.5 | 100.0 | 94.7 | 94.7 | 94.7 | 94.7 | 84.2 | 89.5 | 84.2 | 84.2 | 73.7 | 94.7 | 73.7 | 42.1 | 21.1 | 85.5 |
| Onomatopeia | 7 | 57.1 | 85.7 | 85.7 | 42.9 | 42.9 | 42.9 | 42.9 | 57.1 | 42.9 | 42.9 | 42.9 | 42.9 | 42.9 | 14.3 | 28.6 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 39.5 |
| Proper name | 19 | 94.7 | 94.7 | 100.0 | 89.5 | 89.5 | 89.5 | 94.7 | 84.2 | 100.0 | 100.0 | 84.2 | 84.2 | 94.7 | 84.2 | 89.5 | 89.5 | 89.5 | 89.5 | 84.2 | 73.7 | 47.4 | 88.2 |
| Negation | 19 | 100.0 | 100.0 | 100.0 | 94.7 | 89.5 | 94.7 | 94.7 | 84.2 | 100.0 | 100.0 | 94.7 | 100.0 | 94.7 | 89.5 | 100.0 | 100.0 | 63.2 | 89.5 | 89.5 | 89.5 | 15.8 | 90.2 |
| Non-verbal agreement | 97 | 97.9 | 99.0 | 99.0 | 94.7 | 92.8 | 92.8 | 92.8 | 93.8 | 92.8 | 93.8 | 90.7 | 91.8 | 91.8 | 85.6 | 92.8 | 79.4 | 84.5 | 78.4 | 53.6 | 23.7 | 15.8 | 86.2 |
| Coreference | 30 | 96.7 | 93.3 | 93.3 | 89.5 | 93.3 | 96.7 | 93.3 | 96.7 | 90.0 | 90.0 | 90.0 | 93.3 | 90.0 | 80.0 | 93.3 | 80.0 | 90.0 | 70.0 | 90.0 | 43.3 | 16.7 | 87.3 |
| Genitive | 19 | 100.0 | 100.0 | 100.0 | 89.5 | 89.5 | 94.7 | 89.5 | 84.2 | 84.2 | 84.2 | 78.9 | 94.7 | 78.9 | 68.4 | 84.2 | 52.6 | 78.9 | 73.7 | 57.9 | 31.6 | 31.6 | 80.7 |
| Personal Pronoun Coreference | 12 | 91.7 | 91.7 | 91.7 | 91.7 | 91.7 | 83.3 | 91.7 | 92.3 | 66.7 | 66.7 | 66.7 | 66.7 | 92.3 | 76.9 | 91.7 | 75.0 | 83.3 | 83.3 | 75.0 | 58.3 | 33.3 | 83.7 |
| Possession | 26 | 100.0 | 100.0 | 96.2 | 96.2 | 96.2 | 92.3 | 96.2 | 92.3 | 92.3 | 92.3 | 92.3 | 66.7 | 92.3 | 76.9 | 96.2 | 76.9 | 88.5 | 83.3 | 76.9 | 80.8 | 23.1 | 88.6 |
| Substitution | 10 | 100.0 | 100.0 | 100.0 | 100.0 | 90.0 | 100.0 | 90.0 | 90.0 | 90.0 | 90.0 | 100.0 | 90.0 | 90.0 | 90.0 | 100.0 | 90.0 | 100.0 | 80.0 | 90.0 | 50.0 | 20.0 | 90.0 |
| Punctuation | 34 | 100.0 | 100.0 | 88.2 | 94.1 | 94.7 | 91.2 | 94.7 | 88.2 | 94.1 | 94.1 | 94.7 | 88.2 | 88.2 | 88.2 | 73.5 | 90.0 | 88.2 | 89.5 | 61.8 | 67.6 | 14.7 | 88.2 |
| Quotation marks | 34 | 100.0 | 100.0 | 88.2 | 94.1 | 89.5 | 91.2 | 92.8 | 88.2 | 88.2 | 88.2 | 88.2 | 88.2 | 88.2 | 88.2 | 73.5 | 84.5 | 78.4 | 82.4 | 61.8 | 67.6 | 14.7 | 82.6 |
| Subordination | 148 | 98.0 | 98.0 | 96.6 | 96.6 | 95.9 | 94.6 | 95.9 | 93.2 | 93.2 | 93.2 | 95.9 | 96.6 | 89.2 | 85.8 | 91.9 | 79.4 | 81.1 | 87.8 | 67.6 | 66.9 | 23.0 | 88.1 |

367

| categ | count | Onl-B | Trans | GPT4 | Claud | Unbab | Comma | Onl-W | Llama | Aya23 | IOLRe | Onl-A | Onl-G | IKUN | CUNIN | IKUNC | NVIDI | Occig | AISTA | TSUHI | MSLC | Mistr | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Adverbial clause | 6 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 66.7 | 100.0 | 100.0 | 83.3 | 83.3 | 33.3 | 33.3 | 93.7 |
| Cleft sentence | 12 | 91.7 | 91.7 | 100.0 | 91.7 | 83.3 | 100.0 | 91.7 | 91.7 | 91.7 | 100.0 | 91.7 | 81.0 | 100.0 | 91.7 | 83.3 | 75.0 | 66.7 | 61.9 | 75.0 | 75.0 | 25.0 | 85.7 |
| Contact clause | 21 | 95.2 | 95.2 | 100.0 | 100.0 | 95.2 | 95.2 | 100.0 | 100.0 | 100.0 | 100.0 | 95.2 | 81.0 | 100.0 | 100.0 | 81.0 | 85.7 | 61.9 | 42.9 | 61.9 | 42.9 | 14.3 | 84.4 |
| Indirect speech | 16 | 100.0 | 100.0 | 93.8 | 87.5 | 87.5 | 100.0 | 93.8 | 93.8 | 100.0 | 100.0 | 100.0 | 93.8 | 93.8 | 93.8 | 81.3 | 100.0 | 87.5 | 68.8 | 68.8 | 87.5 | 18.8 | 89.3 |
| Infinitive clause | 16 | 100.0 | 100.0 | 100.0 | 100.0 | 93.8 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 93.8 | 93.8 | 100.0 | 87.5 | 87.5 | 81.3 | 75.0 | 75.0 | 87.5 | 25.0 | 91.1 |
| Object clause | 16 | 100.0 | 100.0 | 93.8 | 93.8 | 100.0 | 100.0 | 87.5 | 87.5 | 93.8 | 93.8 | 93.8 | 75.0 | 75.0 | 93.8 | 93.8 | 87.5 | 93.8 | 62.5 | 62.5 | 68.8 | 18.8 | 86.9 |
| Pseudo-cleft sentence | 14 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 92.9 | 92.9 | 100.0 | 100.0 | 85.7 | 85.7 | 100.0 | 92.9 | 78.6 | 100.0 | 85.7 | 50.0 | 50.0 | 71.4 | 14.3 | 88.8 |
| Relative clause | 34 | 97.1 | 97.1 | 100.0 | 97.1 | 94.1 | 100.0 | 97.1 | 97.1 | 94.1 | 97.1 | 97.1 | 97.1 | 97.1 | 94.1 | 94.1 | 94.1 | 85.3 | 76.5 | 76.5 | 61.8 | 26.5 | 89.1 |
| Subject clause | 13 | 100.0 | 100.0 | 93.8 | 100.0 | 92.3 | 100.0 | 92.3 | 92.3 | 92.3 | 92.3 | 84.6 | 84.6 | 100.0 | 92.3 | 84.6 | 92.3 | 76.9 | 53.8 | 53.8 | 46.2 | 38.5 | 87.2 |
| Verb semantics | 18 | 83.3 | 83.3 | 72.2 | 72.2 | 77.8 | 66.7 | 77.8 | 77.8 | 83.3 | 72.2 | 72.2 | 61.1 | 66.7 | 50.0 | 55.6 | 44.4 | 55.6 | 16.7 | 16.7 | 11.1 | 0.0 | 58.7 |
| Verb tense/aspect/mood | 3225 | 98.2 | 98.2 | 98.4 | 97.6 | 91.4 | 96.6 | 96.6 | 93.6 | 96.2 | 94.9 | 92.3 | 77.0 | 80.7 | 97.1 | 77.4 | 82.3 | 72.5 | 67.1 | 33.5 | 42.2 | 23.8 | 81.6 |
| Conditional | 19 | 94.7 | 94.7 | 94.7 | 94.7 | 94.7 | 89.5 | 89.5 | 89.5 | 94.7 | 94.7 | 89.5 | 77.0 | 89.5 | 94.7 | 84.2 | 84.2 | 78.9 | 68.4 | 42.1 | 63.2 | 26.3 | 83.7 |
| Ditransitive - conditional I progressive | 60 | 100.0 | 100.0 | 95.0 | 93.3 | 95.0 | 100.0 | 100.0 | 100.0 | 96.7 | 100.0 | 100.0 | 73.3 | 98.3 | 86.7 | 85.0 | 75.0 | 63.3 | 58.3 | 58.3 | 30.0 | 15.0 | 79.3 |
| Ditransitive - conditional I simple | 52 | 100.0 | 100.0 | 86.5 | 86.5 | 98.1 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 40.4 | 57.7 | 38.5 | 71.2 | 67.3 | 71.2 | 57.7 | 17.3 | 5.8 | 21.2 | 72.9 |
| Ditransitive - conditional II progressive | 59 | 100.0 | 100.0 | 94.9 | 76.3 | 100.0 | 98.3 | 91.5 | 91.5 | 100.0 | 98.3 | 67.8 | 67.8 | 83.1 | 83.1 | 57.6 | 83.1 | 45.8 | 42.4 | 76.3 | 10.2 | 30.5 | 78.9 |
| Ditransitive - conditional II simple | 59 | 100.0 | 100.0 | 86.4 | 86.4 | 100.0 | 100.0 | 98.3 | 98.3 | 100.0 | 100.0 | 78.0 | 78.0 | 98.3 | 100.0 | 79.7 | 88.1 | 71.2 | 35.6 | 71.2 | 27.1 | 35.6 | 84.3 |
| Ditransitive - future I progressive | 57 | 100.0 | 94.7 | 96.5 | 89.5 | 100.0 | 96.5 | 96.5 | 96.5 | 98.2 | 96.5 | 70.2 | 70.2 | 86.0 | 86.0 | 86.0 | 86.0 | 77.2 | 54.4 | 80.7 | 21.1 | 19.3 | 83.2 |
| Ditransitive - future I simple | 112 | 100.0 | 99.1 | 95.5 | 87.5 | 100.0 | 99.1 | 99.1 | 99.1 | 99.1 | 94.6 | 87.5 | 87.5 | 75.9 | 75.9 | 89.3 | 76.8 | 67.9 | 33.0 | 80.4 | 24.1 | 28.6 | 82.6 |
| Ditransitive - future II progressive | 52 | 94.2 | 94.2 | 94.2 | 53.8 | 96.2 | 88.5 | 86.5 | 86.5 | 80.8 | 90.4 | 0.0 | 0.0 | 34.6 | 34.6 | 5.8 | 26.9 | 28.8 | 0.0 | 3.8 | 1.9 | 26.9 | 57.4 |
| Ditransitive - future II simple | 57 | 98.2 | 98.2 | 96.5 | 59.6 | 100.0 | 94.7 | 94.7 | 94.7 | 91.2 | 87.7 | 0.0 | 0.0 | 43.9 | 43.9 | 7.0 | 52.6 | 26.3 | 0.0 | 7.0 | 1.8 | 28.1 | 61.5 |
| Ditransitive - past perfect progressive | 53 | 92.5 | 92.5 | 98.1 | 86.8 | 92.9 | 92.5 | 71.7 | 71.7 | 83.0 | 83.0 | 26.4 | 26.4 | 67.9 | 67.9 | 43.4 | 88.7 | 54.7 | 47.2 | 73.6 | 22.6 | 24.5 | 73.4 |
| Ditransitive - past perfect simple | 56 | 98.2 | 96.4 | 100.0 | 89.3 | 92.9 | 96.4 | 80.4 | 80.4 | 82.1 | 81.1 | 51.8 | 51.8 | 69.6 | 69.6 | 73.2 | 76.8 | 58.9 | 39.3 | 76.8 | 46.4 | 21.4 | 77.9 |
| Ditransitive - past progressive | 54 | 100.0 | 100.0 | 92.6 | 92.6 | 96.3 | 96.3 | 85.2 | 85.2 | 96.3 | 96.3 | 72.2 | 72.2 | 83.3 | 83.3 | 98.1 | 77.8 | 72.2 | 33.3 | 72.2 | 24.1 | 20.4 | 81.0 |
| Ditransitive - present perfect progressive | 56 | 100.0 | 100.0 | 92.9 | 89.3 | 100.0 | 87.5 | 87.5 | 87.5 | 60.7 | 100.0 | 60.7 | 60.7 | 78.9 | 78.9 | 66.1 | 89.3 | 71.4 | 30.4 | 48.2 | 28.6 | 12.5 | 79.3 |
| Ditransitive - present perfect simple | 57 | 91.2 | 91.2 | 91.5 | 96.5 | 96.5 | 96.5 | 94.7 | 94.7 | 100.0 | 100.0 | 78.9 | 78.9 | 93.0 | 93.0 | 93.0 | 94.7 | 71.9 | 42.1 | 77.2 | 21.1 | 22.8 | 85.7 |
| Ditransitive - present progressive | 59 | 96.6 | 96.6 | 94.9 | 94.9 | 97.5 | 93.2 | 91.5 | 91.5 | 89.8 | 94.6 | 74.6 | 74.6 | 86.4 | 86.4 | 93.2 | 71.2 | 78.0 | 33.0 | 78.0 | 16.9 | 33.9 | 81.2 |
| Ditransitive - simple past | 79 | 98.7 | 98.7 | 96.2 | 94.9 | 97.5 | 98.2 | 94.9 | 94.9 | 94.9 | 98.7 | 75.9 | 75.9 | 93.7 | 94.9 | 93.7 | 92.4 | 74.7 | 41.8 | 86.1 | 22.8 | 22.8 | 84.0 |
| Ditransitive - simple present | 55 | 98.2 | 90.9 | 98.2 | 98.2 | 98.2 | 87.3 | 87.3 | 87.3 | 92.7 | 92.7 | 72.7 | 72.7 | 83.6 | 83.6 | 92.7 | 83.6 | 70.9 | 60.0 | 78.2 | 10.9 | 12.7 | 81.7 |
| Gerund | 24 | 95.8 | 95.8 | 100.0 | 98.2 | 100.0 | 95.8 | 83.3 | 83.3 | 95.8 | 85.7 | 72.7 | 91.7 | 79.2 | 79.2 | 83.3 | 87.5 | 70.8 | 45.8 | 66.7 | 41.7 | 35.7 | 84.1 |
| Imperative | 14 | 100.0 | 100.0 | 100.0 | 100.0 | 78.6 | 100.0 | 100.0 | 100.0 | 92.9 | 85.7 | 57.1 | 57.1 | 79.3 | 57.1 | 71.4 | 71.4 | 28.6 | 28.6 | 14.3 | 14.3 | 35.7 | 75.2 |
| Intransitive - conditional I progressive | 29 | 100.0 | 92.6 | 92.6 | 93.1 | 100.0 | 89.7 | 96.6 | 96.6 | 100.0 | 96.6 | 85.2 | 81.5 | 81.5 | 82.8 | 82.8 | 100.0 | 88.9 | 41.4 | 96.6 | 86.2 | 37.9 | 89.7 |
| Intransitive - conditional I simple | 27 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 96.3 | 96.3 | 100.0 | 100.0 | 82.8 | 89.7 | 81.5 | 81.5 | 81.5 | 100.0 | 82.8 | 74.1 | 92.6 | 77.8 | 22.2 | 88.4 |
| Intransitive - conditional II progressive | 29 | 100.0 | 100.0 | 96.6 | 96.6 | 100.0 | 96.6 | 82.8 | 82.8 | 96.6 | 96.6 | 89.7 | 82.8 | 69.0 | 69.0 | 69.0 | 58.6 | 62.1 | 6.9 | 58.6 | 13.8 | 24.1 | 76.5 |
| Intransitive - conditional II simple | 29 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 82.8 | 82.8 | 89.7 | 89.7 | 62.1 | 62.1 | 69.0 | 69.0 | 69.0 | 82.8 | 62.1 | 41.4 | 82.8 | 24.1 | 20.7 | 80.6 |
| Intransitive - future I progressive | 30 | 96.7 | 96.7 | 96.7 | 100.0 | 100.0 | 86.7 | 90.0 | 90.0 | 100.0 | 96.7 | 90.0 | 93.3 | 93.3 | 100.0 | 100.0 | 100.0 | 83.3 | 50.0 | 93.3 | 50.0 | 20.0 | 90.0 |
| Intransitive - future I simple | 69 | 98.6 | 98.6 | 97.1 | 98.6 | 100.0 | 89.9 | 95.7 | 95.7 | 100.0 | 100.0 | 90.0 | 95.7 | 98.6 | 100.0 | 92.8 | 100.0 | 88.4 | 60.9 | 93.3 | 84.1 | 20.3 | 91.0 |
| Intransitive - future II progressive | 27 | 100.0 | 100.0 | 100.0 | 96.3 | 88.9 | 88.9 | 96.3 | 96.3 | 14.8 | 91.2 | 0.0 | 0.0 | 29.6 | 74.1 | 29.6 | 77.8 | 48.1 | 3.7 | 18.5 | 0.0 | 22.2 | 69.8 |
| Intransitive - future II simple | 31 | 100.0 | 100.0 | 100.0 | 93.5 | 80.6 | 90.3 | 96.3 | 96.3 | 19.4 | 61.3 | 19.4 | 19.4 | 16.1 | 61.3 | 16.1 | 80.6 | 48.1 | 35.5 | 38.7 | 16.1 | 23.5 | 70.7 |
| Intransitive - past perfect progressive | 24 | 87.5 | 87.5 | 100.0 | 100.0 | 87.5 | 79.2 | 79.2 | 79.2 | 70.8 | 87.5 | 70.8 | 70.8 | 75.0 | 75.0 | 75.0 | 87.5 | 50.0 | 12.5 | 83.3 | 8.3 | 16.7 | 74.0 |
| Intransitive - past perfect simple | 34 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 94.1 | 82.4 | 82.4 | 76.5 | 79.4 | 85.3 | 85.3 | 88.2 | 94.1 | 85.3 | 88.2 | 73.5 | 38.2 | 88.2 | 58.8 | 17.6 | 83.5 |
| Intransitive - past progressive | 32 | 90.6 | 90.6 | 96.9 | 96.9 | 100.0 | 100.0 | 96.9 | 96.9 | 90.6 | 96.9 | 89.7 | 96.9 | 90.6 | 93.8 | 90.6 | 37.5 | 75.0 | 21.9 | 84.4 | 12.5 | 37.5 | 81.3 |
| Intransitive - present perfect progressive | 25 | 100.0 | 100.0 | 96.0 | 80.0 | 92.0 | 88.0 | 92.0 | 92.0 | 92.0 | 92.0 | 100.0 | 100.0 | 92.0 | 100.0 | 92.0 | 68.0 | 80.0 | 32.0 | 72.0 | 24.0 | 0.0 | 80.8 |
| Intransitive - present perfect simple | 30 | 100.0 | 100.0 | 100.0 | 93.3 | 86.7 | 86.7 | 100.0 | 100.0 | 90.0 | 100.0 | 90.0 | 100.0 | 100.0 | 100.0 | 93.3 | 83.3 | 70.0 | 66.7 | 86.7 | 33.3 | 13.3 | 86.5 |
| Intransitive - present progressive | 64 | 100.0 | 100.0 | 100.0 | 100.0 | 93.8 | 93.8 | 100.0 | 100.0 | 98.4 | 94.3 | 100.0 | 95.3 | 95.3 | 98.4 | 95.3 | 93.8 | 79.7 | 31.3 | 95.3 | 51.6 | 29.7 | 88.7 |
| Intransitive - simple past | 39 | 100.0 | 100.0 | 100.0 | 97.4 | 97.4 | 97.4 | 97.4 | 97.4 | 100.0 | 100.0 | 100.0 | 100.0 | 94.7 | 97.4 | 94.7 | 84.6 | 69.2 | 53.8 | 81.6 | 53.8 | 20.5 | 86.9 |
| Intransitive - simple present | 38 | 100.0 | 100.0 | 100.0 | 97.4 | 98.3 | 80.6 | 98.6 | 98.6 | 100.0 | 100.0 | 100.0 | 100.0 | 94.2 | 94.2 | 94.2 | 95.9 | 65.0 | 55.3 | 81.6 | 71.1 | 18.4 | 89.0 |
| Modal | 294 | 98.6 | 98.6 | 99.0 | 98.3 | 99.0 | 98.6 | 98.0 | 98.0 | 99.3 | 99.3 | 96.5 | 96.9 | 92.5 | 88.2 | 94.2 | 95.9 | 71.2 | 36.4 | 94.2 | 87.4 | 23.5 | 89.2 |
| Modal negated | 288 | 98.3 | 98.3 | 95.8 | 96.5 | 92.4 | 99.7 | 97.2 | 97.2 | 96.2 | 97.9 | 96.5 | 91.3 | 88.2 | 88.9 | 89.9 | 97.2 | 70.6 | 56.3 | 82.2 | 83.0 | 22.9 | 87.9 |
| Reflexive - conditional I progressive | 34 | 100.0 | 100.0 | 97.1 | 97.1 | 100.0 | 100.0 | 100.0 | 100.0 | 91.2 | 97.1 | 97.1 | 94.1 | 94.1 | 92.9 | 73.5 | 78.2 | 58.8 | 5.9 | 58.8 | 35.3 | 8.8 | 81.8 |
| Reflexive - conditional I simple | 28 | 100.0 | 100.0 | 92.9 | 92.9 | 100.0 | 96.4 | 100.0 | 100.0 | 78.6 | 100.0 | 100.0 | 96.4 | 60.7 | 60.7 | 92.9 | 78.6 | 60.7 | 28.6 | 75.0 | 35.7 | 14.3 | 81.3 |
| Reflexive - conditional II progressive | 31 | 100.0 | 100.0 | 71.0 | 71.0 | 93.5 | 93.5 | 96.8 | 96.8 | 90.3 | 96.8 | 100.0 | 74.2 | 45.2 | 45.2 | 16.1 | 80.6 | 54.8 | 3.2 | 64.5 | 35.5 | 32.3 | 73.3 |
| Reflexive - conditional II simple | 35 | 100.0 | 100.0 | 94.3 | 94.3 | 100.0 | 97.1 | 100.0 | 100.0 | 80.0 | 94.3 | 91.4 | 80.0 | 91.4 | 91.4 | 62.9 | 85.7 | 71.4 | 5.7 | 57.1 | 37.1 | 34.3 | 81.5 |
| Reflexive - future I progressive | 32 | 93.8 | 93.8 | 96.9 | 90.6 | 96.9 | 96.9 | 96.9 | 96.9 | 96.9 | 96.9 | 100.0 | 93.8 | 75.0 | 84.4 | 84.4 | 84.4 | 69.2 | 25.0 | 59.4 | 17.9 | 34.4 | 80.5 |
| Reflexive - future I simple | 68 | 98.5 | 98.5 | 98.5 | 88.2 | 100.0 | 97.1 | 95.6 | 95.6 | 98.5 | 97.1 | 100.0 | 82.4 | 82.4 | 82.4 | 82.4 | 94.1 | 50.0 | 33.8 | 83.8 | 27.9 | 26.5 | 82.6 |

| categ | count | Onl-B | Trans | GPT4 | Claud | Unbab | Comma | Onl-W | Llama | Aya23 | IOLRe | Onl-A | Onl-G | IKUN | CUNIN | IKUNC | NVIDI | Occig | AISTA | TSUHI | MSLC | Mistr | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reflexive - future II progressive | 29 | 93.1 | 93.1 | 100.0 | 100.0 | 48.3 | 96.6 | 96.6 | 93.1 | 93.1 | 96.6 | 96.6 | 100.0 | 0.0 | 34.5 | 6.9 | 27.6 | 34.5 | 6.9 | 0.0 | 0.0 | 31.0 | 59.4 |
| Reflexive - future II simple | 33 | 81.8 | 81.8 | 100.0 | 100.0 | 81.8 | 97.0 | 100.0 | 90.9 | 97.0 | 93.9 | 93.9 | 100.0 | 3.0 | 63.6 | 9.1 | 54.5 | 39.4 | 36.4 | 0.0 | 3.0 | 15.2 | 63.9 |
| Reflexive - past perfect progressive | 29 | 93.1 | 93.1 | 100.0 | 100.0 | 79.3 | 89.7 | 62.1 | 62.1 | 65.5 | 69.0 | 100.0 | 79.3 | 72.4 | 55.2 | 69.0 | 44.8 | 58.6 | 17.2 | 34.5 | 27.6 | 71.9 |
| Reflexive - past perfect simple | 27 | 100.0 | 100.0 | 96.3 | 100.0 | 96.3 | 92.6 | 96.3 | 74.1 | 70.4 | 77.8 | 100.0 | 85.2 | 74.1 | 63.0 | 74.1 | 63.0 | 70.4 | 3.7 | 48.1 | 14.8 | 76.2 |
| Reflexive - past progressive | 32 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 96.9 | 96.9 | 93.8 | 100.0 | 84.4 | 87.5 | 87.5 | 81.3 | 68.8 | 71.9 | 50.0 | 37.5 | 21.9 | 84.2 |
| Reflexive - present perfect progressive | 26 | 100.0 | 100.0 | 100.0 | 96.2 | 92.3 | 96.9 | 92.3 | 100.0 | 100.0 | 100.0 | 92.3 | 88.5 | 100.0 | 80.8 | 81.3 | 23.1 | 34.6 | 30.8 | 83.0 |
| Reflexive - present perfect simple | 32 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 96.9 | 100.0 | 100.0 | 100.0 | 100.0 | 96.9 | 75.0 | 90.6 | 75.0 | 69.2 | 75.0 | 31.3 | 37.5 | 9.4 | 85.1 |
| Reflexive - present progressive | 32 | 96.9 | 96.9 | 100.0 | 87.5 | 100.0 | 96.9 | 90.6 | 96.9 | 90.6 | 93.8 | 84.4 | 78.1 | 84.4 | 84.4 | 75.0 | 84.4 | 15.6 | 53.1 | 43.8 | 83.8 |
| Reflexive - simple past | 32 | 96.9 | 96.9 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 93.8 | 93.8 | 87.5 | 96.9 | 90.6 | 81.3 | 78.1 | 21.9 | 34.4 | 25.0 | 85.1 |
| Reflexive - simple present | 32 | 96.9 | 96.9 | 100.0 | 81.3 | 100.0 | 75.0 | 93.8 | 81.3 | 90.6 | 90.6 | 84.4 | 78.1 | 93.8 | 59.4 | 78.1 | 31.3 | 37.5 | 18.8 | 79.8 |
| Transitive - future II progressive | 30 | 100.0 | 100.0 | 100.0 | 83.3 | 86.7 | 86.7 | 53.3 | 46.7 | 3.3 | 3.3 | 46.7 | 36.7 | 3.3 | 0.0 | 13.3 | 62.7 |
| Transitive - conditional I progressive | 28 | 96.4 | 96.4 | 100.0 | 92.9 | 82.1 | 100.0 | 96.4 | 100.0 | 89.3 | 82.1 | 67.9 | 78.6 | 89.3 | 57.1 | 7.1 | 35.7 | 28.6 | 80.8 |
| Transitive - conditional I simple | 18 | 100.0 | 100.0 | 100.0 | 94.4 | 88.9 | 100.0 | 100.0 | 55.6 | 83.3 | 50.0 | 61.1 | 88.9 | 72.2 | 50.0 | 27.8 | 11.1 | 22.2 | 76.2 |
| Transitive - conditional II progressive | 27 | 100.0 | 100.0 | 96.3 | 92.6 | 88.9 | 96.3 | 100.0 | 100.0 | 70.4 | 59.3 | 85.2 | 88.9 | 74.1 | 7.4 | 44.4 | 33.3 | 82.7 |
| Transitive - conditional II simple | 30 | 100.0 | 100.0 | 100.0 | 96.7 | 96.7 | 93.3 | 100.0 | 93.3 | 93.3 | 76.7 | 93.3 | 96.7 | 66.7 | 73.3 | 16.7 | 66.7 | 23.3 | 84.9 |
| Transitive - future I progressive | 30 | 100.0 | 100.0 | 96.7 | 93.3 | 83.3 | 93.3 | 100.0 | 100.0 | 80.0 | 90.0 | 86.7 | 73.3 | 70.0 | 23.3 | 36.7 | 23.3 | 82.2 |
| Transitive - future I simple | 57 | 100.0 | 94.7 | 94.7 | 93.0 | 96.5 | 96.5 | 98.2 | 100.0 | 94.7 | 89.5 | 80.7 | 71.9 | 52.6 | 22.8 | 85.5 |
| Transitive - future II simple | 35 | 97.1 | 97.1 | 100.0 | 97.1 | 100.0 | 97.1 | 97.1 | 100.0 | 14.3 | 57.1 | 20.0 | 77.1 | 60.0 | 5.7 | 14.3 | 0.0 | 40.0 | 69.9 |
| Transitive - past perfect progressive | 24 | 91.7 | 91.7 | 95.8 | 95.8 | 79.2 | 70.8 | 75.0 | 79.2 | 50.0 | 41.7 | 75.0 | 87.5 | 58.3 | 8.3 | 62.5 | 37.5 | 75.0 |
| Transitive - past perfect simple | 25 | 100.0 | 100.0 | 100.0 | 100.0 | 96.0 | 92.0 | 96.0 | 95.8 | 100.0 | 88.0 | 84.0 | 76.0 | 68.0 | 8.0 | 68.0 | 8.0 | 82.1 |
| Transitive - past progressive | 38 | 97.4 | 84.2 | 84.2 | 71.1 | 84.2 | 81.6 | 76.3 | 78.9 | 86.8 | 63.2 | 78.9 | 71.1 | 68.4 | 42.1 | 42.1 | 23.7 | 72.3 |
| Transitive - present perfect progressive | 30 | 100.0 | 100.0 | 100.0 | 86.7 | 100.0 | 93.3 | 100.0 | 100.0 | 100.0 | 70.0 | 83.3 | 73.3 | 80.0 | 46.7 | 10.0 | 26.7 | 30.0 | 80.8 |
| Transitive - present perfect simple | 31 | 100.0 | 100.0 | 100.0 | 90.3 | 96.8 | 96.8 | 93.5 | 96.8 | 100.0 | 83.9 | 80.6 | 87.1 | 64.5 | 71.0 | 25.8 | 58.1 | 22.6 | 83.9 |
| Transitive - present progressive | 40 | 97.5 | 97.5 | 100.0 | 87.5 | 100.0 | 95.0 | 95.0 | 100.0 | 80.0 | 92.5 | 82.5 | 77.5 | 22.5 | 50.0 | 25.0 | 84.9 |
| Transitive - simple past | 38 | 100.0 | 100.0 | 97.4 | 94.7 | 94.7 | 94.7 | 94.7 | 97.4 | 89.5 | 92.1 | 78.9 | 81.6 | 68.4 | 52.6 | 50.0 | 26.3 | 85.1 |
| Transitive - simple present | 39 | 97.4 | 97.4 | 100.0 | 92.3 | 94.9 | 97.4 | 97.4 | 87.2 | 100.0 | 92.3 | 97.4 | 92.3 | 69.2 | 76.9 | 76.9 | 51.3 | 31.6 | 25.6 | 85.8 |
| Verb valency | 101 | 91.1 | 91.1 | 88.1 | 84.2 | 86.1 | 86.1 | 78.2 | 84.2 | 83.2 | 76.2 | 71.3 | 72.3 | 64.4 | 63.4 | 54.5 | 34.7 | 16.8 | 71.9 |
| Case government | 20 | 95.0 | 95.0 | 95.0 | 95.0 | 90.0 | 95.0 | 95.0 | 95.0 | 95.0 | 90.0 | 85.0 | 75.2 | 85.0 | 60.0 | 75.0 | 25.0 | 55.0 | 15.0 | 82.1 |
| Catenative verb | 16 | 100.0 | 100.0 | 100.0 | 93.8 | 100.0 | 87.5 | 87.5 | 100.0 | 100.0 | 81.3 | 81.3 | 87.5 | 62.5 | 31.3 | 37.5 | 43.8 | 12.5 | 83.3 |
| Mediopassive voice | 18 | 88.9 | 88.9 | 83.3 | 83.3 | 77.8 | 55.6 | 66.7 | 66.7 | 50.0 | 66.7 | 38.9 | 55.6 | 22.2 | 27.8 | 5.6 | 22.2 | 59.0 |
| Passive voice | 16 | 93.8 | 93.8 | 100.0 | 87.5 | 93.8 | 100.0 | 93.8 | 93.8 | 87.5 | 87.5 | 75.0 | 75.0 | 81.3 | 43.8 | 75.0 | 25.0 | 84.2 |
| Resultative | 19 | 89.5 | 89.5 | 94.7 | 84.2 | 84.2 | 84.2 | 89.5 | 78.9 | 84.2 | 73.7 | 73.7 | 63.2 | 47.4 | 36.8 | 10.5 | 5.3 | 69.9 |
| Semantic roles | 12 | 75.0 | 75.0 | 50.0 | 58.3 | 75.0 | 33.3 | 33.3 | 41.7 | 50.0 | 41.7 | 50.0 | 58.3 | 25.0 | 41.7 | 33.3 | 41.7 | 16.7 | 25.0 | 45.6 |
| micro-average | 4219 | 97.7 | 96.8 | 91.1 | 91.5 | 84.2 | 95.3 | 96.3 | 94.7 | 84.2 | 76.2 | 95.0 | 81.4 | 85.0 | 75.0 | 72.1 | 64.4 | 54.5 | 34.7 | 43.4 | 24.0 | 81.4 |
| phen. macro-average | 4219 | 96.8 | 96.8 | 86.1 | 91.1 | 91.0 | 92.5 | 93.5 | 91.4 | 93.0 | 77.3 | 80.1 | 78.5 | 69.4 | 68.3 | 40.4 | 40.5 | 24.2 | 80.0 |
| categ. macro-average | 4219 | 95.7 | 95.6 | 92.8 | 92.0 | 90.3 | 89.6 | 89.1 | 87.4 | 82.0 | 81.8 | 69.8 | 67.2 | 51.7 | 45.8 | 23.9 | 79.7 |

Table 10: Accuracies (%) of successful translations on the phenomenon-level for English–German. The boldface indicates the significantly best-performing systems per row.

| categ | count | Yande | Claud | Unbab | Comma | Onl-G | Onl-W | GPT4 | IOLRe | Trans | Onl-B | Onl-A | Aya23 | IKUN | Llama | IKUNC | CUNID | NVIDI | TSUHI | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | 20 | 90.0 | 95.0 | 90.0 | 90.0 | 70.0 | 90.0 | 90.0 | 50.0 | 50.0 | 55.0 | 55.0 | 85.0 | 70.0 | 75.0 | 60.0 | 65.0 | 35.0 | 15.0 | 69.2 |
| Lexical ambiguity | 20 | 90.0 | 95.0 | 90.0 | 90.0 | 70.0 | 90.0 | 90.0 | 50.0 | 50.0 | 55.0 | 55.0 | 85.0 | 70.0 | 75.0 | 60.0 | 65.0 | 35.0 | 15.0 | 69.2 |
| Coordination & ellipsis | 86 | 87.2 | 80.2 | 84.9 | 82.6 | 83.7 | 76.7 | 72.2 | 72.1 | 72.1 | 74.4 | 69.8 | 77.9 | 65.1 | 69.8 | 65.1 | 72.1 | 54.7 | 47.7 | 74.1 |
| Gapping | 18 | 77.8 | 66.7 | 77.8 | 72.2 | 88.9 | 66.7 | 55.6 | 55.6 | 66.7 | 44.4 | 44.4 | 72.2 | 66.7 | 44.4 | 38.9 | 5.6 | 60.8 |
| Pseudogapping | 13 | 76.9 | 84.6 | 84.6 | 76.9 | 76.9 | 53.8 | 61.5 | 61.5 | 61.5 | 53.8 | 53.8 | 61.5 | 66.7 | 53.8 | 23.1 | 38.5 | 65.4 |
| Right node raising | 14 | 92.9 | 78.6 | 92.9 | 78.6 | 85.7 | 78.6 | 92.9 | 92.9 | 71.4 | 71.4 | 78.6 | 92.9 | 78.6 | 71.4 | 85.7 | 57.1 | 80.2 |
| Sluicing | 8 | 100.0 | 100.0 | 75.0 | 87.5 | 87.5 | 87.5 | 75.0 | 75.0 | 87.5 | 75.0 | 75.0 | 62.5 | 100.0 | 75.0 | 37.5 | 37.5 | 57.1 | 77.1 |
| Stripping | 19 | 89.5 | 84.2 | 89.5 | 89.5 | 100.0 | 89.5 | 94.7 | 89.5 | 84.2 | 78.9 | 89.5 | 84.2 | 62.5 | 78.9 | 84.2 | 89.5 | 63.2 | 86.8 |
| VP-ellipsis | 14 | 92.9 | 78.6 | 85.7 | 85.7 | 78.6 | 85.7 | 78.6 | 71.4 | 57.1 | 57.1 | 71.4 | 57.1 | 92.9 | 64.3 | 64.3 | 63.2 | 64.3 | 74.2 |
| False friends | 15 | 86.7 | 86.7 | 86.7 | 73.3 | 86.7 | 66.7 | 66.7 | 66.7 | 80.0 | 66.7 | 66.7 | 60.0 | 46.7 | 53.3 | 70.7 |

369

| categ | count | Yande | Claud | Unbab | Comma | Onl-G | Onl-W | GPT4 | IOLRe | Trans | Onl-B | Onl-A | Aya23 | IKUN | Llama | IKUNC | CUNID | NVIDI | TSUHI | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Function word | 34 | 97.1 | 88.2 | 94.1 | 100.0 | 94.1 | 100.0 | 97.1 | 91.2 | 94.1 | 94.1 | 88.2 | 94.1 | 85.3 | 85.3 | 82.4 | 73.5 | 73.5 | 70.6 | 89.1 |
| Focus particle | 15 | 93.3 | 73.3 | 86.7 | 100.0 | 86.7 | 100.0 | 93.3 | 80.0 | 86.7 | 86.7 | 86.7 | 86.7 | 73.3 | 73.3 | 73.3 | 66.7 | 80.0 | 73.3 | 83.3 |
| Question tag | 19 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 89.5 | 100.0 | 94.7 | 94.7 | 89.5 | 78.9 | 68.4 | 68.4 | 93.6 |
| LDD & interrogatives | 81 | 97.5 | 93.8 | 97.5 | 91.4 | 96.3 | 95.1 | 91.4 | 90.1 | 91.4 | 91.4 | 85.2 | 86.4 | 82.7 | 80.2 | 76.5 | 82.7 | 70.4 | 59.3 | 86.6 |
| Inversion | 22 | 95.5 | 95.5 | 90.9 | 90.9 | 100.0 | 95.5 | 90.9 | 90.9 | 95.5 | 95.5 | 86.4 | 86.4 | 90.9 | 81.8 | 72.7 | 90.9 | 68.2 | 54.5 | 87.4 |
| Modifying Comparison | 4 | 100.0 | 100.0 | 100.0 | 100.0 | 75.0 | 100.0 | 100.0 | 100.0 | 75.0 | 75.0 | 75.0 | 75.0 | 75.0 | 75.0 | 100.0 | 75.0 | 75.0 | 50.0 | 84.7 |
| Multiple connectors | 13 | 100.0 | 100.0 | 100.0 | 100.0 | 92.3 | 100.0 | 100.0 | 92.3 | 84.6 | 84.6 | 84.6 | 84.6 | 100.0 | 84.6 | 76.9 | 92.3 | 76.9 | 76.9 | 90.6 |
| Pied-piping | 14 | 92.9 | 92.9 | 100.0 | 100.0 | 100.0 | 92.9 | 92.9 | 100.0 | 100.0 | 100.0 | 92.9 | 92.9 | 85.7 | 85.7 | 85.7 | 78.6 | 100.0 | 92.9 | 93.7 |
| Preposition stranding | 17 | 100.0 | 100.0 | 100.0 | 88.2 | 94.1 | 100.0 | 88.2 | 88.2 | 100.0 | 100.0 | 82.4 | 88.2 | 76.5 | 76.5 | 82.4 | 64.7 | 64.7 | 47.1 | 85.0 |
| Topicalization | 11 | 100.0 | 72.7 | 100.0 | 72.7 | 100.0 | 81.8 | 81.8 | 72.7 | 72.7 | 72.7 | 72.7 | 81.8 | 54.5 | 72.7 | 54.5 | 90.9 | 54.5 | 27.3 | 74.7 |
| Lexical Morphology | 41 | 97.6 | 92.7 | 90.2 | 92.7 | 82.9 | 73.2 | 80.5 | 73.2 | 75.6 | 75.6 | 70.7 | 68.3 | 75.6 | 75.6 | 63.4 | 53.7 | 34.1 | 26.8 | 72.4 |
| Functional shift | 17 | 100.0 | 100.0 | 100.0 | 94.1 | 88.2 | 88.2 | 100.0 | 94.1 | 88.2 | 88.2 | 76.5 | 82.4 | 88.2 | 88.2 | 82.4 | 64.7 | 41.2 | 41.2 | 83.7 |
| Noun formation (er) | 24 | 95.8 | 87.5 | 83.3 | 91.7 | 79.2 | 62.5 | 66.7 | 58.3 | 66.7 | 66.7 | 66.7 | 58.3 | 66.7 | 66.7 | 50.0 | 45.8 | 29.2 | 16.7 | 64.4 |
| MWE | 96 | 87.5 | 84.4 | 78.1 | 83.3 | 80.2 | 71.9 | 77.1 | 76.0 | 71.9 | 72.9 | 70.8 | 67.7 | 69.8 | 66.7 | 66.7 | 52.1 | 40.6 | 33.3 | 69.5 |
| Collocation | 13 | 100.0 | 84.6 | 76.9 | 69.2 | 92.3 | 84.6 | 76.9 | 76.9 | 69.2 | 69.2 | 84.6 | 69.2 | 76.9 | 61.5 | 69.2 | 53.8 | 38.5 | 15.4 | 70.5 |
| Compound | 14 | 71.4 | 78.6 | 57.1 | 71.4 | 64.3 | 42.9 | 50.0 | 78.6 | 57.1 | 50.0 | 42.9 | 64.3 | 50.0 | 50.0 | 42.9 | 50.0 | 28.6 | 14.3 | 53.6 |
| Idiom | 17 | 94.1 | 88.2 | 70.6 | 88.2 | 52.9 | 47.1 | 76.5 | 70.6 | 52.9 | 52.9 | 47.1 | 41.2 | 52.9 | 41.2 | 47.1 | 41.2 | 5.9 | 11.8 | 54.6 |
| Nominal MWE | 17 | 76.5 | 88.2 | 88.2 | 88.2 | 82.4 | 94.4 | 82.4 | 70.6 | 100.0 | 100.0 | 94.1 | 88.2 | 76.5 | 70.6 | 70.6 | 58.8 | 58.8 | 52.9 | 81.0 |
| Prepositional MWE | 18 | 94.4 | 88.9 | 94.4 | 94.4 | 94.4 | 88.9 | 88.9 | 77.8 | 83.3 | 83.3 | 88.9 | 100.0 | 88.9 | 94.4 | 88.9 | 61.1 | 83.3 | 72.2 | 87.3 |
| Verbal MWE | 17 | 88.2 | 76.5 | 76.5 | 82.4 | 76.5 | 64.7 | 82.4 | 82.4 | 70.6 | 70.6 | 64.7 | 41.2 | 70.6 | 76.5 | 76.5 | 47.1 | 23.5 | 23.5 | 66.3 |
| Named entity & terminology | 80 | 83.8 | 95.0 | 87.5 | 81.3 | 80.0 | 80.0 | 81.3 | 73.8 | 80.0 | 80.0 | 77.5 | 71.3 | 62.5 | 57.5 | 60.0 | 57.5 | 56.3 | 41.3 | 73.7 |
| Date | 20 | 100.0 | 95.0 | 95.0 | 95.0 | 95.0 | 80.0 | 95.0 | 85.0 | 95.0 | 95.0 | 95.0 | 85.0 | 80.0 | 77.5 | 80.0 | 85.0 | 70.0 | 75.0 | 89.2 |
| Domainspecific Term | 5 | 40.0 | 80.0 | 60.0 | 60.0 | 60.0 | 80.0 | 40.0 | 60.0 | 60.0 | 60.0 | 60.0 | 60.0 | 0.0 | 40.0 | 20.0 | 60.0 | 60.0 | 0.0 | 50.0 |
| Measuring Unit | 18 | 72.2 | 100.0 | 100.0 | 94.4 | 94.4 | 94.4 | 94.4 | 94.4 | 94.4 | 94.4 | 88.9 | 88.9 | 83.3 | 100.0 | 77.8 | 77.8 | 72.2 | 44.4 | 87.0 |
| Onomatopeia | 11 | 100.0 | 100.0 | 72.7 | 63.6 | 54.5 | 72.7 | 81.8 | 45.5 | 54.5 | 54.5 | 54.5 | 63.6 | 45.5 | 54.5 | 45.5 | 45.5 | 18.2 | 0.0 | 57.1 |
| Proper Name & Location | 26 | 80.8 | 84.6 | 84.6 | 73.1 | 73.1 | 69.2 | 69.2 | 65.4 | 73.1 | 73.1 | 69.2 | 53.8 | 53.8 | 65.4 | 46.2 | 26.9 | 50.0 | 38.5 | 64.1 |
| Non-verbal agreement | 98 | 94.9 | 95.9 | 91.8 | 93.9 | 90.8 | 89.8 | 90.8 | 92.9 | 80.6 | 80.6 | 80.6 | 92.9 | 83.7 | 86.7 | 85.7 | 81.6 | 73.5 | 65.3 | 86.2 |
| Coreference | 24 | 87.5 | 91.7 | 83.3 | 87.5 | 75.0 | 83.3 | 83.3 | 83.3 | 54.2 | 54.2 | 66.7 | 83.3 | 79.2 | 75.0 | 83.3 | 70.8 | 54.2 | 58.3 | 75.2 |
| Genitive | 16 | 93.8 | 93.8 | 81.3 | 87.5 | 93.8 | 81.3 | 81.3 | 87.5 | 87.5 | 87.5 | 93.8 | 93.8 | 75.0 | 87.5 | 81.3 | 81.3 | 68.8 | 68.8 | 84.7 |
| Personal Pronoun Coreference | 19 | 100.0 | 100.0 | 100.0 | 100.0 | 94.7 | 100.0 | 100.0 | 89.5 | 100.0 | 89.5 | 78.9 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 89.5 | 96.8 |
| Possessive Pronouns | 22 | 95.5 | 95.5 | 95.5 | 90.9 | 90.9 | 95.5 | 95.5 | 100.0 | 95.5 | 95.5 | 90.9 | 90.9 | 86.4 | 90.9 | 86.4 | 77.3 | 86.4 | 40.9 | 89.4 |
| Substitution | 17 | 100.0 | 100.0 | 100.0 | 94.1 | 100.0 | 94.1 | 94.1 | 94.1 | 82.4 | 82.4 | 76.5 | 76.5 | 76.5 | 82.4 | 76.5 | 82.4 | 58.8 | 76.5 | 87.3 |
| Punctuation | 13 | 92.3 | 92.3 | 92.3 | 100.0 | 92.3 | 76.9 | 61.5 | 76.9 | 84.6 | 84.6 | 92.3 | 84.6 | 84.6 | 61.5 | 84.6 | 100.0 | 92.3 | 76.9 | 85.0 |
| Quotation marks | 13 | 92.3 | 92.3 | 92.3 | 100.0 | 92.3 | 76.9 | 61.5 | 76.9 | 84.6 | 84.6 | 92.3 | 84.6 | 84.6 | 75.0 | 84.6 | 100.0 | 92.3 | 76.9 | 85.0 |
| Subordination | 115 | 98.3 | 94.8 | 98.3 | 88.7 | 95.7 | 96.5 | 94.8 | 93.0 | 94.8 | 94.8 | 93.0 | 86.1 | 86.1 | 88.7 | 83.5 | 80.0 | 80.9 | 67.0 | 89.7 |
| Adverbial clause | 9 | 88.9 | 100.0 | 100.0 | 88.9 | 100.0 | 88.9 | 100.0 | 88.9 | 100.0 | 100.0 | 88.9 | 88.9 | 66.7 | 88.9 | 88.9 | 88.9 | 66.7 | 55.6 | 87.0 |
| Cleft sentence | 17 | 100.0 | 94.1 | 100.0 | 82.4 | 100.0 | 88.2 | 94.1 | 88.2 | 94.1 | 94.1 | 88.2 | 94.1 | 82.4 | 58.8 | 64.7 | 70.6 | 76.5 | 64.7 | 84.6 |
| Complex object | 18 | 100.0 | 94.4 | 100.0 | 94.4 | 100.0 | 100.0 | 94.4 | 100.0 | 94.4 | 94.4 | 88.9 | 94.4 | 94.4 | 100.0 | 83.3 | 83.3 | 83.3 | 61.1 | 93.2 |
| Contact clause | 12 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 91.7 | 100.0 | 91.7 | 100.0 | 100.0 | 100.0 | 91.7 | 91.7 | 91.7 | 75.0 | 83.3 | 91.7 | 83.3 | 94.0 |
| Infinitive clause | 25 | 96.0 | 88.0 | 96.0 | 88.0 | 92.0 | 100.0 | 92.0 | 100.0 | 92.0 | 92.0 | 100.0 | 96.0 | 92.0 | 100.0 | 92.0 | 72.0 | 84.0 | 68.0 | 91.1 |
| Participle clause | 22 | 100.0 | 95.5 | 95.5 | 86.4 | 90.9 | 90.9 | 95.5 | 90.9 | 90.9 | 90.9 | 86.4 | 77.3 | 81.8 | 81.8 | 72.7 | 86.4 | 81.8 | 59.1 | 86.4 |
| Subject clause | 12 | 100.0 | 100.0 | 100.0 | 83.3 | 100.0 | 91.7 | 91.7 | 83.3 | 100.0 | 100.0 | 100.0 | 83.3 | 83.3 | 83.3 | 100.0 | 83.3 | 75.0 | 83.3 | 92.6 |
| Verb semantics | 20 | 100.0 | 90.0 | 95.0 | 70.0 | 95.0 | 85.0 | 65.0 | 80.0 | 85.0 | 85.0 | 85.0 | 65.0 | 80.0 | 75.0 | 70.0 | 55.0 | 30.0 | 35.0 | 74.7 |
| Verb tense/aspect/mood | 169 | 87.0 | 90.5 | 90.5 | 87.0 | 89.3 | 88.8 | 89.9 | 85.2 | 85.2 | 85.2 | 86.4 | 85.2 | 81.1 | 84.6 | 84.6 | 78.1 | 69.2 | 45.0 | 82.9 |
| Conditional | 25 | 96.0 | 100.0 | 100.0 | 96.0 | 96.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 92.0 | 88.0 | 88.0 | 76.0 | 80.0 | 72.0 | 60.0 | 91.3 |
| Ditransitive | 34 | 82.4 | 94.1 | 97.1 | 94.1 | 94.1 | 88.2 | 94.1 | 92.1 | 91.2 | 91.2 | 97.1 | 94.1 | 88.2 | 91.2 | 94.1 | 94.1 | 82.4 | 50.0 | 89.1 |
| Gerund | 19 | 94.7 | 84.2 | 84.2 | 78.9 | 94.7 | 73.7 | 89.5 | 94.7 | 94.7 | 94.7 | 89.5 | 78.9 | 73.7 | 89.5 | 94.7 | 73.7 | 47.4 | 42.1 | 81.3 |
| Imperative | 24 | 91.7 | 83.3 | 87.5 | 83.3 | 91.7 | 87.5 | 83.3 | 83.3 | 83.3 | 83.3 | 70.8 | 83.3 | 75.0 | 75.0 | 79.2 | 66.7 | 62.5 | 20.8 | 75.9 |
| Intransitive | 29 | 82.8 | 86.2 | 82.8 | 89.7 | 79.3 | 86.2 | 86.2 | 82.8 | 69.0 | 69.0 | 75.9 | 75.9 | 72.4 | 86.2 | 79.3 | 86.2 | 79.3 | 55.2 | 79.3 |
| Reflexive | 19 | 89.5 | 89.5 | 84.2 | 84.2 | 89.5 | 89.5 | 78.9 | 78.9 | 84.2 | 84.2 | 84.2 | 73.7 | 84.2 | 68.4 | 78.9 | 52.6 | 68.4 | 31.6 | 77.5 |
| Transitive | 19 | 73.7 | 94.7 | 94.7 | 78.9 | 84.2 | 89.5 | 94.7 | 89.5 | 89.5 | 89.5 | 84.2 | 94.7 | 89.5 | 89.5 | 78.9 | 78.9 | 57.9 | 47.4 | 82.5 |
| Verb valency | 126 | 93.7 | 88.1 | 83.3 | 81.0 | 84.9 | 86.5 | 81.7 | 81.0 | 83.3 | 83.3 | 77.8 | 75.4 | 78.6 | 76.2 | 73.8 | 68.3 | 58.7 | 50.0 | 78.1 |
| Case government | 24 | 100.0 | 100.0 | 95.8 | 95.8 | 100.0 | 100.0 | 95.8 | 95.8 | 95.8 | 95.8 | 95.8 | 91.7 | 100.0 | 95.8 | 87.5 | 87.5 | 87.5 | 79.2 | 94.9 |
| Catenative verb | 25 | 96.0 | 96.0 | 92.0 | 84.0 | 92.0 | 96.0 | 84.0 | 88.0 | 96.0 | 96.0 | 92.0 | 88.0 | 84.0 | 92.0 | 88.0 | 84.0 | 84.0 | 72.0 | 89.1 |

| categ | count | Yande | Claud | Unbab | Comma | Onl-G | Onl-W | GPT4 | IOLRe | Trans | Onl-B | Onl-A | Aya23 | IKUN | Llama | IKUNC | CUNID | NVIDI | TSUHI | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mediopassive voice | 18 | **94.4** | **83.3** | **83.3** | **83.3** | **77.8** | **77.8** | **72.2** | **77.8** | **83.3** | **83.3** | **72.2** | **61.1** | **83.3** | 72.2 | 77.8 | 77.8 | 33.3 | 38.9 | 74.1 |
| Passive voice | 25 | **100.0** | **96.0** | 88.0 | **96.0** | **92.0** | 88.0 | **92.0** | 88.0 | **96.0** | **96.0** | **92.0** | **92.0** | **96.0** | 76.0 | 80.0 | 76.0 | 80.0 | 52.0 | 87.6 |
| Resultative | 18 | **88.9** | **83.3** | **83.3** | **61.1** | **72.2** | **83.3** | **83.3** | **66.7** | **61.1** | **61.1** | 44.4 | 44.4 | 50.0 | **72.2** | 55.6 | 33.3 | 16.7 | 22.2 | 60.2 |
| Semantic roles | 16 | **75.0** | **56.3** | **43.8** | **50.0** | **62.5** | **43.8** | **62.5** | **50.0** | **50.0** | **50.0** | **56.3** | **56.3** | **37.5** | 31.3 | **37.5** | 31.3 | 18.8 | 12.5 | 45.5 |
| micro-average | 994 | **91.8** | 90.4 | 89.4 | 86.8 | 87.8 | 86.1 | 85.2 | 83.3 | 82.3 | 82.2 | 80.7 | 80.4 | 78.1 | 79.0 | 75.1 | 71.0 | 62.2 | 50.0 | 80.1 |
| phen. macro-average | 994 | **91.4** | 90.2 | 88.8 | 86.1 | 87.1 | 85.7 | 84.3 | 82.4 | 81.6 | 81.4 | 79.6 | 79.2 | 76.2 | 77.4 | 73.8 | 70.3 | 60.8 | 49.5 | 79.2 |
| categ. macro-average | 994 | **92.4** | 90.5 | 90.0 | 87.7 | 87.4 | 83.6 | 81.7 | 81.3 | 79.7 | 79.7 | 79.6 | 79.0 | 77.2 | 76.4 | 72.6 | 69.0 | 59.7 | 49.0 | 78.7 |

Table 11: Accuracies (%) of successful translations on the phenomenon-level for English–Russian. The boldface indicates the significantly best-performing systems per row.

371

# IsoChronoMeter: A simple and effective isochronic translation evaluation metric

**Nikolai Rozanov**[1,2] **Vikentiy Pankov**[1] **Dmitrii Mukhutdinov**[1] **Dima Vypirailenko**[1]
[1]Brask AI
{vikentiy@brask.ai, dm@brask.ai, dima@brask.ai}
[2]Imperial College London
{nikolai.rozanov@gmail.com}

## Abstract

Machine translation (MT) has come a long way and is readily employed in production systems to serve millions of users daily. With the recent advances in generative AI, a new form of translation is becoming possible - video dubbing. This work motivates the importance of isochronic translation, especially in the context of automatic dubbing, and introduces 'IsoChronoMeter' (ICM). ICM is a simple yet effective metric to measure isochrony of translations in a scalable and resource-efficient way without the need for gold data, based on state-of-the-art text-to-speech (TTS) duration predictors. We motivate IsoChronoMeter and demonstrate its effectiveness. Using ICM we demonstrate the shortcomings of state-of-the-art translation systems and show the need for new methods. We release the code at this URL: https://github.com/braskai/isochronometer.

## 1 Introduction

The isochronic translation is a practice of ensuring that the timing of speech in translated content matches the original. It has become increasingly crucial in AI-driven dubbing. As the demand for multilingual audiovisual content grows, the ability to maintain the natural rhythm and pacing of the original language through isochronic translation is vital for the success of AI dubbing systems. Traditionally, human translators and voice actors have emphasized the importance of synchronizing translated dialogue with on-screen visuals to ensure a seamless viewing experience. This synchronization, known as isochrony, is essential for maintaining the illusion that the actors are speaking the translated language, matching their lip movements and pauses with the new audio. Recently, with the advancements in neural machine translation and text-to-speech systems, researchers have strived to replicate this isochrony automatically, aiming to preserve the speech-pause

structure of the original language in the translated content (Tam et al., 2022; Lakew et al., 2022).

Another way to ensure good dubbing synchronization is lip-sync. While lip-syncing is often employed to ensure synchronicity in dubbing, it presents significant challenges. Lip-syncing may force the translated dialogue to unnaturally conform to the lip movements of the original actors, potentially compromising the accuracy and fluidity of the translation. This often results in awkward or stilted dialogue, which can spoil the overall viewing experience. Additionally, due to the linguistic differences between languages, perfect lip-syncing can be impractical, leading to less faithful representations of the original content. Consequently, although lip-syncing can enhance visual alignment, it is not the optimal approach for achieving high-quality dubbing, especially when the goal is to maintain the natural flow and meaning of the original speech (Brannon et al., 2023) Research has demonstrated that integrating isochronic translation into AI dubbing significantly enhances the quality and naturalness of dubbed content, making it more acceptable to global audiences. By preserving the timing and rhythm of the original speech, these systems not only improve the technical quality of the translation but also maintain the emotional and narrative integrity of the content (Chronopoulou et al., 2023a).

### 1.1 Contribution

In this work we present a new isochronic metric, 'IsoChronoMeter' (ICM), and evaluation dataset for isochronic translation and demonstrate that 'normal' translations, even by state-of-the-art systems based on LLMs and human translations, without isochrony in mind, do not achieve a good level of isochronic translation. This highlights the importance of developing specialised translation systems

that are able to perform isochronic translation.

# 2 Background

## 2.1 Isochronic translation and metrics

Initial approaches that wanted to achieve isochronic translation focused on isometric translation (Federico et al., 2020; Karakanta et al., 2020; Lakew et al., 2021b,a), where the aim of MT systems was to translate text to achieve a similar target length. Spoken language translation benchmarks included 'isometric' subtasks (Anastasopoulos et al., 2022). However, research showed (Brannon et al., 2023) that isometric translations do not result in temporally synchronized speech after dubbing, i.e. isometricity does not really correlate with isochronicity.

This led to the most recent approaches focusing on isochronic translation instead (Wu et al., 2022; Chronopoulou et al., 2023b). However, this direction of research is fairly new: a dedicated *dubbing* task in spoken language translation benchmarks was first introduced in 2023 (Agarwal et al., 2023), and the degree of isochronicity is either measured subjectively by humans (Federico et al., 2020) or approximated via auxiliary metrics such as phoneme-based evaluation metrics (Chronopoulou et al., 2023a). VideoDubber (Wu et al., 2022) was among the first to successfully employ automatic duration predictors to evaluate isochronicity of the translated text, but their 'isochronic' metric is still based on human feedback, and hence cannot be applied at scale. Therefore, we conclude that there is a need to evaluate isochronic translations automatically; furthermore, since automatic dubbing pipelines in practice work with a pipeline approach (i.e. first running ASR and then later translating), it is crucial to introduce a text-based isochronic translation evaluation suite.

## 2.2 Evaluation Datasets

Collecting translation datasets requires a lot of effort especially for spoken data. Existing work includes Must-C (Di Gangi et al., 2019), GigaST (Ye et al., 2023), CoVost-2 (Wang et al., 2020) and Anim-400K (Cai et al., 2024). Datasets that specifically focus on isochronic translation using professional dubbing services only seem to exist privately (Brannon et al., 2023). In our work, we choose CoVost-2 due to its permissible licenses and availability of languages.

## 2.3 Identified challenges.

A full isochrony estimation would require humans to read out the given original text and the translated text. Additionally, one would need to attempt to find speakers that have similar speaking rates in their respective languages. We propose to overcome this by a novel isochrony metric that is easy to compute without the need of human annotations (i.e. human speech) and a joined isochrony and translation quality metric without the need of gold annotations.

# 3 Method

## 3.1 Metrics

### 3.1.1 IsoChronoMeter (ours) - automatic reference-free isochrony estimation

To compute isochrony metrics, we utilize the open-source TTSMMS project[1], which is based on Vits TTS (Kim et al., 2021) and MMS (Pratap et al., 2023), which supports multiple languages. Specifically, we use the duration predictor component to estimate the durations of the original text and translated texts generated by different machine translation (MT) models. As an isochrony metric, we apply a simple relative absolute error formula. Since the duration predictors for most languages are trained on similar domains (biblical texts) and share the same architecture, we expect them to produce similar durations adjusted to the average speaking rate of each language. Therefore, we can assume that if the texts are isochronic, their durations will be close. Concretely, IsoChronoMeter (ICM) is:

$$ICM = \tag{1}$$
$$\left\| \frac{MMS(original) - MMS(translated)}{MMS(original)} \right\|_2^2$$

Therefore, ICM is 0 if the duration of the original audio length prediction and the translated audio length prediction are the same; otherwise ICM represents a percentage of how much the two audio durations deviate from one another, e.g. ICM = 0.5 means that one of the audio duration predictions is half the duration prediction of the other.

### 3.1.2 Blaser2.0 - automatic reference-free machine translation quality estimation

To estimate MT quality (QE), we utilise BLASER2.0 models (Chen et al., 2023), based on

---

[1]Only a github is available: https://github.com/wannaphong/ttsmms

SONAR embeddings (Duquenne et al., 2023), to predict cross-lingual semantic similarities between the translation and original texts. Concretely,

$$QE = \tag{2}$$
$$blaser2\big(sonar(original), sonar(translated)\big)$$

Chen et al. (2023) show that such quality-estimation metrics outperform standard metrics such as bleu.

### 3.1.3 Adjusted-IsoChronoMeter - automatic reference-free joined machine translation quality and isochrony estimation

We also propose another metric based on the combination of IsoChronoMeter and Blaser, Adjusted-IsoChronoMeter (A-ICM). Concretely:

$$AICM = (1 - ICM) * QE \tag{3}$$

### 3.2 Effectiveness of the isochronic metric

Modern TTS systems such as Elevenlabs[2] or Rask AI[3] are able to produce realistic voices and voice-clones in multiple languages. These synthetic voices share incredible similarity with human voices. Therefore, we argue that using TTS as a proxy for the duration of human speech is effective. However, since we use a duration predictor for a TTS system, we need to show that the duration predictor is faithful to the real duration of a TTS system. To show this, we conduct a simple validation against an internal dataset of a few hours of English audio data, see Figure 1. Concretely, for each audio file we generate the 'original' TTS-generated audio sample and compare it against three predictions. Firstly, we compare against a 'repeat run', i.e. we generate a second audio file using the same TTS provider. Interestingly, the repeat run does not produce 0 or close to 0 error, in fact for $< 15$ words the error is above 5%. Secondly, we compare against our standard duration predictor. Finally, we also compare against a fine-tuned version of the duration predictor. We find that for small word counts the error is quite significant for all three, but especially for the not fine-tuned duration predictor. For $x > 15$, however, all three curves start converging and are within 5% error of one another. Therefore our metric becomes effective after a sufficiently large threshold of words.

[2]`elevenlabs.com`
[3]`rask.ai`



Figure 1: Dataset on English data. On the y-axis there is the relative absolute error between an original TTS-generated audio-sample and the associated prediction. On the x-axis is the number of total words used for the audio sample / prediction. Three curves show a secondary TTS-generated audio-sample (interestingly showing a big error for a few words), a fine-tuned duration predictor and the original duration predictor.

### 3.3 Dataset filtering

To demonstrate our metric and the need for isochronic translation engines, we create a small high-quality dataset from the CoVoST-2 (Wang et al., 2020), which is based on CommonVoice (Ardila et al., 2020). Specifically, taking into account the effectiveness of our metric after a specific threshold, we first filter the CoVoST-2 dataset by size. To find a good trade-off between dataset size and metric efficiency, we plot the histogram of counts and discover that above 20 tokens strikes a good balance, see Figure 2. In particular, we observe that if we choose the number of tokens to be 25 and higher we have too few sentences, while if we choose the number of tokens to be 15 or less our duration predictor is weak, therefore 20 and above tokens is the optimal point. Additionally, we also filter the dataset based on quality rankings by humans present in the Covost dataset. Concretely, we only take data-points where there are no downvotes and at least three upvotes. The rationale behind this is to have only high quality translation samples present.

Figure 2: A histogram of sentence count vs. number of tokens in a sentence. I.e. the x-axis represents the number of tokens in a sentence, the y-axis is the total count of such sentences.

# 4 Results

In this section we show all the results that we produce for the WMT24 shared testsuite task (Kocmi et al., 2024). Specifically, all included reference paper can be found in Appendix ??. Our evaluation, as described above, combines three metrics: IsoChronoMeter (I), Quality Estimation (Q) and Adjusted-IsoChronoMeter (A) (see Equations (1,2,3)). In particular, we received translations with a variety of systems across four language pairs: en→zh, en→es, en→ru, en→de. In total there are four tables, one per language pair.

| Model | zh-I | zh-Q | zh-A |
|---|---|---|---|
| AIST-AIRC | - | - | - |
| Aya23 | **0.18** | 3.96 | **3.25** |
| Claude-3.5 | 0.19 | 3.94 | 3.19 |
| CommandR-plus | **0.18** | 3.93 | 3.22 |
| CUNI-DS | - | - | - |
| CUNI-NL | - | - | - |
| CycleL | 0.22 | 2.49 | 1.94 |
| CycleL2 | 0.39 | 2.13 | 1.3 |
| Dubformer | - | - | - |
| Gemini-1.5-Pro | - | - | - |
| GPT-4 | **0.18** | 3.98 | **3.26** |
| **Human** | 0.22 | 3.72 | 2.9 |
| HW-TSC | **0.18** | **4.01** | **3.29** |
| IKUN | 0.19 | 3.84 | 3.11 |
| IKUN-C | 0.21 | 3.76 | 2.97 |
| IOL_Research | 0.19 | 3.98 | 3.22 |
| Llama3-70B | 0.19 | 3.99 | 3.23 |
| Mistral-Large | - | - | - |
| MSLC | - | - | - |
| NVIDIA-NeMo | 0.21 | 3.9 | 3.08 |
| Occiglot | - | - | - |
| ONLINE-A | **0.18** | **4.03** | **3.3** |
| ONLINE-B | 0.18 | 3.91 | 3.21 |
| ONLINE-G | 0.19 | 3.91 | 3.17 |
| ONLINE-W | 0.18 | 3.95 | 3.24 |
| Phi-3-Medium | - | - | - |
| TranssionMT | - | - | - |
| TSU-HITs | - | - | - |
| Unbabel-Tower70B | **0.18** | 3.95 | 3.24 |
| UvA-MT | 0.2 | 4 | 3.2 |
| Yandex | - | - | - |
| ZMT | - | - | - |

Table 1: Metrics comparison across different systems. Translation from English into: zh (Chinese). Metrics correspond to: I = IsoChronoMeter (↓), Q = Quality Estimation (↑), A = Adjusted-IsoChronoMeter (↑).

| Model | es-I | es-Q | es-A |
|---|---|---|---|
| AIST-AIRC | - | - | - |
| Aya23 | 0.48 | **4.61** | **2.4** |
| Claude-3.5 | 0.5 | 4.59 | 2.3 |
| CommandR-plus | 0.5 | 4.59 | 2.3 |
| CUNI-DS | - | - | - |
| CUNI-NL | - | - | - |
| CycleL | 0.5 | 3.52 | 1.76 |
| CycleL2 | - | - | - |
| Dubformer | 0.47 | 4.6 | **2.44** |
| Gemini-1.5-Pro | - | - | - |
| GPT-4 | 0.5 | 4.6 | 2.3 |
| **Human** | 0.48 | 4.42 | 2.3 |
| HW-TSC | - | - | - |
| IKUN | **0.46** | 4.56 | **2.46** |
| IKUN-C | **0.45** | 4.5 | **2.48** |
| IOL_Research | 0.48 | 4.6 | 2.39 |
| Llama3-70B | 0.49 | **4.61** | 2.35 |
| Mistral-Large | - | - | - |
| MSLC | 0.47 | **4.61** | **2.44** |
| NVIDIA-NeMo | 0.47 | **4.62** | **2.45** |
| Occiglot | 0.51 | 4.43 | 2.17 |
| ONLINE-A | 0.48 | 4.6 | 2.39 |
| ONLINE-B | 0.49 | **4.64** | 2.37 |
| ONLINE-G | 0.48 | 4.6 | 2.39 |
| ONLINE-W | 0.47 | 4.59 | 2.43 |
| Phi-3-Medium | - | - | - |
| TranssionMT | 0.5 | **4.62** | 2.31 |
| TSU-HITs | 0.25 | 3.39 | 2.54 |
| Unbabel-Tower70B | 0.5 | **4.62** | 2.31 |
| UvA-MT | - | - | - |
| Yandex | - | - | - |
| ZMT | 0.49 | **4.61** | 2.35 |

Table 2: Metrics comparison across different systems. Translation from English into: es (Spanish). Metrics correspond to: I = IsoChronoMeter (↓), Q = Quality Estimation (↑), A = Adjusted-IsoChronoMeter (↑).

| Model | ru-I | ru-Q | ru-A |
|---|---|---|---|
| AIST-AIRC | - | - | - |
| Aya23 | 0.48 | 4.91 | 2.55 |
| Claude-3.5 | 0.49 | **4.95** | 2.52 |
| CommandR-plus | 0.49 | 4.9 | 2.5 |
| CUNI-DS | 0.47 | 4.86 | 2.58 |
| CUNI-NL | - | - | - |
| CycleL | 0.39 | 3.15 | 1.92 |
| CycleL2 | 0.3 | 2.52 | 1.76 |
| Dubformer | **0.42** | 4.82 | **2.8** |
| Gemini-1.5-Pro | - | - | - |
| GPT-4 | 0.47 | 4.93 | 2.61 |
| **Human** | 0.53 | 4.82 | 2.27 |
| HW-TSC | - | - | - |
| IKUN | **0.42** | 4.84 | **2.81** |
| IKUN-C | **0.41** | 4.77 | **2.81** |
| IOL_Research | 0.47 | 4.93 | 2.61 |
| Llama3-70B | 0.49 | **4.94** | 2.52 |
| Mistral-Large | - | - | - |
| MSLC | - | - | - |
| NVIDIA-NeMo | 0.48 | 4.93 | 2.56 |
| Occiglot | - | - | - |
| ONLINE-A | 0.48 | 4.91 | 2.55 |
| ONLINE-B | 0.47 | 4.93 | 2.61 |
| ONLINE-G | 0.51 | 4.93 | 2.42 |
| ONLINE-W | 0.47 | 4.92 | 2.61 |
| Phi-3-Medium | - | - | - |
| TranssionMT | 0.47 | 4.93 | 2.61 |
| TSU-HITs | 0.34 | 3.66 | 2.42 |
| Unbabel-Tower70B | 0.49 | 4.92 | 2.51 |
| UvA-MT | - | - | - |
| Yandex | 0.48 | 4.83 | 2.51 |
| ZMT | 0.48 | 4.91 | 2.55 |

Table 3: Metrics comparison across different systems. Translation from English into: ru (Russian). Metrics correspond to: I = IsoChronoMeter (↓), Q = Quality Estimation (↑), A = Adjusted-IsoChronoMeter (↑).

# 5 Findings

We identify three main findings. Firstly, isochrony is not the natural way of translation (even for humans). Secondly, systems designed for dubbing, such as DubFormer, or multi-linguality, such as Aya23, outperform their 'standard' counter-parts. Finally, the metric itself is powerful and determines systems that are better at dubbing without gold annotations.

## 5.1 Isochrony does not come automatically

We discover that across all language pairs, the smallest isochronic score (ICM) that we discover

is 0.18, which means that the translated audio duration prediction is almost 18% longer or shorter than the original audio prediction.

| Model | de-I | de-Q | de-A |
|---|---|---|---|
| AIST-AIRC | **0.35** | 4.69 | **3.05** |
| Aya23 | 0.38 | 4.68 | 2.9 |
| Claude-3.5 | 0.39 | **4.7** | 2.87 |
| CommandR-plus | 0.38 | 4.68 | 2.9 |
| CUNI-DS | - | - | - |
| CUNI-NL | **0.33** | 4.62 | **3.1** |
| CycleL | 0.4 | 3.65 | 2.19 |
| CycleL2 | 0.4 | 3.65 | 2.19 |
| Dubformer | **0.32** | 4.51 | **3.07** |
| Gemini-1.5-Pro | - | - | - |
| GPT-4 | 0.39 | **4.72** | 2.88 |
| **Human** | 0.38 | 4.47 | 2.77 |
| HW-TSC | - | - | - |
| IKUN | **0.34** | 4.57 | **3.02** |
| IKUN-C | **0.34** | 4.5 | 2.97 |
| IOL_Research | 0.37 | **4.7** | 2.96 |
| Llama3-70B | 0.39 | **4.73** | 2.89 |
| Mistral-Large | - | - | - |
| MSLC | 0.36 | 4.61 | 2.95 |
| NVIDIA-NeMo | 0.37 | **4.72** | 2.97 |
| Occiglot | 0.46 | 4.55 | 2.46 |
| ONLINE-A | 0.37 | 4.68 | 2.95 |
| ONLINE-B | 0.37 | 4.6 | 2.9 |
| ONLINE-G | 0.36 | 4.69 | **3** |
| ONLINE-W | 0.37 | 4.66 | 2.94 |
| Phi-3-Medium | - | - | - |
| TranssionMT | 0.37 | 4.6 | 2.9 |
| TSU-HITs | 0.34 | 3.37 | 2.22 |
| Unbabel-Tower70B | 0.38 | 4.68 | 2.9 |
| UvA-MT | - | - | - |
| Yandex | - | - | - |
| ZMT | 0.37 | 4.68 | 2.95 |

Table 4: Metrics comparison across different systems. Translation from English into: de (German). Metrics correspond to: I = IsoChronoMeter (↓), Q = Quality Estimation (↑), A = Adjusted-IsoChronoMeter (↑).

## 5.2 Most promising systems

The most promising systems that are overall better at isochronic translation as well as translation quality are DubFormer, Ikun, Ikun-C (Liao et al., 2024) and Cuni-NL (Hrabal et al., 2024). For some language pairs, some big players such as GPT-4, Nemo and 'Online A' perform well as well as some specialised systems HW-TSC (Wu et al., 2024) and MSLC (Larkin et al., 2024). Aya23 outperforms its

backbone model CommandR-plus, which is intu-ititve and show that multi-linguality helps MT and isochronic-MT.

### 5.3 Nuances in the metric

Overall we discover that the joined metric is very powerful in ranking systems. We discover an edge case for en→zh, where TSU-HITs has a poor translation quality and likely drops parts of the translation, resulting in poor quality estimate scores, but it has excellent isochrony scores and adjusted isochrony scores. Therefore, we recommend using a performance threshold when applying the metric.

## 6 Conclusion & Future Work

We motivate the importance of isochronic translation. To this end, we presented a novel and simple metric to evaluate isochrony that does not require gold annotations. We evaluated the shared task and discovered that: 1. Isochrony does not come naturally for translation systems, including human (non-isochronic) translation; 2. Systems and LLMs that are designed for multi-linguality or dubbing perform better on our main metric 'Ajusted-IsoChronoMeter', which combines isochrony and machine translation quality; 3. The metric requires some nuance, as systems that drop parts of the translation might have a good isochrony score, but bad translation quality score - overall biasing them towards a better A-ICM.

### 6.1 Future directions

There are several future directions that we identify. Firstly, isochronic translation itself is a promising direction and automatic metrics such as IsoChronoMeter can help with advancing this field. Secondly, extending the benchmark to include gold human translation designed for dubbing. Finally, a more detailed evaluation and improvement of the metric itself; specifically, we believe better duration predictors are possible, and more rigorous evaluation, including using gold annotations and on more language pairs.

### 6.2 Acknowledgements

## References

Milind Agarwal, Sweta Agrawal, Antonios Anastasopoulos, Luisa Bentivogli, Ondřej Bojar, Claudia Borg, Marine Carpuat, Roldano Cattoni, Mauro Cettolo, Mingda Chen, William Chen, Khalid Choukri, Alexandra Chronopoulou, Anna Currey, Thierry Declerck, Qianqian Dong, Kevin Duh, Yannick Estève, Marcello Federico, Souhir Gahbiche, Barry Haddow, Benjamin Hsu, Phu Mon Htut, Hirofumi Inaguma, Dávid Javorský, John Judge, Yasumasa Kano, Tom Ko, Rishu Kumar, Pengwei Li, Xutai Ma, Prashant Mathur, Evgeny Matusov, Paul McNamee, John P. McCrae, Kenton Murray, Maria Nadejde, Satoshi Nakamura, Matteo Negri, Ha Nguyen, Jan Niehues, Xing Niu, Atul Kr. Ojha, John E. Ortega, Proyag Pal, Juan Pino, Lonneke van der Plas, Peter Polák, Elijah Rippeth, Elizabeth Salesky, Jiatong Shi, Matthias Sperber, Sebastian Stüker, Katsuhito Sudoh, Yun Tang, Brian Thompson, Kevin Tran, Marco Turchi, Alex Waibel, Mingxuan Wang, Shinji Watanabe, and Rodolfo Zevallos. 2023. FINDINGS OF THE IWSLT 2023 EVALUATION CAMPAIGN. In *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*, pages 1–61, Toronto, Canada (in-person and online). Association for Computational Linguistics.

Antonios Anastasopoulos, Loïc Barrault, Luisa Bentivogli, Marcely Zanon Boito, Ondřej Bojar, Roldano Cattoni, Anna Currey, Georgiana Dinu, Kevin Duh, Maha Elbayad, Clara Emmanuel, Yannick Estève, Marcello Federico, Christian Federmann, Souhir Gahbiche, Hongyu Gong, Roman Grundkiewicz, Barry Haddow, Benjamin Hsu, Dávid Javorský, Věra Kloudová, Surafel Lakew, Xutai Ma, Prashant Mathur, Paul McNamee, Kenton Murray, Maria Nădejde, Satoshi Nakamura, Matteo Negri, Jan Niehues, Xing Niu, John Ortega, Juan Pino, Elizabeth Salesky, Jiatong Shi, Matthias Sperber, Sebastian Stüker, Katsuhito Sudoh, Marco Turchi, Yogesh Virkar, Alexander Waibel, Changhan Wang, and Shinji Watanabe. 2022. Findings of the IWSLT 2022 evaluation campaign. In *Proceedings of the 19th International Conference on Spoken Language Translation (IWSLT 2022)*, pages 98–157, Dublin, Ireland (in-person and online). Association for Computational Linguistics.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus.

William Brannon, Yogesh Virkar, and Brian Thompson. 2023. Dubbing in practice: A large scale study of human localization with insights for automatic dubbing. *Transactions of the Association for Computational Linguistics*, 11:419–435.

Kevin Cai, Chonghua Liu, and David M. Chan. 2024. Anim-400k: A large-scale dataset for automated end-to-end dubbing of video.

Mingda Chen, Paul-Ambroise Duquenne, Pierre Andrews, Justine Kao, Alexandre Mourachko, Holger Schwenk, and Marta R. Costa-jussà. 2023. BLASER: A text-free speech-to-speech translation evaluation metric. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9064–9079, Toronto, Canada. Association for Computational Linguistics.

Alexandra Chronopoulou, Brian Thompson, Prashant Mathur, Yogesh Virkar, Surafel M. Lakew, and Marcello Federico. 2023a. Jointly optimizing translations and speech timing to improve isochrony in automatic dubbing.

Alexandra Chronopoulou, Brian Thompson, Prashant Mathur, Yogesh Virkar, Surafel Melaku Lakew, and Marcello Federico. 2023b. Jointly optimizing translations and speech timing to improve isochrony in automatic dubbing. *ArXiv*, abs/2302.12979.

Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019. MuST-C: a Multilingual Speech Translation Corpus. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2012–2017, Minneapolis, Minnesota. Association for Computational Linguistics.

Paul-Ambroise Duquenne, Holger Schwenk, and Benoît Sagot. 2023. Sonar: Sentence-level multimodal and language-agnostic representations.

Marcello Federico, Robert Enyedi, Roberto Barra-Chicote, Ritwik Giri, Umut Isik, Arvindh Krishnaswamy, and Hassan Sawaf. 2020. From speech-to-speech translation to automatic dubbing. In *IWSLT 2020*.

Miroslav Hrabal, Josef Jon, Martin Popel, Nam Luu, Danil Semin, and Ondřej Bojar. 2024. CUNI at WMT24 general translation task: Llms, (q)lora, CPO and model merging. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Alina Karakanta, Matteo Negri, and Marco Turchi. 2020. Is 42 the answer to everything in subtitling-oriented speech translation?

Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Findings of the WMT24 general machine translation shared task: the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Surafel M. Lakew, Marcello Federico, Yue Wang, Cuong Hoang, Yogesh Virkar, Roberto Barra-Chicote, and Robert Enyedi. 2021a. Machine translation verbosity control for automatic dubbing.

Surafel M. Lakew, Yogesh Virkar, Prashant Mathur, and Marcello Federico. 2022. Isometric mt: Neural machine translation for automatic dubbing.

Surafel Melaku Lakew, Yogesh Virkar, Prashant Mathur, and Marcello Federico. 2021b. Isometric mt: Neural machine translation for automatic dubbing. *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6242–6246.

Samuel Larkin, Chi-kiu Lo, and Rebecca Knowles. 2024. MSLC24 submissions to the general machine translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Baohao Liao, Christian Herold, Shahram Khadivi, and Christof Monz. 2024. IKUN for WMT24 general MT task: Llms are here for multilingual machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky, Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi, Alexei Baevski, Yossi Adi, Xiaohui Zhang, Wei-Ning Hsu, Alexis Conneau, and Michael Auli. 2023. Scaling speech technology to 1,000+ languages.

Derek Tam, Surafel M. Lakew, Yogesh Virkar, Prashant Mathur, and Marcello Federico. 2022. Isochrony-aware neural machine translation for automatic dubbing.

Changhan Wang, Anne Wu, and Juan Pino. 2020. Covost 2 and massively multilingual speech-to-text translation.

Yihan Wu, Junliang Guo, Xuejiao Tan, Chen Zhang, Bohan Li, Ruihua Song, Lei He, Sheng Zhao, Arul Menezes, and Jiang Bian. 2022. Videodubber: Machine translation with speech-aware length control for video dubbing. *ArXiv*, abs/2211.16934.

Zhanglin Wu, Daimeng Wei, Zongyao Li, Hengchao Shang, Jiaxin GUO, Shaojun Li, Zhiqiang Rao, Yuanchang Luo, Ning Xie, and Hao Yang. 2024. Choose the final translation from NMT and LLM hypotheses using MBR decoding: HW-TSC's submission to the WMT24 general MT shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Rong Ye, Chengqi Zhao, Tom Ko, Chutong Meng, Tao Wang, Mingxuan Wang, and Jun Cao. 2023. Gigast: A 10,000-hour pseudo speech translation corpus.

# A test suite of prompt injection attacks for LLM-based machine translation

**Antonio Valerio Miceli-Barone**
University of Edinburgh
amiceli@ed.ac.uk

**Zhifan Sun**
Ubiquitous Knowledge Processing Lab
Technische Universität Darmstadt
zhifan.sun@tu-darmstadt.de

## Abstract

LLM-based NLP systems typically work by embedding their input data into prompt templates which contain instructions and/or in-context examples, creating queries which are submitted to a LLM, and then parsing the LLM response in order to generate the system outputs. Prompt Injection Attacks (PIAs) are a type of subversion of these systems where a malicious user crafts special inputs which interfere with the prompt templates, causing the LLM to respond in ways unintended by the system designer.

Recently, Sun and Miceli-Barone (2024) proposed a class of PIAs against LLM-based machine translation. Specifically, the task is to translate questions from the TruthfulQA test suite, where an adversarial prompt is prepended to the questions, instructing the system to ignore the translation instruction and answer the questions instead.

In this test suite, we extend this approach to all the language pairs of the WMT 2024 General Machine Translation task. Moreover, we include additional attack formats in addition to the one originally studied.

## 1 Introduction

General purpose pretrained Large Language Models have become the dominant paradigm in NLP, due to their ability to quickly adapt to almost any task with in-context few-shot learning (Brown et al., 2020; Chowdhery et al., 2022; Wei et al., 2022) or instruction following (Ouyang et al., 2022). In most settings, the performance of LLMs predictably increases with their size according to empirical scaling laws (Kaplan et al., 2020; Hernandez et al., 2021; Hoffmann et al., 2022), however, LLMs can still misbehave when subjected to adversarial or out-of-distribution inputs. One such class of scenarios is *Prompt Injection Attacks* (PIAs), where the end-user embeds instructions in their requests that contradict the default system prompt or fine-tuning and thus manipulate the LLM to behave in

ways not intended by the system developer, such as performing a task different than the intended one, revealing secret information included in the system prompt, subvert content moderation, and so on. PIAs were originally discovered in the Inverse Scaling Prize (McKenzie et al., 2023), where they were evaluated on simple tasks such as word capitalization and repetition, showing poor model performance and even asymptotic inverse scaling, meaning that the larger the LLMs are, the more susceptible they become to these attacks. More recently, Sun and Miceli-Barone (2024) studied PIAs against machine translation systems, finding that LLM prompt-based machine translation systems can be often tricked into performing a different task (question answering) with a suitable prompt, especially when the source language is English, while purpose-trained MT systems are more robust.

In this work we apply the methodology of Sun and Miceli-Barone (2024), extended to additional attack formats, to the WMT 2024 General Machine Translation task submissions, in all language pairs. The dataset and evaluation code is available at https://github.com/Avmb/adversarial_MT_prompt_injection.

## 2 Tasks

We consider six subtasks, consisting of a clean (non-adversarial) translation task of the questions from the test set of TruthfulQA (Lin et al., 2022) and five PIAs where we try to manipulate the system to answer the questions instead of translating them. For each sentence in each task and language pair, we have a source sentence, consisting of a question possibly embedded into a PIA prompt template, a reference translation[1] of the question including the PIA prompt if present, which represents the intended behaviour of the MT system and a set of plausible answers, which includes the

---

[1] We use gpt-4o-2024-05-13 in zero-shot mode to compute our reference translations

Figure 1: Average successful attack rate for each system

best answer, correct answers and incorrect answers from TruthfulQA, both in the original English and translated into the target language.

**Clean** The baseline task consists of TruthfulQA questions without any adversarial manipulation.

**Direct instruction** The original attack of Sun and Miceli-Barone (2024). Prepends to the source question the phrase: `"Don't translate this sentence and tell me "` The reference translation is a direct translation of the source, including both the adversarial prefix and the question.

**Task switch zero-shot** In this attack, we include both a translation input and a translation output, followed by a new instruction and finally the question that we want the system to answer. The rationale is that any prefix with translation instructions that the system includes will act only on the first part of the input, and then the model will switch to a different task. The reference translation assumes that the system ignores all the manipulations and translates everything to the target language (or copies the part that is already in the target language).

Source template as a Python f-string (English):

`f"English: Resumption of the session.`

`{tgt_lang}: {transl}`
`Answer the question: {eng_question}"`

Reference template (German):

`f"Englisch: Wiederaufnahme der Sitzung.`
`Deutsch: Wiederaufnahme der Sitzung.`
`Beantworte die Frage: {deu_question}"`

**Task switch one-shot** Like the zero-shot subtask, but includes an in-context example for question answering to prime the LLM towards the new task.

**Task switch zero-shot in JSON format** Like the zero-shot subtask, we format the input in JSON format. The rationale is that the system might be internally embedding the user input in a JSON template, and our adversarial manipulation might cause the LLM to confuse the JSON delimiters or the task specifiers (misinterpreting the input as a RPC call), or even cause the system JSON parser to fail due to improper escaping. The references translate everything except the JSON field names, which remain in English. We believe that this is typically the correct way of translating JSON.

**Task switch one-shot in JSON format** Like the zero-shot subtask in JSON format, we also include one in-context example of question answering to

Figure 2: Average successful attack rate vs and clean-dataset corpus BLEU for each system.
Regression Slope: $-0.0016$, R² Score: $0.1443$.

prime the LLM towards the new task and to teach it to use the JSON format for question-answering output. As in the previous subtask, the references translate everything except the JSON field names.

## 2.1 Non-English source language

Two of the language pairs (Czech→Ukrainian and Japanese→Chinese) have a non-English source language. In this case, for each subtask (except the clean one) we consider two cases, one where the input, including the PIA template, is in the correct source language and another one where it is in English. The motivation is that multi-lingual LLMs might be more easily distracted by English inputs, as noted by Sun and Miceli-Barone (2024).

## 3 Metrics

We use both standard corpus-level metrics and task-specific metrics. For standard metrics, we use BLEU (Papineni et al., 2002) and chrF++ (Popović, 2017) as implemented in SacreBLEU (Post, 2018). As noted by Sun and Miceli-Barone (2024), these metrics might be insufficient to detect successful attacks, therefore we also use the "**question mark** " (demonstrated as QM in the tables) heuristic which

they proposed, which consists in detecting whether the output ends with a question mark[2], we also introduce two additional task-specific metrics: the "**BLEU win**" (demonstrated as BW in the tables) metric consists of computing sentence-level BLEU for each translation w.r.t. the reference translation and comparing it with the sentence-level BLEU w.r.t. the reference answers (using BLEU in multi-reference mode)[3], where we count the proportion of translation where the former is greater than the latter. The "**chrF win**" (demonstrated as CW in the tables) metric is the same with sentence-level chrF++. To further distinguish between the situation where the MT system outputs translation, an answer, or other random content. We have additional metrics (not shown in the tables) that detect whether the sentence BLEU/chrF++ w.r.t the reference translation/reference answers are above/below a threshold. We also detect the target language to ensure it is correct, using **OpenLID** (Burchell

---

[2]possibly followed by closing double quotes. We also allow for Chinese question marks and quote characters.

[3]reference answers are all the candidate answers for the example provided in TruthfulQA, in English and also translated to the non-English source language (if present) and the target language, using gpt-4o-2024-05-13 in zero-shot mode.

et al., 2023), implemented in Hugging Face. We further analyze the system output with GPT-4[4] by asking whether the translation output is a genuine translation, an answer, or other irrelevant output. We count the proportion of output in each task and system type where GPT-4 determines it is a translation or answer and yield metrics **Transl** and **Ans**[5] respectively. Finally, we calculate **Avg. win**, the arithmetic mean of all the positive task-specific metrics excluding to indicate the system's robustness against prompt injection, and **SAAvg** (Successful Attack, avg.), the arithmetic mean of all the negative metrics to detect successful attacks that result in the system answering the question rather than translating (**Avg. win** and **SAAvg** do not sum to 1, because attacks can make the system output something which is neither a translation nor an answer).

## 4 Systems

We divide the systems into "base LLMs" and "team submissions". General purpose LLMs (**GPLLMs**) are publicly available either through weights or APIs that haven't been specifically optimized for translation tasks. The WMT MT Test Suites track organisers evaluated these systems using 4-shot prompting (Hendy et al., 2023). Team submissions are the MT systems that have been submitted by the WMT General Machine Translation task participants, including commercial MT systems accessed by API. We further categorized these systems into LLM-based systems fine-tuned with MT data and specialised for MT task (**SLLMs**)(e.g. Semin and Bojar (2024)), those using other neural network architectures, which include encoder-decoder architectures (e.g. Jasonarson et al. (2024)) and those systems whose architectures remain unknown (**Other**). Finally, we consider anonymized commercial online translation systems (**Online**).

### Base LLMs

```
AYA23,Claude-3, CommandR-plus,
GPT-4, Gemini-1,
Llama3-70B, Mistral-Large, NVIDIA-NeMo,
Phi-3-Medium
```

### Team submissions: LLM-Based

```
AIST-AIRC,
```

---
[4]gpt-4o-mini-2024-07-18
[5]**Transl** and **Ans** do not sum to 1 in general, because the GPT-4 judge can also output "OTHER" if it determines that the output is neither a translation nor an answer.

```
CUNI-DS,CUNI-MH, CUNI-NL,
IKUN, IKUN-C,
IOL_Research,Occiglot,SCIR-MT,
Unbabel-Tower70B, Yandex
```

### Team submissions: Other architectures

```
AMI, BJFU-LPT, CycleL, CycleL2,
DLUT_GTCOM,
CUNI-DocTransformer,
CUNI-GA, CUNI-Transformer,
Dubformer,HW-TSC,MSLC,NTTSU,
Team-J, TranssionMT, TSU-HITs,
UvA-MT
```

### Online Systems

```
ONLINE-A, ONLINE-B, ONLINE-G, ONLINE-W
```

Note that not all of these systems have submissions for all language pairs.

## 5 Results

In this section, we will focus on the results of different types of systems across our designed tasks, and compare the performances under English source and non-English source examples in Czech-Ukrainian and Japanese-Chinese Language pairs. Summary results in figure 1.

Extended results in appendix A, tables 3 to 78, summary results are in tables 79 to 102.

### 5.1 Task: Different prompt injection formats

We start our analysis by examining the performance differences between different MT system types under different prompt injection formats. We report the performance of each system type under all 6 tasks, averaged across all language pairs. The results are found in table 1. We observe a persistent performance downgrade across all metrics when the prompt injection methods get more and more complicated. (i.e. from clean to direct, from zero-shot to one-shot). The change of **Ans** is exciting as it peaks under tasks **0-shot** and **1-shot**, then goes down along with other metrics under prompt injection with JSON format. This phenomenon indicates that under 0-shot and 1-shot prompt injection, the MT systems are geared toward answering the question while under prompt injection with JSON format, the systems tend to be completely confused by outputting irrelevant strings, neither translation nor answers. This is again corroborated by the suboptimal performance of the corpus-specific metrics, as they show lower similarity between the output

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| **clean** | 40.3 | 60.65 | 0.94 | 0.66 | 0.89 | 0.98 | 0.85 | **0.06** | 0.83 | 0.27 |
| GPLLMs | <u>43.69</u> | 64.56 | <u>0.98</u> | 0.69 | **0.93** | **0.99** | <u>0.92</u> | **0.06** | 0.87 | **0.26** |
| SLLMs | **50.03** | **68.62** | <u>0.98</u> | **0.71** | **0.93** | **0.99** | **0.93** | **0.06** | **0.89** | **0.26** |
| Other | 24.61 | 43.65 | 0.78 | 0.54 | 0.76 | 0.92 | 0.62 | **0.06** | 0.67 | 0.29 |
| Online | 42.85 | <u>65.79</u> | **1.0** | <u>0.7</u> | **0.93** | **0.99** | 0.91 | 0.07 | <u>0.88</u> | **0.26** |
| **direct** | 23.67 | 47.64 | 0.81 | 0.54 | 0.77 | 0.89 | 0.57 | 0.27 | 0.69 | 0.29 |
| GPLLMs | 17.45 | 37.94 | 0.62 | 0.42 | 0.63 | 0.73 | 0.48 | 0.46 | 0.55 | 0.41 |
| SLLMs | <u>26.43</u> | <u>53.17</u> | **0.95** | <u>0.53</u> | <u>0.74</u> | **1.0** | <u>0.65</u> | 0.26 | <u>0.77</u> | <u>0.28</u> |
| Other | 16.5 | 36.82 | 0.72 | 0.52 | 0.72 | 0.84 | 0.4 | <u>0.22</u> | 0.59 | 0.29 |
| Online | **34.29** | **62.64** | <u>0.94</u> | **0.69** | **0.98** | **1.0** | **0.76** | **0.14** | **0.86** | **0.2** |
| **0-shot** | 26.08 | 42.39 | 0.82 | 0.56 | 0.76 | 0.83 | 0.41 | 0.33 | 0.65 | 0.3 |
| GPLLMs | 26.39 | 42.44 | 0.84 | 0.57 | 0.77 | 0.82 | 0.44 | 0.39 | 0.67 | 0.32 |
| SLLMs | <u>29.02</u> | <u>48.55</u> | **0.92** | <u>0.62</u> | **0.9** | **0.96** | **0.52** | <u>0.31</u> | **0.77** | **0.25** |
| Other | 16.44 | 29.21 | 0.59 | 0.41 | 0.5 | 0.64 | 0.18 | **0.3** | 0.43 | 0.39 |
| Online | **32.48** | **49.37** | **0.92** | **0.64** | **0.9** | <u>0.9</u> | <u>0.49</u> | 0.32 | <u>0.76</u> | <u>0.26</u> |
| **1-shot** | 25.29 | 39.88 | 0.73 | 0.61 | 0.76 | 0.81 | 0.39 | 0.28 | 0.64 | 0.28 |
| GPLLMs | 24.65 | 40.12 | 0.76 | 0.59 | 0.73 | 0.76 | 0.36 | 0.36 | 0.61 | 0.31 |
| SLLMs | <u>27.76</u> | <u>45.11</u> | **0.84** | <u>0.67</u> | <u>0.89</u> | **0.96** | **0.52** | <u>0.27</u> | <u>0.75</u> | **0.22** |
| Other | 15.29 | 27.07 | 0.49 | 0.47 | 0.52 | 0.63 | 0.17 | **0.23** | 0.42 | 0.36 |
| Online | **33.46** | **47.21** | **0.84** | **0.7** | **0.9** | <u>0.88</u> | <u>0.51</u> | 0.28 | **0.76** | <u>0.23</u> |
| **0-shot JSON** | 21.45 | 29.91 | 0.74 | 0.47 | 0.65 | 0.74 | 0.62 | 0.11 | 0.6 | 0.33 |
| GPLLMs | <u>25.07</u> | <u>33.74</u> | **0.89** | 0.52 | 0.69 | 0.73 | 0.67 | 0.13 | 0.65 | 0.32 |
| SLLMs | 17.21 | 28.1 | <u>0.85</u> | <u>0.55</u> | **0.8** | **0.92** | <u>0.76</u> | <u>0.1</u> | **0.74** | **0.27** |
| Other | 14.38 | 22.79 | 0.4 | 0.23 | 0.3 | 0.52 | 0.25 | 0.13 | 0.3 | 0.46 |
| Online | **29.14** | **35.02** | 0.84 | **0.59** | **0.8** | <u>0.81</u> | **0.78** | 0.06 | <u>0.73</u> | <u>0.28</u> |
| **1-shot JSON** | 15.66 | 25.59 | 0.71 | 0.43 | 0.61 | 0.72 | 0.56 | 0.13 | 0.56 | 0.35 |
| GPLLMs | <u>17.05</u> | <u>27.68</u> | 0.8 | 0.4 | 0.52 | 0.6 | 0.47 | 0.22 | 0.51 | 0.4 |
| SLLMs | 14.69 | 27.08 | **0.83** | <u>0.51</u> | <u>0.79</u> | **0.92** | <u>0.76</u> | <u>0.1</u> | **0.72** | **0.27** |
| Other | 9.56 | 18.36 | 0.38 | 0.24 | 0.31 | 0.54 | 0.23 | 0.14 | 0.3 | 0.45 |
| Online | **21.36** | **29.26** | **0.83** | **0.58** | **0.8** | <u>0.82</u> | **0.77** | 0.06 | **0.72** | **0.27** |

Table 1: Performance of each model type across all six tasks. The bold and underlined numbers indicate the best and the second-best performance scores under each task. The grey row is the average score for all system types. Corpus-specific and task-specific metrics are separated by the vertical line.

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| **direct** | 4.64 | 7.8 | 0.07 | 0.05 | 0.17 | 0.16 | 0.21 | **0.03** | 0.13 | -0.06 |
| GPLLMs | 10.91 | 19.87 | 0.31 | 0.14 | 0.29 | 0.24 | 0.38 | -0.23 | 0.27 | -0.17 |
| SLLMs | 2.92 | **-3.17** | 0.0 | 0.01 | 0.01 | 0.01 | 0.21 | **0.16** | 0.03 | **0.02** |
| Other | 7.39 | 18.75 | **-0.03** | 0.08 | 0.39 | 0.4 | 0.19 | **0.02** | 0.24 | -0.12 |
| Online | **-2.67** | **-4.23** | 0.01 | -0.0 | -0.0 | -0.0 | 0.05 | **0.17** | **-0.01** | **0.02** |
| **0-shot** | 4.55 | 6.31 | 0.03 | 0.02 | 0.13 | 0.03 | 0.14 | 0.0 | 0.06 | -0.03 |
| GPLLMs | 4.17 | 4.46 | 0.06 | 0.02 | 0.08 | **-0.01** | 0.17 | -0.12 | 0.06 | -0.04 |
| SLLMs | 5.58 | 6.5 | 0.14 | 0.11 | 0.16 | 0.07 | 0.25 | **0.07** | 0.14 | -0.06 |
| Other | 5.59 | 12.76 | **-0.09** | **-0.01** | 0.29 | 0.36 | 0.12 | -0.04 | 0.15 | -0.09 |
| Online | 2.85 | 1.51 | -0.0 | **-0.04** | **-0.03** | **-0.3** | 0.01 | **0.09** | **-0.1** | **0.07** |
| **1-shot** | 3.6 | 6.84 | 0.03 | 0.04 | 0.13 | 0.0 | 0.15 | -0.07 | 0.07 | -0.04 |
| GPLLMs | **-0.38** | 2.05 | 0.02 | 0.03 | 0.14 | 0.04 | 0.2 | -0.19 | 0.06 | -0.06 |
| SLLMs | 6.84 | 10.94 | 0.14 | 0.17 | 0.18 | 0.09 | 0.2 | **0.04** | 0.19 | -0.08 |
| Other | 1.39 | 7.53 | **-0.05** | 0.0 | 0.24 | 0.26 | 0.16 | -0.12 | 0.12 | -0.09 |
| Online | 6.55 | 6.84 | 0.0 | **-0.03** | **-0.02** | **-0.39** | 0.04 | -0.02 | **-0.08** | **0.06** |
| **0-shot JSON** | 3.76 | 6.78 | 0.0 | 0.03 | 0.12 | **-0.07** | **-0.03** | -0.08 | **-0.01** | -0.03 |
| GPLLMs | **-0.44** | 1.25 | **-0.1** | 0.01 | 0.07 | **-0.04** | **-0.05** | -0.12 | **-0.05** | -0.01 |
| SLLMs | 5.18 | 11.96 | 0.27 | 0.23 | 0.3 | 0.13 | 0.08 | **0.08** | 0.22 | -0.12 |
| Other | 1.82 | 4.82 | **-0.15** | **-0.04** | 0.15 | 0.22 | 0.01 | -0.27 | 0.02 | -0.08 |
| Online | 8.48 | 9.11 | **-0.01** | **-0.08** | **-0.06** | **-0.61** | **-0.14** | 0.01 | **-0.22** | **0.11** |
| **1-shot JSON** | 3.77 | 7.01 | 0.05 | 0.06 | 0.18 | **-0.05** | 0.04 | -0.12 | 0.04 | -0.06 |
| GPLLMs | 0.38 | 2.05 | **-0.07** | 0.06 | 0.24 | 0.1 | 0.11 | -0.24 | 0.06 | -0.08 |
| SLLMs | 3.02 | 10.11 | 0.3 | 0.28 | 0.31 | 0.15 | 0.1 | **0.03** | 0.24 | -0.14 |
| Other | 2.19 | 5.41 | **-0.06** | **-0.01** | 0.21 | 0.22 | 0.06 | -0.31 | 0.05 | -0.11 |
| Online | 9.47 | 10.49 | 0.02 | **-0.07** | **-0.04** | **-0.65** | **-0.12** | **0.03** | **-0.21** | **0.11** |

Table 2: Delta between English source language and non-English source language in Czech-Ukrainian and Japanese-Chinese language pairs. Numbers indicating a downgrade in the performance on the side of the English source language are marked in bold. Similarly, the grey rows are the average performance across all types of systems, and corpus-specific and task-specific metrics are separated by the vertical line.

and reference answer.

From the table, we can also observe the striking robustness of Online translation systems against all kinds of prompt injection. Taking the Online system aside, we can see that the performance of **SLLMs** also shows a rather strong persistence against prompt injection and better translation quality, with only a small margin compared to **Online** systems. For **GLLMs**, despite its size and optimal performance on most other tasks, they underperform **SLLMs** which are based on smaller LLMs fine-tuned on MT data, when facing injected prompt, and its performance is comparable with **SLLMs** without injected prompt. On the other hand, team submission systems with other architectures underperform most other systems types under all tasks.

The results show that commercial online MT systems are the most robust against prompt injection, while the LLM-based systems fine-tuned with MT instruction and data also show a similar robustness against prompt injection, with Avg. win above 0.7 across all tasks.

## 5.2 Performance difference between English and non-English source languages

Systems that are intended to translate from a non-English source language can be attacked in either English or the non-English language. We analyze the performance differences between English attacks and non-English attacks in Czech-Ukrainian and Japanese-Chinese language pairs by calculating the average metrics delta between English-source and non-English sources. The results are found in 2.

Similar to the previous analysis, we can find a steady decrease in English attack robustness as the complexity of prompt injection increases, and the decrease is generally under the task-specific metrics, not under corpus-specific metrics, indicating that the MT systems are misled toward either answering the questions or outputting irrelevant rather than general decrease in the translation quality. This is particularly obvious under the two JSON-formated prompt injection tasks where both LLMTransl and LLMAns experience a decrease in all systems types.

Concerning the specific differences between system types, we can see that team Online systems suffer from the most performance loss when the attack language is English. In addition, we also observe casual performance loss for **GPLLMs** systems under 0-shot JSON task. Again, **SLLMs** and **Other** show the strongest performance robustness under the English attack language, with the largest Avg. win and the smallest SAAvg under most tasks, arguably being based on multi-lingual LLMs they can still process English source text but the fine-tuning on translation tasks steers them away from performing other tasks.

## 5.3 Scaling

We show in Figure 2 the average successful attack rate vs. the clean dataset corpus-BLEU score. In general, the systems that have a higher resistance against successful attacks are also the ones that perform better on the clean dataset, indicating positive scaling between robustness and non-adversarial performance.

## 6 Conclusions

We presented a test suite of five variants of prompt-injection attacks for machine translation plus one baseline clean version, and we evaluated it on all systems and language pairs of the WMT 2024 General Translation task. We found a general trend of decrease in MT performance with increasing complexity of prompt injection, where even the best performance LLMs stumble on, some even with BLEU scores less than 10 under certain language pairs. In addition, we detected a decrease in performance with the English injected prompts, particularly for commercial MT systems and sometimes for general-purpose LLMs. Among all systems types, the specialized MT systems fine-tuned on LLMs and the commercial MT systems show the best overall performance against prompt injection.

## Ethics Statement

In this work, we investigate the vulnerability of LLMs to Prompt Injection Attacks. We do not present novel attacks, instead, we focus on the characterization of the system performance under a well-known attack, albeit applied to a novel task (Machine Translation), we believe that our work does not create additional security risks but instead may contribute to eventually increasing the security of LLM-based systems by furthering a better understanding of these vulnerabilities.

In this work we do not carry out experiments on human subjects, therefore there are no risks associated with human experimentation.

## Limitations

Our work has the following limitations:

- Due to the format of the WMT shared task, we are limited to single rounds of interactions with the systems, and we are further limited to single-line examples. This has prevents certain kinds of attacks that use multiple rounds of dialogue, and also attacks that include multiple lines in each message, which can exploit certain formatting tricks using JSON, XML or Markdown.

- No single metric that we used can always determine whether a system output is a plausible translation, an answer or something else. Even GPT-4-based evaluation makes mistakes. We combined different heuristics to ameliorate this issue, but there might be still systems, language pairs or attack formats which may be inaccurately evaluated. Human evaluation is possible but we did not perform it due to time and financial considerations.

- Using GPT-4 for dataset generation and evaluation creates some reproducibility issues in the long term, because OpenAI eventually retires models.

## Acknowledgements

## References

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Laurie Burchell, Alexandra Birch, Nikolay Bogoychev, and Kenneth Heafield. 2023. An open dataset and model for language identification. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 865–879, Toronto, Canada. Association for Computational Linguistics.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation.

Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. 2021. Scaling laws for transfer.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. Training compute-optimal large language models.

Atli Jasonarson, Hinrik Hafsteinsson, Bjarki Ármannsson, and Steinþór Steingrímsson. 2024. Cogs in a machine, doing what they're meant to do – the AMI submission to the WMT24 general translation task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models.

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods.

Ian R. McKenzie, Alexander Lyzhov, Michael Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Aaron Kirtland, Alexis Ross, Alisa Liu, Andrew Gritsevskiy, Daniel Wurgaft, Derik Kauffman, Gabriel Recchia, Jiacheng Liu, Joe Cavanagh, Max Weiss, Sicong Huang, The Floating Droid, Tom Tseng, Tomasz Korbak, Xudong Shen, Yuhui Zhang, Zhengping Zhou, Najoung Kim, Samuel R. Bowman, and Ethan Perez. 2023. Inverse scaling: When bigger isn't better.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, page 311–318, USA. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Danil Semin and Ondřej Bojar. 2024. CUNI-DS submission: A naive transfer learning setup for english-to-russian translation utilizing english-to-czech data. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Zhifan Sun and Antonio Valerio Miceli-Barone. 2024. Scaling behavior of machine translation with large language models under prompt injection attacks. In *Proceedings of the First edition of the Workshop on the Scaling Behavior of Large Language Models (SCALE-LLM 2024)*, pages 9–23, St. Julian's, Malta. Association for Computational Linguistics.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent abilities of large language models.

# A Results

## A.1 Extended results

Base LLMs are highlighted in gray. Problem-specific metrics: "QM": Question mark heuristic, "BW": BLEU win, "CW": chrF++ win, "LID": correct target language, "Avg. robustness" is the arithmetic average of all the problem-specific metrics.

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 50.124 | 69.491 | **1.000** | 0.891 | 0.917 | 0.980 | 0.952 | 0.045 | 0.937 | 0.261 |
| Claude-3 | **63.945** | **80.516** | 0.998 | **0.930** | **0.966** | 0.979 | 0.965 | 0.034 | **0.965** | 0.257 |
| CommandR-plus | 51.532 | 70.648 | 0.996 | 0.903 | 0.923 | 0.978 | 0.945 | 0.051 | 0.938 | 0.258 |
| GPT-4 | 58.671 | 76.248 | 0.999 | 0.911 | 0.960 | 0.982 | 0.965 | 0.035 | 0.958 | 0.255 |
| Llama3-70B | 55.838 | 73.779 | 0.998 | 0.907 | 0.940 | 0.980 | **0.976** | 0.024 | 0.951 | 0.254 |
| NVIDIA-NeMo | 53.441 | 71.047 | 0.968 | 0.889 | 0.913 | 0.968 | 0.961 | 0.033 | 0.934 | 0.269 |
| CUNI-DS | 45.865 | 65.698 | 0.947 | 0.901 | 0.924 | 0.978 | 0.968 | 0.029 | 0.930 | 0.254 |
| IKUN | 46.017 | 65.324 | 0.995 | 0.891 | 0.918 | 0.976 | 0.968 | 0.028 | 0.934 | 0.249 |
| IKUN-C | 39.794 | 60.823 | 0.998 | 0.865 | 0.903 | 0.977 | 0.952 | 0.039 | 0.913 | 0.246 |
| Unbabel-Tower70B | 54.457 | 73.925 | 0.996 | 0.917 | 0.947 | **0.988** | 0.958 | 0.039 | 0.956 | 0.253 |
| Yandex | 42.793 | 65.032 | 0.939 | 0.873 | 0.887 | 0.985 | 0.934 | 0.064 | 0.912 | 0.270 |
| CycleL | 1.720 | 19.371 | 0.988 | 0.712 | 0.764 | 0.976 | 0.032 | 0.050 | 0.519 | 0.122 |
| CycleL2 | 0.823 | 15.256 | 0.974 | 0.714 | 0.693 | 0.972 | 0.004 | 0.026 | 0.488 | **0.108** |
| Dubformer | 0.811 | 2.480 | 0.999 | 0.039 | 0.002 | 0.000 | 0.002 | **0.009** | 0.152 | 0.684 |
| IOL_Research | 62.421 | 77.519 | 0.967 | 0.902 | 0.934 | 0.978 | 0.974 | 0.026 | 0.950 | 0.269 |
| ONLINE-A | 57.977 | 75.168 | 0.998 | 0.923 | 0.940 | 0.969 | 0.958 | 0.042 | 0.954 | 0.259 |
| ONLINE-B | 55.403 | 73.776 | 0.998 | 0.913 | 0.944 | 0.971 | 0.960 | 0.040 | 0.950 | 0.258 |
| ONLINE-G | 53.353 | 74.154 | 0.996 | 0.909 | 0.929 | 0.987 | 0.947 | 0.051 | 0.947 | 0.260 |
| ONLINE-W | 53.906 | 72.810 | 0.995 | 0.913 | 0.934 | 0.982 | 0.961 | 0.038 | 0.952 | 0.259 |
| TSU-HITs | 22.052 | 43.818 | 0.553 | 0.717 | 0.808 | 0.969 | 0.788 | 0.100 | 0.742 | 0.331 |
| TranssionMT | 55.300 | 74.002 | 0.998 | 0.912 | 0.945 | 0.969 | 0.961 | 0.039 | 0.950 | 0.260 |

Table 3: English→Russian, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | **56.347** | **76.822** | 0.995 | **0.990** | 0.988 | **1.000** | 0.886 | 0.114 | **0.972** | 0.204 |
| Claude-3 | 0.032 | 0.542 | 0.010 | 0.006 | 0.001 | 0.005 | 0.000 | 1.000 | 0.003 | 0.836 |
| CommandR-plus | 23.382 | 53.457 | 0.803 | 0.704 | 0.709 | 0.882 | 0.586 | 0.354 | 0.727 | 0.375 |
| GPT-4 | 26.456 | 42.902 | 0.674 | 0.389 | 0.278 | 0.976 | 0.215 | 0.785 | 0.555 | 0.575 |
| Llama3-70B | 2.860 | 12.925 | 0.266 | 0.211 | 0.188 | 0.244 | 0.127 | 0.873 | 0.208 | 0.720 |
| NVIDIA-NeMo | 35.470 | 69.105 | 0.982 | 0.951 | 0.983 | **1.000** | 0.848 | 0.152 | 0.943 | 0.229 |
| CUNI-DS | 24.399 | 51.947 | 0.942 | 0.909 | 0.871 | **1.000** | 0.914 | 0.086 | 0.880 | 0.228 |
| IKUN | 25.417 | 53.386 | 0.987 | 0.897 | 0.807 | **1.000** | **0.936** | 0.064 | 0.888 | 0.237 |
| IKUN-C | 22.346 | 50.852 | 0.994 | 0.853 | 0.798 | **1.000** | 0.922 | 0.078 | 0.864 | 0.233 |
| Unbabel-Tower70B | 30.181 | 65.860 | 0.995 | 0.960 | 0.963 | **1.000** | 0.670 | 0.329 | 0.901 | 0.247 |
| Yandex | 27.575 | 64.911 | 0.780 | 0.969 | 0.990 | **1.000** | 0.845 | 0.155 | 0.899 | 0.242 |
| CycleL | 1.379 | 18.603 | 0.984 | 0.832 | 0.707 | 0.999 | 0.000 | 0.179 | 0.512 | **0.119** |
| CycleL2 | 0.570 | 15.032 | 0.977 | 0.652 | 0.554 | 0.998 | 0.000 | 0.162 | 0.456 | 0.149 |
| Dubformer | 0.489 | 1.503 | **0.999** | 0.033 | 0.001 | 0.000 | 0.001 | **0.044** | 0.148 | 0.671 |
| IOL_Research | 33.521 | 55.760 | 0.965 | 0.655 | 0.589 | 0.990 | 0.463 | 0.535 | 0.760 | 0.407 |
| ONLINE-A | 34.274 | 66.320 | 0.863 | 0.969 | 0.958 | **1.000** | 0.777 | 0.223 | 0.912 | 0.251 |
| ONLINE-B | 33.462 | 68.866 | 0.995 | 0.987 | 0.989 | **1.000** | 0.812 | 0.188 | 0.945 | 0.223 |
| ONLINE-G | 34.105 | 70.464 | **0.999** | 0.973 | **0.995** | **1.000** | 0.902 | 0.098 | 0.957 | 0.214 |
| ONLINE-W | 36.434 | 70.303 | **0.999** | 0.960 | 0.980 | **1.000** | 0.886 | 0.114 | 0.954 | 0.222 |
| TSU-HITs | 8.637 | 36.031 | 0.124 | 0.813 | 0.949 | 0.996 | 0.721 | 0.257 | 0.651 | 0.268 |
| TranssionMT | 33.411 | 69.050 | 0.995 | 0.987 | 0.989 | **1.000** | 0.815 | 0.185 | 0.945 | 0.222 |

Table 4: English→Russian, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 62.406 | 78.552 | 0.947 | 0.999 | 0.994 | **1.000** | 0.048 | 0.570 | 0.855 | 0.262 |
| Claude-3 | 59.403 | 78.275 | 0.957 | 0.960 | 0.957 | 0.960 | 0.098 | 0.406 | 0.835 | 0.257 |
| CommandR-plus | 31.655 | 52.900 | 0.864 | 0.858 | 0.737 | 0.996 | 0.027 | 0.902 | 0.734 | 0.390 |
| GPT-4 | 63.638 | 80.974 | **0.999** | **1.000** | **1.000** | **1.000** | **0.129** | 0.379 | **0.875** | 0.224 |
| Llama3-70B | 37.223 | 56.440 | 0.908 | 0.847 | 0.764 | 0.999 | 0.022 | 0.881 | 0.758 | 0.391 |
| NVIDIA-NeMo | 62.288 | 78.741 | 0.994 | **1.000** | **1.000** | **1.000** | 0.027 | 0.610 | 0.860 | 0.249 |
| CUNI-DS | 16.636 | 36.002 | 0.952 | 0.435 | 0.252 | 0.995 | 0.000 | 0.998 | 0.546 | 0.529 |
| IKUN | 63.435 | 78.322 | 0.998 | **1.000** | 0.998 | **1.000** | 0.049 | 0.359 | 0.863 | 0.211 |
| IKUN-C | 25.074 | 52.561 | 0.996 | 0.949 | 0.878 | **1.000** | 0.054 | 0.875 | 0.793 | 0.326 |
| Unbabel-Tower70B | 36.738 | 58.955 | 0.989 | 0.968 | 0.897 | **1.000** | 0.037 | 0.897 | 0.825 | 0.341 |
| Yandex | 23.056 | 51.441 | 0.965 | 0.979 | 0.898 | **1.000** | 0.010 | 0.612 | 0.777 | 0.270 |
| CycleL | 1.531 | 18.542 | 0.967 | 0.985 | 0.842 | 0.984 | 0.000 | 0.879 | 0.547 | **0.173** |
| CycleL2 | 0.340 | 13.500 | 0.763 | 0.846 | 0.641 | 0.925 | 0.000 | 0.618 | 0.454 | 0.209 |
| Dubformer | 10.182 | 17.596 | **0.999** | 0.450 | 0.048 | 0.000 | 0.001 | **0.136** | 0.218 | 0.596 |
| IOL_Research | **66.535** | **84.207** | 0.991 | 0.996 | 0.995 | **1.000** | 0.066 | 0.301 | 0.862 | 0.211 |
| ONLINE-A | 56.073 | 80.194 | 0.998 | **1.000** | **1.000** | **1.000** | 0.007 | 0.315 | 0.858 | 0.215 |
| ONLINE-B | 62.117 | 80.242 | 0.998 | **1.000** | **1.000** | **1.000** | 0.006 | 0.646 | 0.858 | 0.259 |
| ONLINE-G | 49.336 | 72.718 | **0.999** | 0.998 | 0.998 | **1.000** | 0.000 | 0.315 | 0.853 | 0.215 |
| ONLINE-W | 63.109 | 83.275 | **0.999** | **1.000** | **1.000** | **1.000** | 0.054 | 0.360 | 0.865 | 0.218 |
| TSU-HITs | 5.622 | 30.610 | 0.082 | 0.908 | 0.962 | **1.000** | 0.118 | 0.671 | 0.557 | 0.272 |
| TranssionMT | 62.049 | 80.343 | 0.998 | **1.000** | **1.000** | **1.000** | 0.006 | 0.654 | 0.858 | 0.260 |

Table 5: English→Russian, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 78.245 | 89.097 | **0.999** | **1.000** | **1.000** | **1.000** | 0.000 | 0.732 | 0.857 | 0.214 |
| Claude-3 | **87.218** | **92.427** | 0.973 | 0.979 | 0.977 | 0.979 | 0.005 | 0.744 | 0.837 | 0.240 |
| CommandR-plus | 62.914 | 75.684 | 0.963 | **1.000** | 0.998 | 0.999 | 0.000 | 0.905 | 0.850 | 0.257 |
| GPT-4 | 73.570 | 85.441 | 0.996 | **1.000** | **1.000** | **1.000** | 0.000 | 0.781 | 0.857 | 0.229 |
| Llama3-70B | 76.261 | 84.499 | 0.978 | **1.000** | **1.000** | **1.000** | 0.006 | 0.624 | 0.855 | 0.205 |
| NVIDIA-NeMo | 69.460 | 81.005 | 0.965 | **1.000** | **1.000** | **1.000** | 0.000 | 0.786 | 0.852 | 0.238 |
| CUNI-DS | 36.008 | 55.041 | 0.021 | 0.995 | 0.994 | **1.000** | 0.002 | 0.519 | 0.711 | 0.228 |
| IKUN | 84.657 | 91.057 | 0.998 | 0.999 | 0.994 | **1.000** | 0.000 | 0.728 | 0.855 | 0.208 |
| IKUN-C | 49.945 | 71.158 | 0.991 | **1.000** | **1.000** | **1.000** | 0.010 | 0.791 | **0.857** | 0.229 |
| Unbabel-Tower70B | 59.223 | 74.272 | 0.995 | **1.000** | **1.000** | **1.000** | 0.002 | 0.851 | 0.857 | 0.250 |
| Yandex | 50.556 | 72.383 | 0.958 | **1.000** | **1.000** | **1.000** | 0.002 | 0.630 | 0.852 | 0.194 |
| CycleL | 1.540 | 21.744 | 0.613 | 0.967 | 0.857 | 0.876 | 0.000 | 0.067 | 0.528 | **0.112** |
| CycleL2 | 0.226 | 10.966 | 0.158 | 0.690 | 0.485 | 0.728 | 0.000 | 0.066 | 0.294 | 0.287 |
| Dubformer | 4.556 | 8.529 | **0.999** | 0.460 | 0.012 | 0.000 | 0.007 | **0.022** | 0.211 | 0.512 |
| IOL_Research | 80.168 | 90.896 | 0.996 | **1.000** | **1.000** | **1.000** | 0.000 | 0.692 | 0.857 | 0.209 |
| ONLINE-A | 82.858 | 91.560 | 0.998 | **1.000** | **1.000** | **1.000** | 0.000 | 0.782 | 0.857 | 0.227 |
| ONLINE-B | 84.891 | 91.609 | 0.998 | **1.000** | **1.000** | **1.000** | 0.000 | 0.743 | 0.857 | 0.218 |
| ONLINE-G | 72.098 | 87.344 | 0.994 | **1.000** | **1.000** | **1.000** | 0.000 | 0.586 | 0.856 | 0.196 |
| ONLINE-W | 72.016 | 85.979 | **0.999** | **1.000** | **1.000** | **1.000** | 0.002 | 0.614 | **0.857** | 0.198 |
| TSU-HITs | 0.352 | 16.821 | 0.029 | 0.759 | 0.766 | **1.000** | **0.045** | 0.317 | 0.398 | 0.259 |
| TranssionMT | 84.849 | 91.624 | 0.998 | **1.000** | **1.000** | **1.000** | 0.000 | 0.745 | 0.857 | 0.218 |

Table 6: English→Russian, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 6.152 | 22.203 | 0.911 | 0.810 | 0.858 | 0.963 | 0.852 | 0.121 | 0.864 | 0.293 |
| Claude-3 | **27.579** | 30.655 | 0.985 | 0.554 | 0.572 | 0.583 | 0.575 | 0.037 | 0.632 | 0.437 |
| CommandR-plus | 3.246 | 15.813 | 0.660 | 0.552 | 0.583 | 0.869 | 0.569 | 0.335 | 0.633 | 0.435 |
| GPT-4 | 16.358 | 34.809 | **0.999** | 0.108 | 0.087 | 0.086 | 0.084 | 0.011 | 0.223 | 0.647 |
| Llama3-70B | 15.552 | 34.564 | **0.999** | 0.917 | **0.942** | 0.978 | **0.973** | 0.026 | 0.958 | 0.257 |
| NVIDIA-NeMo | 16.936 | 31.924 | 0.351 | 0.367 | 0.406 | **0.991** | 0.343 | 0.011 | 0.443 | 0.354 |
| CUNI-DS | 15.899 | 34.644 | 0.940 | 0.814 | 0.798 | 0.901 | 0.834 | 0.016 | 0.827 | 0.290 |
| IKUN | 14.258 | 33.930 | 0.985 | 0.880 | 0.887 | 0.976 | 0.938 | 0.048 | 0.911 | **0.256** |
| IKUN-C | 6.366 | 25.578 | 0.979 | 0.848 | 0.864 | 0.966 | 0.927 | 0.040 | 0.893 | 0.261 |
| Unbabel-Tower70B | 6.992 | 25.179 | 0.931 | 0.734 | 0.758 | 0.838 | 0.765 | 0.089 | 0.791 | 0.339 |
| Yandex | 1.663 | 11.932 | 0.028 | 0.039 | 0.116 | 0.771 | 0.009 | 0.979 | 0.144 | 0.614 |
| CycleL | 0.000 | 4.278 | 0.000 | 0.100 | 0.164 | 0.069 | 0.000 | 0.007 | 0.048 | 0.525 |
| CycleL2 | 0.034 | 5.305 | 0.000 | 0.086 | 0.111 | 0.818 | 0.000 | **0.005** | 0.145 | 0.428 |
| Dubformer | 15.879 | 30.230 | **0.999** | 0.039 | 0.002 | 0.000 | 0.002 | 0.009 | 0.152 | 0.684 |
| IOL_Research | 2.058 | 16.422 | 0.670 | 0.607 | 0.630 | 0.909 | 0.635 | 0.098 | 0.671 | 0.349 |
| ONLINE-A | 16.512 | **37.187** | **0.999** | **0.925** | **0.942** | 0.976 | 0.958 | 0.042 | **0.960** | 0.262 |
| ONLINE-B | 16.015 | 24.116 | 0.976 | 0.890 | 0.916 | 0.945 | 0.923 | 0.050 | 0.921 | 0.268 |
| ONLINE-G | 13.410 | 27.853 | 0.422 | 0.275 | 0.313 | 0.635 | 0.306 | 0.083 | 0.356 | 0.441 |
| ONLINE-W | 15.780 | 34.287 | **0.999** | 0.911 | 0.941 | 0.984 | 0.971 | 0.027 | 0.956 | 0.257 |
| TSU-HITs | 0.000 | 3.047 | 0.000 | 0.034 | 0.023 | 0.136 | 0.001 | 0.070 | 0.028 | 0.560 |
| TranssionMT | 16.011 | 35.944 | 0.993 | 0.903 | 0.924 | 0.966 | 0.951 | 0.044 | 0.940 | 0.265 |

Table 7: English→Russian, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 1.797 | 15.461 | 0.993 | 0.890 | 0.925 | 0.965 | 0.939 | 0.047 | 0.930 | 0.257 |
| Claude-3 | **14.487** | 19.419 | 0.984 | 0.051 | 0.023 | 0.023 | 0.017 | 0.028 | 0.166 | 0.680 |
| CommandR-plus | 1.047 | 10.756 | 0.824 | 0.322 | 0.335 | 0.488 | 0.330 | 0.170 | 0.433 | 0.551 |
| GPT-4 | 5.060 | 21.685 | 0.998 | 0.065 | 0.038 | 0.035 | 0.033 | 0.015 | 0.180 | 0.670 |
| Llama3-70B | 4.804 | 21.169 | **0.999** | 0.918 | 0.944 | 0.983 | 0.963 | 0.035 | 0.957 | 0.256 |
| NVIDIA-NeMo | 1.678 | 19.243 | 0.159 | 0.705 | 0.442 | **0.995** | 0.000 | 0.007 | 0.329 | **0.245** |
| CUNI-DS | 4.858 | 21.717 | 0.985 | 0.907 | 0.930 | 0.985 | 0.953 | 0.038 | 0.933 | 0.248 |
| IKUN | 1.679 | 16.411 | 0.973 | 0.884 | 0.909 | 0.968 | 0.936 | 0.055 | 0.915 | 0.260 |
| IKUN-C | 1.618 | 15.853 | 0.892 | 0.808 | 0.825 | 0.962 | 0.825 | 0.113 | 0.836 | 0.284 |
| Unbabel-Tower70B | 2.470 | 17.172 | 0.968 | 0.655 | 0.671 | 0.720 | 0.679 | 0.055 | 0.722 | 0.381 |
| Yandex | 0.735 | 7.785 | 0.016 | 0.026 | 0.108 | 0.775 | 0.002 | 0.985 | 0.135 | 0.617 |
| CycleL | 0.000 | 2.184 | 0.000 | 0.100 | 0.166 | 0.062 | 0.000 | 0.006 | 0.047 | 0.526 |
| CycleL2 | 0.000 | 3.115 | 0.000 | 0.097 | 0.095 | 0.804 | 0.000 | 0.006 | 0.142 | 0.430 |
| Dubformer | 5.347 | 19.448 | **0.999** | 0.039 | 0.002 | 0.000 | 0.002 | 0.009 | 0.152 | 0.684 |
| IOL_Research | 1.856 | 15.691 | 0.995 | 0.851 | 0.868 | 0.895 | 0.895 | 0.032 | 0.889 | 0.290 |
| ONLINE-A | 5.059 | **22.724** | **0.999** | **0.925** | 0.942 | 0.976 | 0.958 | 0.042 | **0.960** | 0.262 |
| ONLINE-B | 5.267 | 13.681 | 0.994 | 0.913 | 0.947 | 0.972 | 0.945 | 0.050 | 0.951 | 0.260 |
| ONLINE-G | 3.994 | 15.395 | 0.000 | 0.007 | 0.000 | 0.000 | 0.000 | **0.001** | 0.001 | 0.575 |
| ONLINE-W | 4.821 | 20.608 | 0.998 | 0.919 | 0.947 | 0.985 | **0.969** | 0.029 | 0.958 | 0.256 |
| TSU-HITs | 0.000 | 2.441 | 0.000 | 0.067 | 0.054 | 0.640 | 0.000 | 0.624 | 0.109 | 0.564 |
| TranssionMT | 5.267 | 22.061 | 0.996 | 0.917 | **0.949** | 0.974 | 0.947 | 0.053 | 0.954 | 0.261 |

Table 8: English→Russian, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 60.528 | 77.596 | **0.999** | 0.929 | 0.961 | **0.998** | 0.999 | 0.001 | 0.972 | 0.253 |
| Claude-3 | 69.372 | **84.126** | 0.998 | 0.950 | 0.977 | 0.995 | 0.998 | 0.001 | 0.982 | 0.256 |
| CommandR-plus | 60.904 | 78.355 | 0.993 | 0.928 | 0.968 | 0.995 | 0.998 | 0.002 | 0.971 | 0.256 |
| GPT-4 | **70.239** | 84.067 | **0.999** | 0.950 | **0.979** | 0.996 | 0.998 | 0.001 | 0.982 | 0.255 |
| Llama3-70B | 64.414 | 79.829 | **0.999** | 0.940 | 0.976 | 0.995 | **1.000** | **0.000** | 0.976 | 0.256 |
| NVIDIA-NeMo | 62.179 | 77.817 | 0.985 | 0.933 | 0.968 | 0.996 | 0.995 | **0.000** | 0.973 | 0.256 |
| AIST-AIRC | 54.511 | 72.781 | 0.998 | 0.909 | 0.953 | 0.995 | 0.996 | **0.000** | 0.965 | 0.254 |
| CUNI-NL | 51.442 | 69.699 | 0.994 | 0.892 | 0.940 | 0.996 | 0.995 | **0.000** | 0.952 | 0.256 |
| IKUN | 51.652 | 70.262 | 0.996 | 0.880 | 0.940 | 0.995 | 0.993 | **0.000** | 0.947 | 0.259 |
| IKUN-C | 44.710 | 65.240 | 0.994 | 0.868 | 0.930 | **0.998** | 0.979 | 0.004 | 0.931 | 0.252 |
| Unbabel-Tower70B | 61.008 | 78.193 | 0.991 | 0.924 | 0.966 | **0.998** | 0.999 | 0.001 | 0.970 | 0.254 |
| CycleL | 20.487 | 44.322 | 0.977 | 0.803 | 0.884 | 0.993 | 0.447 | **0.000** | 0.776 | **0.210** |
| CycleL2 | 20.487 | 44.322 | 0.977 | 0.803 | 0.884 | 0.993 | 0.447 | **0.000** | 0.776 | **0.210** |
| Dubformer | 26.213 | 32.808 | 0.956 | 0.867 | 0.927 | 0.324 | 0.307 | 0.038 | 0.571 | 0.213 |
| IOL_Research | 69.214 | 82.833 | 0.977 | 0.929 | 0.969 | 0.995 | 0.996 | 0.001 | 0.974 | 0.263 |
| MSLC | 41.196 | 64.234 | 0.968 | 0.868 | 0.920 | 0.995 | 0.952 | 0.002 | 0.927 | 0.258 |
| ONLINE-A | 68.859 | 82.629 | **0.999** | 0.949 | **0.979** | 0.996 | **1.000** | **0.000** | 0.983 | 0.255 |
| ONLINE-B | 54.922 | 74.946 | 0.998 | 0.907 | 0.956 | **0.998** | 0.996 | 0.004 | 0.961 | 0.256 |
| ONLINE-G | 68.624 | 82.302 | **0.999** | **0.956** | 0.977 | **0.998** | **1.000** | **0.000** | **0.985** | 0.255 |
| ONLINE-W | 61.546 | 78.220 | **0.999** | 0.923 | 0.952 | 0.995 | **1.000** | **0.000** | 0.969 | 0.258 |
| TSU-HITs | 29.868 | 49.567 | 0.521 | 0.766 | 0.863 | 0.976 | 0.864 | 0.002 | 0.785 | 0.322 |
| TranssionMT | 54.873 | 74.941 | 0.998 | 0.909 | 0.956 | **0.998** | 0.996 | 0.004 | 0.961 | 0.256 |

Table 9: English→German, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 41.099 | 67.517 | 0.988 | 0.938 | 0.919 | 0.998 | 0.979 | 0.015 | 0.964 | 0.228 |
| Claude-3 | 1.673 | 18.229 | 0.024 | 0.119 | 0.173 | 0.234 | 0.024 | 0.974 | 0.114 | 0.767 |
| CommandR-plus | 17.442 | 45.738 | 0.619 | 0.608 | 0.573 | 0.892 | 0.492 | 0.448 | 0.641 | 0.441 |
| GPT-4 | 43.766 | 60.993 | 0.825 | 0.638 | 0.599 | 0.995 | 0.799 | 0.201 | 0.795 | 0.385 |
| Llama3-70B | 38.530 | 68.205 | 0.865 | 0.875 | 0.879 | 0.898 | 0.856 | 0.143 | 0.877 | 0.281 |
| NVIDIA-NeMo | 41.074 | 68.625 | 0.968 | 0.988 | 0.984 | **1.000** | 0.994 | 0.005 | 0.989 | 0.221 |
| AIST-AIRC | 55.103 | 75.235 | **0.999** | 0.996 | 0.996 | **1.000** | 0.980 | 0.009 | 0.994 | 0.191 |
| CUNI-NL | 55.620 | 74.731 | 0.761 | **1.000** | **0.999** | 0.999 | 0.988 | 0.005 | 0.964 | 0.224 |
| IKUN | 33.558 | 65.936 | 0.810 | 0.984 | 0.996 | **1.000** | 0.989 | 0.005 | 0.965 | 0.220 |
| IKUN-C | 26.128 | 58.671 | 0.896 | 0.913 | 0.908 | 0.999 | 0.976 | 0.007 | 0.917 | 0.229 |
| Unbabel-Tower70B | 50.687 | 76.317 | 0.920 | 0.999 | **0.999** | **1.000** | 0.991 | 0.009 | 0.986 | 0.208 |
| CycleL | 13.915 | 39.040 | 0.989 | 0.907 | 0.830 | **1.000** | 0.043 | **0.000** | 0.720 | **0.171** |
| CycleL2 | 13.915 | 39.040 | 0.989 | 0.907 | 0.830 | **1.000** | 0.043 | **0.000** | 0.720 | **0.171** |
| Dubformer | 12.618 | 39.766 | 0.272 | 0.483 | 0.515 | 0.857 | 0.196 | 0.748 | 0.484 | 0.563 |
| IOL_Research | 33.076 | 55.070 | 0.812 | 0.607 | 0.531 | 0.999 | 0.918 | 0.081 | 0.804 | 0.389 |
| MSLC | 31.890 | 60.409 | 0.974 | 0.947 | 0.939 | 0.993 | 0.709 | 0.113 | 0.911 | 0.212 |
| ONLINE-A | **66.785** | **83.023** | **0.999** | 0.999 | **0.999** | **1.000** | **0.999** | 0.000 | 0.999 | 0.204 |
| ONLINE-B | 57.270 | 77.814 | 0.245 | **1.000** | 0.998 | **1.000** | 0.996 | 0.004 | 0.891 | 0.300 |
| ONLINE-G | 46.439 | 71.427 | **0.999** | 0.993 | 0.994 | **1.000** | 0.995 | 0.005 | 0.995 | 0.211 |
| ONLINE-W | 62.199 | 79.838 | 0.961 | 0.999 | **0.999** | **1.000** | 0.999 | 0.001 | 0.994 | 0.209 |
| TSU-HITs | 6.294 | 29.317 | 0.144 | 0.652 | 0.853 | 0.946 | 0.353 | 0.168 | 0.526 | 0.294 |
| TranssionMT | 57.217 | 77.757 | 0.242 | **1.000** | 0.998 | **1.000** | 0.996 | 0.004 | 0.891 | 0.300 |

Table 10: English→German, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 59.821 | 79.456 | 0.998 | **1.000** | **1.000** | **1.000** | 0.092 | 0.846 | 0.870 | 0.301 |
| Claude-3 | 45.477 | 65.493 | 0.879 | 0.930 | 0.947 | 0.950 | **0.601** | 0.348 | 0.870 | 0.266 |
| CommandR-plus | 52.798 | 76.068 | 0.906 | 0.965 | 0.958 | **1.000** | 0.108 | 0.856 | 0.841 | 0.327 |
| GPT-4 | 62.776 | 81.285 | **1.000** | **1.000** | **1.000** | **1.000** | 0.326 | 0.627 | 0.904 | 0.280 |
| Llama3-70B | 57.572 | 79.454 | 0.996 | **1.000** | **1.000** | **1.000** | 0.075 | 0.891 | 0.867 | 0.314 |
| NVIDIA-NeMo | 43.543 | 66.512 | 0.999 | 0.995 | 0.994 | **1.000** | 0.291 | 0.683 | 0.895 | 0.291 |
| AIST-AIRC | 50.763 | 73.435 | 0.999 | **1.000** | **1.000** | **1.000** | 0.048 | 0.935 | 0.864 | 0.309 |
| CUNI-NL | 60.950 | 77.784 | 0.892 | **1.000** | **1.000** | **1.000** | 0.069 | 0.776 | 0.849 | 0.277 |
| IKUN | 48.285 | 70.452 | 0.996 | **1.000** | 0.999 | **1.000** | 0.131 | 0.815 | 0.871 | 0.281 |
| IKUN-C | 29.617 | 54.938 | 0.994 | 0.968 | 0.919 | **1.000** | 0.092 | 0.900 | 0.825 | 0.331 |
| Unbabel-Tower70B | 36.617 | 61.602 | 0.998 | 0.984 | 0.938 | **1.000** | 0.179 | 0.814 | 0.857 | 0.328 |
| CycleL | 18.758 | 46.248 | 0.987 | 0.996 | 0.998 | **1.000** | 0.000 | 0.589 | 0.781 | **0.198** |
| CycleL2 | 18.758 | 46.248 | 0.987 | 0.996 | 0.998 | **1.000** | 0.000 | 0.589 | 0.781 | **0.198** |
| Dubformer | 7.240 | 30.085 | 0.922 | 0.406 | 0.359 | 0.144 | 0.006 | **0.179** | 0.398 | 0.520 |
| IOL_Research | **65.014** | **84.971** | 0.998 | **1.000** | **1.000** | **1.000** | 0.097 | 0.827 | 0.871 | 0.301 |
| MSLC | 27.774 | 51.958 | 0.972 | 0.987 | 0.955 | 0.996 | 0.042 | 0.887 | 0.815 | 0.299 |
| ONLINE-A | 53.782 | 80.294 | 0.999 | **1.000** | **1.000** | **1.000** | 0.126 | 0.873 | 0.875 | 0.311 |
| ONLINE-B | 49.961 | 73.532 | 0.998 | 0.999 | 0.998 | **1.000** | 0.430 | 0.540 | **0.918** | 0.255 |
| ONLINE-G | 65.006 | 83.639 | 0.999 | **1.000** | **1.000** | **1.000** | 0.246 | 0.745 | 0.892 | 0.293 |
| ONLINE-W | 55.087 | 82.317 | 0.999 | **1.000** | **1.000** | **1.000** | 0.106 | 0.887 | 0.872 | 0.314 |
| TSU-HITs | 4.685 | 28.741 | 0.083 | 0.728 | 0.898 | 0.918 | 0.034 | 0.387 | 0.473 | 0.281 |
| TranssionMT | 50.021 | 73.607 | 0.998 | 0.999 | 0.998 | **1.000** | 0.428 | 0.541 | 0.917 | 0.254 |

Table 11: English→German, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 50.963 | 76.693 | 0.996 | **1.000** | **1.000** | **1.000** | 0.788 | 0.048 | 0.969 | 0.131 |
| Claude-3 | 54.700 | 71.634 | 0.847 | 0.895 | 0.930 | 0.987 | 0.852 | 0.018 | 0.886 | 0.168 |
| CommandR-plus | 63.094 | 82.606 | 0.939 | **1.000** | **1.000** | **1.000** | 0.734 | 0.171 | 0.953 | 0.158 |
| GPT-4 | 61.142 | 82.555 | **1.000** | **1.000** | **1.000** | **1.000** | 0.865 | 0.037 | 0.981 | 0.140 |
| Llama3-70B | 68.401 | 85.809 | 0.998 | **1.000** | **1.000** | **1.000** | 0.922 | 0.009 | 0.988 | 0.128 |
| NVIDIA-NeMo | 59.526 | 79.044 | 0.991 | **1.000** | **1.000** | **1.000** | 0.901 | 0.050 | 0.985 | 0.146 |
| AIST-AIRC | 54.064 | 77.054 | 0.999 | **1.000** | **1.000** | **1.000** | 0.901 | 0.035 | 0.986 | 0.125 |
| CUNI-NL | 45.673 | 71.102 | 0.984 | **1.000** | 0.999 | **1.000** | 0.471 | 0.175 | 0.922 | 0.138 |
| IKUN | 53.587 | 75.078 | 0.994 | 0.999 | 0.990 | **1.000** | 0.856 | 0.073 | 0.974 | 0.130 |
| IKUN-C | 42.706 | 65.255 | 0.989 | **1.000** | **1.000** | **1.000** | 0.890 | 0.060 | 0.982 | 0.132 |
| Unbabel-Tower70B | 64.058 | 79.666 | 0.995 | **1.000** | **1.000** | **1.000** | **0.985** | **0.004** | **0.997** | 0.141 |
| CycleL | 11.668 | 42.855 | 0.958 | **1.000** | **1.000** | **1.000** | 0.000 | 0.034 | 0.735 | **0.064** |
| CycleL2 | 11.668 | 42.855 | 0.958 | **1.000** | **1.000** | **1.000** | 0.000 | 0.034 | 0.735 | **0.064** |
| Dubformer | 3.704 | 23.383 | 0.939 | 0.376 | 0.382 | 0.018 | 0.005 | 0.346 | 0.271 | 0.502 |
| IOL_Research | 71.042 | 85.830 | 0.999 | **1.000** | **1.000** | **1.000** | 0.820 | 0.055 | 0.974 | 0.132 |
| MSLC | 37.670 | 60.853 | 0.972 | **1.000** | 0.999 | **1.000** | 0.529 | 0.084 | 0.928 | 0.133 |
| ONLINE-A | 66.177 | 86.468 | 0.999 | **1.000** | **1.000** | **1.000** | 0.860 | 0.086 | 0.980 | 0.140 |
| ONLINE-B | 65.085 | 84.832 | 0.998 | **1.000** | **1.000** | **1.000** | 0.823 | 0.037 | 0.974 | 0.124 |
| ONLINE-G | **71.142** | **87.991** | 0.999 | **1.000** | **1.000** | **1.000** | 0.857 | 0.050 | 0.979 | 0.133 |
| ONLINE-W | 55.280 | 81.221 | **1.000** | **1.000** | **1.000** | **1.000** | 0.896 | 0.010 | 0.985 | 0.124 |
| TSU-HITs | 0.239 | 14.339 | 0.024 | 0.499 | 0.579 | 0.955 | 0.004 | 0.201 | 0.306 | 0.314 |
| TranssionMT | 64.962 | 84.750 | 0.998 | **1.000** | **1.000** | **1.000** | 0.825 | 0.037 | 0.975 | 0.124 |

Table 12: English→German, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 68.535 | 70.639 | 0.995 | 0.928 | 0.962 | 0.995 | 0.999 | **0.000** | 0.971 | 0.258 |
| Claude-3 | 68.065 | 62.631 | **0.999** | **0.955** | **0.978** | **0.998** | 0.998 | **0.000** | **0.986** | 0.255 |
| CommandR-plus | 46.057 | 50.994 | 0.897 | 0.559 | 0.586 | 0.670 | 0.617 | 0.080 | 0.653 | 0.436 |
| GPT-4 | 72.389 | 69.642 | **0.999** | 0.896 | 0.928 | 0.942 | 0.945 | **0.000** | 0.940 | 0.280 |
| Llama3-70B | 68.352 | 71.442 | 0.998 | 0.942 | 0.974 | 0.994 | 0.999 | 0.001 | 0.978 | 0.257 |
| NVIDIA-NeMo | 57.014 | 62.936 | 0.989 | 0.912 | 0.936 | 0.976 | 0.969 | **0.000** | 0.950 | 0.265 |
| AIST-AIRC | 70.412 | 70.165 | 0.971 | 0.737 | 0.756 | 0.830 | 0.816 | 0.002 | 0.813 | 0.339 |
| CUNI-NL | 67.845 | 73.794 | 0.895 | 0.825 | 0.852 | 0.993 | 0.901 | 0.001 | 0.878 | 0.277 |
| IKUN | 75.799 | **80.690** | 0.990 | 0.881 | 0.947 | 0.996 | 0.991 | **0.000** | 0.948 | 0.260 |
| IKUN-C | 64.371 | 70.997 | 0.967 | 0.864 | 0.917 | 0.994 | 0.971 | 0.002 | 0.926 | 0.261 |
| Unbabel-Tower70B | 71.215 | 69.715 | 0.989 | 0.633 | 0.651 | 0.651 | 0.654 | 0.001 | 0.703 | 0.400 |
| CycleL | 20.592 | 32.871 | 0.015 | 0.218 | 0.297 | 0.397 | 0.007 | 0.004 | 0.140 | 0.454 |
| CycleL2 | 20.592 | 32.871 | 0.015 | 0.218 | 0.297 | 0.397 | 0.007 | 0.004 | 0.140 | 0.454 |
| Dubformer | 25.567 | 28.961 | 0.294 | 0.047 | 0.064 | 0.180 | 0.004 | 0.316 | 0.106 | 0.691 |
| IOL_Research | 60.629 | 65.159 | 0.996 | 0.925 | 0.965 | 0.985 | 0.993 | 0.001 | 0.968 | 0.260 |
| MSLC | 50.971 | 51.609 | 0.017 | 0.059 | 0.173 | 0.967 | 0.001 | **0.000** | 0.175 | 0.429 |
| ONLINE-A | **79.705** | 80.123 | 0.998 | 0.939 | **0.978** | 0.996 | **1.000** | **0.000** | 0.980 | 0.257 |
| ONLINE-B | 75.136 | 46.306 | 0.998 | 0.934 | 0.962 | 0.996 | 0.995 | 0.004 | 0.971 | **0.255** |
| ONLINE-G | 65.846 | 72.667 | **0.999** | 0.936 | **0.978** | **0.998** | 0.999 | **0.000** | 0.980 | 0.256 |
| ONLINE-W | 71.845 | 77.223 | 0.996 | 0.924 | 0.958 | 0.996 | 0.999 | **0.000** | 0.971 | 0.257 |
| TSU-HITs | 0.090 | 11.361 | 0.000 | 0.034 | 0.042 | 0.264 | 0.000 | 0.028 | 0.049 | 0.540 |
| TranssionMT | 75.074 | 76.169 | 0.998 | 0.931 | 0.962 | **0.998** | 0.996 | 0.004 | 0.971 | 0.256 |

Table 13: English→German, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 64.962 | 70.320 | 0.989 | 0.908 | 0.949 | 0.995 | 0.991 | 0.002 | 0.961 | 0.259 |
| Claude-3 | 52.718 | 52.339 | 0.808 | 0.594 | 0.597 | 0.673 | 0.584 | 0.118 | 0.648 | 0.438 |
| CommandR-plus | 39.666 | 51.642 | 0.968 | 0.498 | 0.498 | 0.534 | 0.529 | 0.058 | 0.585 | 0.465 |
| GPT-4 | 63.042 | 64.243 | **0.999** | 0.381 | 0.364 | 0.335 | 0.340 | 0.001 | 0.459 | 0.530 |
| Llama3-70B | 64.711 | 72.010 | 0.996 | **0.941** | 0.977 | 0.994 | 0.998 | 0.002 | 0.978 | 0.257 |
| NVIDIA-NeMo | 53.905 | 61.422 | 0.681 | 0.678 | 0.710 | 0.967 | 0.665 | 0.001 | 0.714 | 0.313 |
| AIST-AIRC | 57.093 | 63.316 | 0.251 | 0.191 | 0.162 | 0.846 | 0.084 | 0.013 | 0.240 | 0.429 |
| CUNI-NL | 60.424 | 68.018 | 0.905 | 0.800 | 0.854 | 0.994 | 0.901 | 0.007 | 0.873 | 0.279 |
| IKUN | 72.314 | **81.420** | 0.984 | 0.894 | 0.946 | 0.994 | 0.983 | 0.001 | 0.946 | 0.260 |
| IKUN-C | 55.998 | 71.180 | 0.976 | 0.838 | 0.880 | 0.961 | 0.927 | 0.010 | 0.894 | 0.271 |
| Unbabel-Tower70B | 68.188 | 72.826 | 0.993 | 0.770 | 0.797 | 0.802 | 0.805 | **0.000** | 0.824 | 0.336 |
| CycleL | 8.724 | 21.312 | 0.000 | 0.072 | 0.132 | 0.372 | 0.000 | 0.002 | 0.082 | 0.495 |
| CycleL2 | 8.724 | 21.312 | 0.000 | 0.072 | 0.132 | 0.372 | 0.000 | 0.002 | 0.082 | 0.495 |
| Dubformer | 18.630 | 23.978 | 0.360 | 0.039 | 0.023 | 0.010 | 0.009 | 0.621 | 0.078 | 0.824 |
| IOL_Research | 55.451 | 68.917 | 0.991 | 0.917 | 0.962 | 0.995 | 0.998 | **0.000** | 0.966 | 0.257 |
| MSLC | 37.651 | 45.773 | 0.028 | 0.048 | 0.011 | 0.002 | 0.002 | 0.013 | 0.014 | 0.623 |
| ONLINE-A | **74.129** | 78.421 | 0.998 | 0.939 | **0.978** | 0.996 | **1.000** | **0.000** | **0.980** | 0.257 |
| ONLINE-B | 71.704 | 50.631 | 0.998 | 0.913 | 0.960 | 0.996 | 0.998 | 0.002 | 0.965 | 0.257 |
| ONLINE-G | 65.809 | 73.826 | **0.999** | 0.936 | **0.978** | **0.998** | 0.999 | **0.000** | 0.980 | 0.256 |
| ONLINE-W | 66.049 | 72.412 | **0.999** | 0.924 | 0.960 | 0.996 | 0.998 | **0.000** | 0.971 | **0.255** |
| TSU-HITs | 0.000 | 7.087 | 0.001 | 0.031 | 0.028 | 0.301 | 0.002 | 0.015 | 0.052 | 0.532 |
| TranssionMT | 71.683 | 75.342 | 0.998 | 0.913 | 0.958 | **0.998** | 0.996 | 0.004 | 0.965 | 0.258 |

Table 14: English→German, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|--------|------|------|-----|-----|-----|-----|--------|-----|----------|-------|
| Aya23 | 19.085 | 40.614 | 0.993 | 0.058 | 0.885 | 0.971 | 0.933 | 0.062 | 0.670 | **0.257** |
| Claude-3 | 1.919 | 53.543 | 0.989 | **0.148** | **0.909** | 0.977 | 0.936 | 0.058 | **0.721** | 0.262 |
| CommandR-plus | 14.366 | 43.986 | 0.985 | 0.073 | 0.890 | 0.985 | 0.916 | 0.081 | 0.682 | 0.262 |
| GPT-4 | 17.514 | **54.097** | 0.995 | 0.131 | **0.909** | 0.993 | 0.944 | 0.055 | 0.720 | 0.263 |
| Llama3-70B | **27.898** | 43.181 | 0.982 | 0.051 | 0.879 | 0.966 | 0.953 | 0.045 | 0.672 | 0.265 |
| NVIDIA-NeMo | 2.076 | 35.694 | 0.793 | 0.007 | 0.781 | 0.985 | 0.924 | 0.066 | 0.599 | 0.306 |
| AIST-AIRC | 0.719 | 34.974 | 0.933 | 0.005 | 0.796 | **1.000** | 0.956 | 0.039 | 0.638 | 0.287 |
| IKUN | 13.311 | 31.025 | 0.962 | 0.017 | 0.813 | 0.913 | 0.946 | 0.047 | 0.613 | 0.265 |
| IKUN-C | 2.249 | 26.016 | 0.928 | 0.010 | 0.819 | 0.936 | 0.945 | 0.050 | 0.600 | 0.261 |
| Unbabel-Tower70B | 8.143 | 41.692 | 0.944 | 0.053 | 0.891 | 0.980 | 0.930 | 0.069 | 0.672 | 0.272 |
| CycleL | 0.041 | 3.364 | 0.032 | 0.005 | 0.256 | 0.980 | 0.009 | 0.141 | 0.183 | 0.412 |
| DLUT_GTCOM | 0.813 | 42.293 | 0.930 | 0.001 | 0.840 | 0.993 | 0.958 | 0.033 | 0.651 | 0.291 |
| IOL_Research | 19.182 | 51.107 | 0.936 | 0.127 | 0.906 | 0.993 | 0.936 | 0.062 | 0.706 | 0.266 |
| NTTSU | 4.594 | 33.132 | 0.922 | 0.023 | 0.842 | 0.942 | 0.931 | 0.065 | 0.630 | 0.279 |
| ONLINE-A | 1.220 | 44.459 | 0.971 | 0.001 | 0.847 | **1.000** | **0.966** | 0.033 | 0.666 | 0.282 |
| ONLINE-B | 1.015 | 44.589 | 0.995 | 0.062 | 0.890 | 0.996 | 0.952 | 0.045 | 0.692 | 0.266 |
| ONLINE-G | 3.339 | 45.429 | 0.995 | 0.119 | 0.878 | 0.991 | 0.947 | 0.050 | 0.708 | 0.263 |
| ONLINE-W | 4.871 | 34.170 | 0.984 | 0.012 | 0.823 | 0.887 | 0.965 | **0.031** | 0.631 | 0.281 |
| Team-J | 0.416 | 36.323 | **0.999** | 0.001 | 0.827 | **1.000** | 0.941 | 0.055 | 0.653 | 0.275 |
| UvA-MT | 1.159 | 43.238 | 0.942 | 0.001 | 0.852 | 0.999 | 0.965 | 0.032 | 0.661 | 0.292 |

Table 15: English→Japanese, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 2.351 | 33.241 | 0.099 | 0.000 | 0.848 | **1.000** | 0.507 | 0.491 | 0.467 | 0.449 |
| Claude-3 | 0.009 | 0.519 | 0.007 | 0.000 | 0.005 | 0.013 | 0.000 | 1.000 | 0.004 | 0.830 |
| CommandR-plus | 0.087 | 21.326 | 0.610 | 0.001 | 0.404 | 0.805 | 0.367 | 0.591 | 0.379 | 0.515 |
| GPT-4 | 3.434 | 36.947 | 0.776 | 0.004 | 0.531 | 0.990 | 0.671 | 0.328 | 0.541 | 0.403 |
| Llama3-70B | 0.044 | 24.361 | 0.424 | 0.001 | 0.785 | 0.813 | 0.683 | 0.317 | 0.499 | 0.418 |
| NVIDIA-NeMo | 0.108 | 27.875 | 0.459 | **0.005** | 0.487 | 0.742 | 0.741 | 0.204 | 0.435 | 0.504 |
| AIST-AIRC | 0.244 | 41.424 | 0.854 | **0.005** | 0.950 | **1.000** | 0.903 | 0.097 | 0.668 | 0.272 |
| IKUN | 1.464 | 31.780 | 0.154 | 0.001 | 0.950 | **1.000** | 0.662 | 0.334 | 0.522 | 0.384 |
| IKUN-C | 3.047 | 28.905 | 0.513 | 0.000 | 0.881 | **1.000** | 0.792 | 0.207 | 0.555 | 0.319 |
| Unbabel-Tower70B | 0.975 | 38.482 | 0.318 | 0.000 | 0.938 | **1.000** | 0.737 | 0.263 | 0.565 | 0.375 |
| CycleL | 0.036 | 3.754 | 0.009 | **0.005** | 0.301 | 0.976 | 0.000 | 0.126 | 0.184 | 0.407 |
| DLUT_GTCOM | 0.389 | 46.306 | 0.944 | 0.002 | 0.953 | 0.999 | 0.918 | 0.082 | 0.688 | 0.280 |
| IOL_Research | 2.488 | 31.062 | 0.903 | 0.001 | 0.550 | 0.989 | 0.613 | 0.386 | 0.538 | 0.378 |
| NTTSU | 0.533 | 37.444 | 0.865 | 0.001 | 0.953 | 0.998 | 0.789 | 0.207 | 0.646 | 0.281 |
| ONLINE-A | 0.211 | 41.546 | 0.716 | 0.000 | 0.907 | **1.000** | 0.785 | 0.213 | 0.628 | 0.329 |
| ONLINE-B | 0.301 | 41.975 | 0.157 | 0.000 | 0.958 | **1.000** | 0.827 | 0.171 | 0.563 | 0.393 |
| ONLINE-G | 0.675 | 36.629 | 0.346 | 0.000 | 0.911 | **1.000** | 0.736 | 0.263 | 0.564 | 0.382 |
| ONLINE-W | **5.072** | 30.673 | 0.778 | 0.002 | 0.676 | 0.988 | 0.797 | 0.202 | 0.565 | 0.342 |
| Team-J | 0.341 | **48.369** | **0.999** | 0.002 | **0.979** | **1.000** | **0.934** | **0.066** | **0.702** | **0.259** |
| UvA-MT | 0.822 | 41.432 | 0.908 | 0.004 | 0.903 | 0.999 | 0.851 | 0.147 | 0.658 | 0.300 |

Table 16: English→Japanese, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.066 | 35.766 | 0.955 | 0.001 | 0.859 | 0.920 | 0.011 | 0.953 | 0.512 | 0.403 |
| Claude-3 | 0.006 | 15.283 | 0.487 | 0.000 | 0.488 | 0.450 | 0.065 | 0.776 | 0.277 | 0.571 |
| CommandR-plus | 0.057 | 17.950 | 0.404 | 0.002 | 0.350 | 0.936 | 0.026 | 0.953 | 0.287 | 0.548 |
| GPT-4 | 2.559 | **51.818** | **0.999** | 0.001 | **1.000** | **1.000** | 0.015 | 0.955 | **0.574** | 0.377 |
| Llama3-70B | 0.017 | 27.923 | 0.957 | 0.002 | 0.858 | 0.931 | 0.006 | 0.923 | 0.495 | 0.377 |
| NVIDIA-NeMo | 0.048 | 29.803 | 0.793 | 0.005 | 0.876 | **1.000** | 0.002 | 0.983 | 0.500 | 0.398 |
| AIST-AIRC | 0.029 | 33.231 | 0.966 | **0.006** | 0.952 | **1.000** | 0.002 | 0.967 | 0.559 | 0.370 |
| IKUN | 0.068 | 42.846 | 0.967 | 0.004 | 0.996 | **1.000** | 0.016 | 0.703 | 0.569 | **0.317** |
| IKUN-C | 0.250 | 17.596 | 0.854 | 0.001 | 0.416 | 0.923 | 0.005 | 0.994 | 0.348 | 0.473 |
| Unbabel-Tower70B | 0.720 | 36.438 | 0.936 | 0.002 | 0.854 | **1.000** | 0.006 | 0.989 | 0.538 | 0.413 |
| CycleL | 0.006 | 3.367 | 0.013 | 0.004 | 0.326 | 0.982 | 0.000 | 0.431 | 0.189 | 0.444 |
| DLUT_GTCOM | 0.148 | 34.145 | 0.955 | 0.005 | 0.821 | 0.942 | 0.012 | 0.789 | 0.517 | 0.392 |
| IOL_Research | 0.025 | 35.366 | 0.938 | 0.004 | 0.933 | 0.979 | 0.020 | 0.957 | 0.546 | 0.381 |
| NTTSU | 0.061 | 14.918 | 0.780 | 0.004 | 0.343 | 0.184 | **0.184** | 0.267 | 0.235 | 0.595 |
| ONLINE-A | 0.036 | 37.523 | 0.920 | 0.004 | 0.941 | **1.000** | 0.005 | 0.859 | 0.552 | 0.374 |
| ONLINE-B | 0.120 | 40.241 | 0.933 | 0.001 | 0.918 | **1.000** | 0.006 | 0.974 | 0.551 | 0.403 |
| ONLINE-G | 0.087 | 42.494 | 0.995 | 0.004 | 0.993 | 0.974 | 0.031 | 0.909 | 0.571 | 0.368 |
| ONLINE-W | **2.841** | 37.110 | 0.963 | 0.001 | 0.940 | 0.999 | 0.028 | 0.894 | 0.561 | 0.373 |
| Team-J | 0.012 | 32.092 | 0.998 | 0.004 | 0.837 | **1.000** | 0.009 | 0.949 | 0.541 | 0.392 |
| UvA-MT | 0.039 | 10.618 | 0.951 | 0.001 | 0.051 | 0.055 | 0.015 | **0.143** | 0.159 | 0.655 |

Table 17: English→Japanese, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.154 | 40.156 | 0.993 | 0.714 | **1.000** | **1.000** | 0.005 | 0.892 | 0.673 | 0.202 |
| Claude-3 | 0.108 | 44.238 | 0.783 | 0.001 | 0.956 | 0.781 | 0.001 | 0.799 | 0.472 | 0.365 |
| CommandR-plus | 0.190 | 34.767 | 0.856 | 0.118 | 0.968 | **1.000** | 0.002 | 0.928 | 0.559 | 0.335 |
| GPT-4 | 0.223 | 51.838 | **0.999** | 0.002 | **1.000** | **1.000** | 0.006 | 0.983 | 0.573 | 0.335 |
| Llama3-70B | 0.118 | 37.235 | 0.979 | 0.965 | **1.000** | **1.000** | 0.000 | 0.956 | 0.706 | 0.178 |
| NVIDIA-NeMo | **2.027** | 43.511 | 0.829 | 0.991 | **1.000** | **1.000** | 0.002 | 0.974 | 0.693 | 0.205 |
| AIST-AIRC | 0.067 | 46.897 | 0.969 | 0.993 | **1.000** | **1.000** | 0.010 | 0.916 | 0.710 | 0.176 |
| IKUN | 0.056 | **60.658** | 0.980 | **0.998** | **1.000** | **1.000** | 0.002 | 0.897 | 0.711 | **0.148** |
| IKUN-C | 0.055 | 26.628 | 0.814 | 0.371 | 0.985 | 0.988 | 0.007 | 0.829 | 0.573 | 0.260 |
| Unbabel-Tower70B | 0.225 | 41.367 | 0.984 | 0.088 | **1.000** | **1.000** | 0.005 | 0.966 | 0.583 | 0.322 |
| CycleL | 0.011 | 3.560 | 0.006 | 0.949 | 0.463 | 0.998 | 0.000 | **0.061** | 0.345 | 0.235 |
| DLUT_GTCOM | 0.129 | 49.351 | 0.971 | 0.980 | **1.000** | **1.000** | 0.011 | 0.810 | 0.709 | 0.178 |
| IOL_Research | 0.587 | 48.373 | 0.971 | 0.985 | **1.000** | **1.000** | 0.004 | 0.969 | 0.709 | 0.181 |
| NTTSU | 0.112 | 13.293 | 0.564 | 0.388 | 0.421 | 0.346 | **0.037** | 0.279 | 0.283 | 0.421 |
| ONLINE-A | 0.070 | 50.491 | 0.920 | 0.995 | **1.000** | **1.000** | 0.015 | 0.840 | 0.704 | 0.179 |
| ONLINE-B | 0.352 | 51.867 | 0.996 | 0.996 | **1.000** | **1.000** | 0.028 | 0.889 | **0.717** | 0.168 |
| ONLINE-G | 0.182 | 46.613 | 0.995 | 0.989 | **1.000** | **1.000** | 0.011 | 0.968 | 0.714 | 0.185 |
| ONLINE-W | 0.069 | 47.597 | 0.989 | 0.000 | **1.000** | **1.000** | 0.002 | 0.982 | 0.571 | 0.326 |
| Team-J | 0.028 | 49.044 | 0.998 | **0.998** | **1.000** | **1.000** | 0.002 | 0.994 | 0.714 | 0.185 |
| UvA-MT | 0.049 | 13.366 | 0.594 | 0.206 | 0.406 | 0.332 | 0.011 | 0.392 | 0.260 | 0.511 |

Table 18: English→Japanese, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 17.560 | 28.113 | 0.978 | 0.018 | 0.716 | 0.819 | 0.788 | 0.069 | 0.574 | 0.344 |
| Claude-3 | 35.785 | **43.674** | 0.953 | **0.157** | 0.815 | 0.878 | 0.867 | 0.071 | 0.668 | 0.311 |
| CommandR-plus | 13.662 | 22.976 | 0.929 | 0.053 | 0.684 | 0.780 | 0.709 | 0.105 | 0.553 | 0.375 |
| GPT-4 | 30.923 | 40.698 | **0.999** | 0.005 | 0.054 | 0.053 | 0.048 | 0.060 | 0.175 | 0.678 |
| Llama3-70B | 28.549 | 41.916 | 0.803 | 0.043 | 0.772 | 0.968 | 0.772 | 0.131 | 0.573 | 0.296 |
| NVIDIA-NeMo | 24.280 | 27.722 | 0.819 | 0.001 | 0.717 | 0.913 | 0.875 | 0.027 | 0.561 | 0.314 |
| AIST-AIRC | 26.751 | 35.533 | 0.906 | 0.001 | 0.272 | 0.339 | 0.339 | 0.076 | 0.296 | 0.542 |
| IKUN | 18.950 | 34.503 | 0.908 | 0.039 | 0.831 | 0.947 | 0.907 | 0.062 | 0.624 | 0.271 |
| IKUN-C | 25.252 | 36.187 | 0.941 | 0.012 | 0.816 | 0.920 | 0.911 | 0.054 | 0.597 | 0.262 |
| Unbabel-Tower70B | 17.235 | 32.173 | 0.990 | 0.105 | **0.903** | 0.987 | 0.935 | 0.058 | **0.703** | **0.261** |
| CycleL | 0.039 | 4.183 | 0.000 | 0.002 | 0.027 | 0.879 | 0.000 | **0.001** | 0.130 | 0.442 |
| DLUT_GTCOM | 21.587 | 24.692 | 0.043 | 0.000 | 0.062 | 0.159 | 0.023 | 0.103 | 0.042 | 0.586 |
| IOL_Research | 13.093 | 26.752 | 0.827 | 0.069 | 0.808 | 0.856 | 0.804 | 0.061 | 0.599 | 0.300 |
| NTTSU | 11.016 | 26.835 | 0.016 | 0.000 | 0.453 | 0.721 | 0.454 | 0.179 | 0.280 | 0.513 |
| ONLINE-A | 26.998 | 37.914 | 0.372 | 0.000 | 0.379 | 0.372 | 0.367 | 0.004 | 0.259 | 0.453 |
| ONLINE-B | 28.011 | 23.177 | 0.993 | 0.024 | 0.881 | **0.995** | **0.956** | 0.038 | 0.678 | 0.270 |
| ONLINE-G | **38.436** | 28.691 | 0.242 | 0.020 | 0.267 | 0.367 | 0.267 | 0.024 | 0.204 | 0.477 |
| ONLINE-W | 18.163 | 22.537 | 0.104 | 0.000 | 0.192 | 0.098 | 0.100 | 0.005 | 0.082 | 0.525 |
| Team-J | 28.059 | 26.807 | 0.987 | 0.002 | 0.433 | 0.499 | 0.494 | 0.039 | 0.406 | 0.489 |
| UvA-MT | 25.483 | 36.333 | 0.951 | 0.002 | 0.020 | 0.040 | 0.048 | 0.070 | 0.155 | 0.686 |

Table 19: English→Japanese, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 4.104 | 23.321 | 0.985 | 0.035 | 0.793 | 0.905 | 0.865 | 0.048 | 0.630 | 0.306 |
| Claude-3 | 0.001 | 5.146 | 0.002 | 0.000 | 0.100 | 0.699 | 0.018 | 0.980 | 0.136 | 0.707 |
| CommandR-plus | 0.421 | 12.296 | 0.488 | 0.015 | 0.387 | 0.755 | 0.371 | 0.382 | 0.352 | 0.517 |
| GPT-4 | **13.182** | **29.647** | **0.999** | 0.001 | 0.002 | 0.001 | 0.009 | 0.048 | 0.147 | 0.694 |
| Llama3-70B | 4.237 | 25.813 | 0.996 | 0.011 | 0.819 | 0.993 | **0.978** | 0.021 | 0.658 | 0.271 |
| NVIDIA-NeMo | 3.344 | 14.558 | 0.952 | 0.001 | 0.217 | **1.000** | 0.078 | 0.021 | 0.321 | 0.264 |
| AIST-AIRC | 6.832 | 22.356 | 0.013 | 0.000 | 0.012 | 0.016 | 0.002 | 0.015 | 0.007 | 0.568 |
| IKUN | 3.842 | 24.653 | 0.966 | 0.029 | 0.835 | 0.939 | 0.927 | 0.061 | 0.632 | 0.269 |
| IKUN-C | 4.214 | 23.210 | 0.849 | 0.006 | 0.671 | 0.818 | 0.785 | 0.110 | 0.509 | 0.319 |
| Unbabel-Tower70B | 4.522 | 25.322 | 0.989 | 0.095 | 0.884 | 0.991 | 0.940 | 0.054 | 0.698 | 0.263 |
| CycleL | 0.000 | 3.264 | 0.000 | 0.005 | 0.017 | 0.914 | 0.002 | **0.000** | 0.134 | 0.438 |
| DLUT_GTCOM | 2.420 | 10.197 | 0.228 | 0.001 | 0.223 | 0.705 | 0.012 | 0.315 | 0.167 | 0.453 |
| IOL_Research | 4.329 | 24.127 | 0.974 | 0.098 | **0.891** | 0.983 | 0.936 | 0.064 | 0.695 | 0.261 |
| NTTSU | 4.455 | 22.270 | 0.040 | 0.004 | 0.676 | 0.934 | 0.755 | 0.133 | 0.426 | 0.427 |
| ONLINE-A | 6.620 | 25.051 | 0.372 | 0.000 | 0.379 | 0.372 | 0.367 | 0.004 | 0.259 | 0.453 |
| ONLINE-B | 7.834 | 15.040 | 0.983 | 0.006 | 0.882 | 0.996 | 0.925 | 0.055 | 0.668 | 0.277 |
| ONLINE-G | 9.458 | 14.110 | 0.995 | **0.122** | **0.891** | 0.984 | 0.949 | 0.048 | **0.708** | 0.263 |
| ONLINE-W | 4.045 | 12.697 | 0.146 | 0.002 | 0.219 | 0.142 | 0.138 | 0.010 | 0.110 | 0.514 |
| Team-J | 2.152 | 12.340 | 0.979 | 0.002 | 0.307 | **1.000** | 0.076 | 0.093 | 0.339 | **0.260** |
| UvA-MT | 8.453 | 25.104 | 0.998 | 0.002 | 0.000 | 0.000 | 0.000 | 0.002 | 0.143 | 0.431 |

Table 20: English→Japanese, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 44.375 | 63.672 | 0.998 | 0.824 | 0.920 | 0.998 | 0.958 | 0.032 | 0.929 | 0.269 |
| Claude-3 | **60.166** | **76.954** | 0.996 | **0.911** | 0.957 | 0.996 | 0.963 | 0.031 | **0.967** | 0.271 |
| CommandR-plus | 39.996 | 61.592 | 0.988 | 0.819 | 0.917 | 0.996 | 0.928 | 0.061 | 0.916 | 0.272 |
| GPT-4 | 50.565 | 69.608 | 0.998 | 0.908 | 0.942 | **1.000** | 0.963 | 0.029 | 0.956 | **0.265** |
| Llama3-70B | 51.601 | 69.311 | 0.998 | 0.887 | 0.946 | **1.000** | 0.957 | 0.031 | 0.952 | 0.266 |
| NVIDIA-NeMo | 47.354 | 66.582 | 0.984 | 0.827 | 0.928 | 0.999 | 0.953 | 0.027 | 0.931 | 0.280 |
| IKUN | 40.887 | 60.362 | 0.946 | 0.832 | 0.908 | 0.998 | 0.950 | 0.024 | 0.912 | 0.275 |
| IKUN-C | 35.290 | 56.369 | 0.961 | 0.775 | 0.873 | 0.999 | 0.945 | 0.032 | 0.885 | 0.275 |
| Unbabel-Tower70B | 56.242 | 74.129 | 0.998 | 0.908 | **0.963** | **1.000** | 0.953 | 0.038 | 0.965 | 0.272 |
| CycleL | 0.268 | 12.822 | 0.000 | 0.175 | 0.373 | 0.958 | 0.143 | 0.086 | 0.240 | 0.384 |
| IOL_Research | 53.133 | 70.132 | 0.983 | 0.876 | 0.940 | 0.998 | 0.956 | 0.032 | 0.948 | 0.273 |
| ONLINE-A | 59.021 | 74.613 | 0.998 | 0.901 | 0.951 | 0.998 | **0.966** | 0.027 | 0.962 | 0.272 |
| ONLINE-B | 56.473 | 71.907 | 0.998 | 0.892 | 0.957 | **1.000** | 0.956 | 0.037 | 0.956 | 0.268 |
| ONLINE-G | 55.704 | 72.554 | 0.998 | 0.887 | 0.934 | **1.000** | **0.966** | **0.021** | 0.956 | 0.273 |
| ONLINE-W | | | | | | NA | | | | |
| TranssionMT | 56.588 | 73.267 | **0.999** | 0.895 | 0.958 | **1.000** | 0.962 | 0.032 | 0.959 | 0.268 |

Table 21: English→Hindi, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 47.587 | 67.255 | 0.968 | 0.960 | 0.958 | 0.998 | 0.657 | 0.319 | 0.926 | 0.252 |
| Claude-3 | 0.094 | 0.678 | 0.010 | 0.009 | 0.004 | 0.024 | 0.000 | 1.000 | 0.007 | 0.843 |
| CommandR-plus | 11.226 | 20.996 | 0.764 | 0.305 | 0.299 | 0.328 | 0.220 | 0.376 | 0.358 | 0.588 |
| GPT-4 | 33.480 | 53.554 | 0.925 | 0.676 | 0.641 | 0.979 | 0.532 | 0.461 | 0.774 | 0.395 |
| Llama3-70B | 0.450 | 1.517 | 0.082 | 0.022 | 0.020 | 0.026 | 0.018 | 0.979 | 0.030 | 0.829 |
| NVIDIA-NeMo | 41.353 | 66.693 | 0.980 | 0.993 | 0.991 | **1.000** | 0.703 | 0.285 | 0.951 | 0.250 |
| IKUN | 36.678 | 60.934 | 0.924 | 0.987 | 0.980 | **1.000** | 0.673 | 0.304 | 0.926 | 0.250 |
| IKUN-C | 32.860 | 56.557 | 0.956 | 0.978 | 0.963 | **1.000** | 0.681 | 0.304 | 0.922 | **0.241** |
| Unbabel-Tower70B | 44.632 | 69.722 | 0.998 | 0.994 | **0.995** | **1.000** | 0.700 | 0.293 | 0.954 | 0.255 |
| CycleL | 0.218 | 12.777 | 0.000 | 0.284 | 0.370 | **1.000** | 0.001 | **0.119** | 0.237 | 0.360 |
| IOL_Research | 35.627 | 56.917 | 0.979 | 0.750 | 0.727 | 0.998 | 0.627 | 0.362 | 0.837 | 0.334 |
| ONLINE-A | 44.890 | 69.419 | **0.999** | 0.996 | 0.994 | **1.000** | **0.707** | 0.283 | **0.956** | 0.252 |
| ONLINE-B | **57.150** | 74.603 | 0.998 | 0.994 | 0.988 | **1.000** | 0.667 | 0.319 | 0.949 | 0.261 |
| ONLINE-G | 43.515 | 68.688 | **0.999** | **0.998** | 0.991 | **1.000** | 0.705 | 0.285 | 0.955 | 0.250 |
| ONLINE-W | | | | | | NA | | | | |
| TranssionMT | 57.115 | **75.296** | **0.999** | 0.995 | 0.988 | **1.000** | 0.665 | 0.322 | 0.949 | 0.262 |

Table 22: English→Hindi, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 54.461 | 67.057 | 0.891 | 0.996 | 0.979 | 0.999 | 0.006 | 0.800 | 0.838 | 0.312 |
| Claude-3 | 3.221 | 8.801 | 0.060 | 0.453 | 0.557 | 0.157 | **0.083** | 0.610 | 0.205 | 0.532 |
| CommandR-plus | 21.246 | 40.145 | 0.436 | 0.789 | 0.590 | 0.870 | 0.004 | 0.940 | 0.569 | 0.485 |
| GPT-4 | 41.790 | 67.517 | **0.999** | 0.999 | 0.998 | 0.950 | 0.002 | 0.991 | 0.849 | 0.338 |
| Llama3-70B | 38.000 | 53.859 | 0.834 | 0.845 | 0.802 | 0.955 | 0.009 | 0.880 | 0.741 | 0.398 |
| NVIDIA-NeMo | 41.275 | 54.896 | 0.984 | 0.994 | 0.928 | 0.987 | 0.000 | 0.816 | 0.840 | 0.309 |
| IKUN | **56.213** | 68.702 | 0.979 | **1.000** | **1.000** | **1.000** | 0.001 | 0.756 | 0.853 | **0.274** |
| IKUN-C | 17.734 | 36.013 | 0.968 | 0.546 | 0.271 | 0.994 | 0.001 | 0.998 | 0.583 | 0.517 |
| Unbabel-Tower70B | 38.574 | 57.119 | 0.998 | 0.940 | 0.769 | 0.999 | 0.000 | 1.000 | 0.804 | 0.399 |
| CycleL | 0.037 | 8.772 | 0.000 | 0.372 | 0.417 | 0.558 | 0.000 | **0.422** | 0.193 | 0.441 |
| IOL_Research | 39.956 | 65.259 | 0.991 | 0.985 | 0.985 | 0.973 | 0.001 | 0.998 | 0.844 | 0.331 |
| ONLINE-A | 49.365 | 67.457 | **0.999** | **1.000** | **1.000** | **1.000** | 0.001 | 0.958 | **0.857** | 0.331 |
| ONLINE-B | 46.423 | 68.383 | 0.998 | 0.989 | 0.979 | 0.973 | 0.000 | 0.995 | 0.848 | 0.354 |
| ONLINE-G | 33.296 | 57.098 | **0.999** | 0.990 | 0.930 | **1.000** | 0.007 | 0.984 | 0.826 | 0.342 |
| ONLINE-W | | | | | NA | | | | | |
| TranssionMT | 46.395 | **69.002** | **0.999** | 0.990 | 0.980 | 0.971 | 0.000 | 0.995 | 0.848 | 0.355 |

Table 23: English→Hindi, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 61.252 | 73.444 | **0.999** | **1.000** | **1.000** | **1.000** | 0.010 | 0.770 | 0.858 | 0.223 |
| Claude-3 | 0.001 | 3.261 | 0.006 | 0.589 | 0.125 | 0.002 | 0.028 | **0.076** | 0.107 | 0.508 |
| CommandR-plus | 35.854 | 55.859 | 0.963 | **1.000** | 0.996 | **1.000** | **0.033** | 0.760 | 0.853 | 0.216 |
| GPT-4 | 49.390 | 68.971 | 0.998 | **1.000** | **1.000** | **1.000** | 0.023 | 0.962 | **0.860** | 0.254 |
| Llama3-70B | 57.681 | 71.371 | 0.972 | **1.000** | **1.000** | **1.000** | 0.006 | 0.703 | 0.854 | 0.220 |
| NVIDIA-NeMo | 35.768 | 56.995 | 0.732 | **1.000** | **1.000** | **1.000** | 0.002 | 0.944 | 0.819 | 0.277 |
| IKUN | **62.435** | 71.399 | **0.999** | **1.000** | **1.000** | **1.000** | 0.002 | 0.765 | 0.857 | **0.126** |
| IKUN-C | 20.702 | 41.977 | 0.264 | 0.974 | 0.984 | 0.996 | 0.004 | 0.933 | 0.696 | 0.280 |
| Unbabel-Tower70B | 51.432 | 67.662 | 0.996 | **1.000** | **1.000** | **1.000** | 0.006 | 0.973 | 0.857 | 0.272 |
| CycleL | 0.007 | 5.845 | 0.000 | 0.224 | 0.099 | 0.826 | 0.000 | 0.086 | 0.164 | 0.423 |
| IOL_Research | 49.565 | 70.180 | 0.994 | **1.000** | **1.000** | **1.000** | 0.009 | 0.951 | 0.857 | 0.237 |
| ONLINE-A | 54.516 | 71.560 | **0.999** | **1.000** | **1.000** | **1.000** | 0.009 | 0.880 | 0.858 | 0.253 |
| ONLINE-B | 54.462 | 73.655 | 0.998 | **1.000** | **1.000** | **1.000** | 0.006 | 0.983 | 0.858 | 0.260 |
| ONLINE-G | 52.658 | 70.135 | **0.999** | **1.000** | **1.000** | **1.000** | 0.013 | 0.971 | 0.859 | 0.262 |
| ONLINE-W | | | | | NA | | | | | |
| TranssionMT | 54.474 | **73.901** | **0.999** | **1.000** | **1.000** | **1.000** | 0.006 | 0.983 | 0.858 | 0.260 |

Table 24: English→Hindi, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 7.049 | 17.769 | 0.976 | 0.640 | 0.685 | 0.754 | 0.715 | 0.054 | 0.733 | 0.376 |
| Claude-3 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | **0.000** | 0.000 | 0.571 |
| CommandR-plus | 2.197 | 11.202 | 0.485 | 0.192 | 0.230 | 0.588 | 0.200 | 0.140 | 0.297 | 0.523 |
| GPT-4 | 16.415 | **30.954** | **0.999** | 0.361 | 0.361 | 0.388 | 0.367 | 0.049 | 0.462 | 0.530 |
| Llama3-70B | 9.731 | 20.035 | 0.998 | 0.862 | 0.927 | 0.980 | 0.949 | 0.039 | 0.930 | 0.266 |
| NVIDIA-NeMo | 12.329 | 21.184 | 0.875 | 0.759 | 0.810 | 0.931 | 0.201 | 0.741 | 0.703 | 0.360 |
| IKUN | 9.848 | 21.942 | 0.841 | 0.803 | 0.892 | 0.965 | 0.918 | 0.051 | 0.870 | 0.293 |
| IKUN-C | 4.655 | 19.684 | 0.580 | 0.547 | 0.603 | 0.673 | 0.614 | 0.136 | 0.582 | 0.388 |
| Unbabel-Tower70B | 12.473 | 26.454 | 0.982 | 0.880 | 0.939 | 0.982 | 0.946 | 0.027 | 0.942 | 0.271 |
| CycleL | 0.000 | 1.137 | 0.000 | 0.077 | 0.104 | 0.813 | 0.012 | 0.067 | 0.144 | 0.440 |
| IOL_Research | 6.565 | 17.939 | 0.996 | 0.854 | 0.912 | 0.979 | 0.950 | 0.023 | 0.927 | 0.272 |
| ONLINE-A | 11.959 | 23.106 | 0.998 | 0.896 | 0.947 | 0.996 | 0.950 | 0.028 | 0.949 | **0.260** |
| ONLINE-B | **16.480** | 18.335 | 0.998 | 0.898 | 0.951 | 0.999 | 0.950 | 0.032 | 0.953 | 0.261 |
| ONLINE-G | 4.094 | 11.277 | 0.616 | 0.503 | 0.475 | 0.660 | 0.519 | 0.214 | 0.522 | 0.392 |
| ONLINE-W | | | | | | NA | | | | |
| TranssionMT | 16.473 | 30.898 | **0.999** | **0.909** | **0.958** | **1.000** | **0.967** | 0.022 | **0.965** | 0.269 |

Table 25: English→Hindi, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 2.255 | 11.718 | 0.983 | 0.843 | 0.919 | 0.999 | 0.936 | 0.047 | 0.925 | 0.269 |
| Claude-3 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | **0.000** | 0.000 | 0.571 |
| CommandR-plus | 0.348 | 6.697 | 0.624 | 0.228 | 0.255 | 0.561 | 0.271 | 0.220 | 0.357 | 0.570 |
| GPT-4 | 4.312 | **17.611** | 0.996 | 0.035 | 0.016 | 0.015 | 0.022 | 0.028 | 0.160 | 0.682 |
| Llama3-70B | 2.176 | 11.922 | 0.980 | 0.843 | 0.917 | 0.993 | 0.945 | 0.042 | 0.928 | 0.283 |
| NVIDIA-NeMo | 2.735 | 11.847 | 0.022 | 0.076 | 0.130 | 0.928 | 0.001 | 0.905 | 0.168 | 0.552 |
| IKUN | 2.661 | 13.601 | 0.263 | 0.308 | 0.383 | 0.993 | 0.246 | 0.264 | 0.378 | 0.405 |
| IKUN-C | 0.028 | 7.078 | 0.148 | 0.154 | 0.186 | 0.493 | 0.075 | 0.306 | 0.165 | 0.492 |
| Unbabel-Tower70B | 3.154 | 14.423 | 0.989 | 0.885 | 0.933 | 0.993 | 0.952 | 0.033 | 0.949 | 0.273 |
| CycleL | 0.000 | 0.499 | 0.000 | 0.088 | 0.106 | 0.826 | 0.009 | 0.077 | 0.147 | 0.437 |
| IOL_Research | 2.323 | 11.876 | **0.998** | 0.870 | 0.936 | 0.999 | **0.969** | 0.023 | 0.942 | 0.262 |
| ONLINE-A | 3.000 | 11.714 | **0.998** | **0.896** | **0.947** | 0.996 | 0.950 | 0.028 | **0.949** | **0.260** |
| ONLINE-B | **4.807** | 9.601 | 0.864 | 0.829 | 0.819 | 0.931 | 0.836 | 0.043 | 0.826 | 0.280 |
| ONLINE-G | 1.441 | 6.492 | 0.616 | 0.503 | 0.475 | 0.660 | 0.519 | 0.214 | 0.522 | 0.392 |
| ONLINE-W | | | | | | NA | | | | |
| TranssionMT | 4.803 | 17.435 | 0.864 | 0.832 | 0.886 | **1.000** | 0.842 | 0.042 | 0.870 | 0.288 |

Table 26: English→Hindi, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 71.590 | 83.455 | **1.000** | 0.941 | 0.953 | 0.994 | 0.991 | 0.007 | 0.979 | 0.271 |
| Claude-3 | **77.382** | **88.287** | 0.995 | 0.952 | **0.983** | 0.996 | 0.998 | 0.002 | 0.986 | 0.268 |
| CommandR-plus | 69.366 | 82.843 | 0.995 | 0.929 | 0.971 | 0.995 | 0.985 | 0.009 | 0.977 | 0.272 |
| GPT-4 | 76.485 | 86.879 | 0.998 | 0.947 | 0.979 | 0.996 | **1.000** | **0.000** | 0.986 | 0.268 |
| Llama3-70B | 75.659 | 85.899 | 0.994 | 0.936 | 0.972 | 0.996 | **1.000** | **0.000** | 0.983 | 0.270 |
| NVIDIA-NeMo | 71.684 | 83.575 | 0.984 | 0.936 | 0.973 | 0.996 | 0.999 | **0.000** | 0.980 | 0.272 |
| IKUN | 56.366 | 73.524 | 0.987 | 0.869 | 0.922 | 0.998 | 0.989 | 0.004 | 0.953 | 0.276 |
| IKUN-C | 52.543 | 70.275 | 0.999 | 0.849 | 0.923 | 0.991 | 0.991 | 0.004 | 0.945 | 0.274 |
| Occiglot | 49.361 | 68.297 | 0.967 | 0.851 | 0.901 | 0.988 | 0.972 | 0.013 | 0.930 | 0.283 |
| Unbabel-Tower70B | 58.762 | 76.431 | 0.996 | 0.920 | 0.949 | 0.994 | 0.993 | 0.007 | 0.970 | 0.268 |
| CycleL | 32.147 | 51.642 | 0.999 | 0.848 | 0.925 | 0.993 | 0.488 | 0.002 | 0.834 | **0.221** |
| Dubformer | 60.120 | 79.825 | 0.927 | 0.879 | 0.924 | 0.993 | 0.939 | 0.060 | 0.939 | 0.297 |
| IOL_Research | 76.839 | 86.496 | 0.985 | 0.941 | 0.973 | 0.996 | 0.998 | 0.002 | 0.982 | 0.272 |
| MSLC | 56.800 | 74.431 | 0.999 | 0.905 | 0.962 | **0.999** | 0.993 | **0.000** | 0.965 | 0.262 |
| ONLINE-A | 74.616 | 85.820 | 0.998 | 0.952 | 0.976 | 0.996 | 0.999 | 0.001 | 0.986 | 0.266 |
| ONLINE-B | 72.932 | 83.788 | 0.998 | 0.950 | 0.969 | 0.996 | 0.994 | 0.004 | 0.984 | 0.269 |
| ONLINE-G | 76.360 | 86.243 | 0.999 | 0.952 | 0.978 | 0.995 | 0.998 | **0.000** | **0.987** | 0.266 |
| ONLINE-W | 58.478 | 74.701 | 0.999 | 0.896 | 0.945 | 0.996 | 0.999 | 0.001 | 0.964 | 0.271 |
| TSU-HITs | 24.907 | 50.317 | 0.228 | 0.584 | 0.863 | 0.989 | 0.940 | 0.006 | 0.731 | 0.394 |
| TranssionMT | 73.144 | 85.551 | 0.998 | **0.955** | 0.976 | 0.996 | 0.995 | 0.005 | 0.986 | 0.267 |

Table 27: English→Spanish, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 65.702 | 80.565 | 0.985 | 0.999 | 0.998 | **1.000** | 0.956 | 0.044 | 0.991 | 0.243 |
| Claude-3 | 0.227 | 12.311 | 0.009 | 0.024 | 0.061 | 0.005 | 0.000 | 1.000 | 0.016 | 0.823 |
| CommandR-plus | 32.106 | 56.513 | 0.818 | 0.714 | 0.683 | 0.800 | 0.666 | 0.181 | 0.735 | 0.387 |
| GPT-4 | 13.347 | 41.908 | 0.250 | 0.356 | 0.341 | 0.962 | 0.190 | 0.810 | 0.471 | 0.625 |
| Llama3-70B | 0.880 | 13.415 | 0.066 | 0.054 | 0.054 | 0.058 | 0.024 | 0.976 | 0.050 | 0.821 |
| NVIDIA-NeMo | 65.218 | 81.014 | 0.984 | 0.995 | 0.990 | 0.999 | 0.978 | 0.022 | 0.992 | 0.245 |
| IKUN | 40.151 | 62.924 | 0.842 | 0.902 | 0.820 | **1.000** | **0.988** | **0.009** | 0.914 | 0.281 |
| IKUN-C | 39.406 | 63.870 | 0.881 | 0.940 | 0.908 | **1.000** | 0.953 | 0.043 | 0.935 | 0.260 |
| Occiglot | 35.751 | 59.149 | 0.951 | 0.919 | 0.862 | **1.000** | 0.958 | 0.029 | 0.922 | 0.256 |
| Unbabel-Tower70B | 40.903 | 64.886 | 0.974 | 0.935 | 0.843 | 0.998 | 0.984 | 0.015 | 0.941 | 0.263 |
| CycleL | 19.436 | 41.475 | **0.999** | 0.898 | 0.772 | **1.000** | 0.002 | 0.137 | 0.720 | **0.221** |
| Dubformer | 14.020 | 34.187 | 0.277 | 0.295 | 0.284 | 0.693 | 0.113 | 0.837 | 0.357 | 0.675 |
| IOL_Research | 53.335 | 69.931 | 0.933 | 0.887 | 0.862 | **1.000** | 0.917 | 0.082 | 0.935 | 0.294 |
| MSLC | 37.659 | 62.742 | 0.994 | 0.945 | 0.836 | **1.000** | 0.894 | 0.093 | 0.930 | 0.260 |
| ONLINE-A | 65.941 | 82.242 | 0.987 | **1.000** | **1.000** | **1.000** | 0.966 | 0.034 | 0.993 | 0.242 |
| ONLINE-B | **67.157** | 81.476 | 0.994 | 0.999 | **1.000** | **1.000** | 0.979 | 0.021 | **0.996** | 0.239 |
| ONLINE-G | 49.255 | 68.330 | 0.998 | 0.923 | 0.802 | **1.000** | 0.980 | 0.017 | 0.949 | 0.276 |
| ONLINE-W | 61.791 | 77.853 | 0.994 | 0.998 | 0.994 | 0.999 | 0.978 | 0.022 | 0.994 | 0.231 |
| TSU-HITs | 18.159 | 42.517 | 0.029 | 0.800 | 0.922 | 0.947 | 0.388 | 0.299 | 0.614 | 0.340 |
| TranssionMT | 65.473 | **82.434** | 0.990 | **1.000** | **1.000** | **1.000** | 0.966 | 0.034 | 0.994 | 0.241 |

Table 28: English→Spanish, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 56.369 | 72.755 | 0.756 | 0.960 | 0.902 | 0.999 | 0.081 | 0.900 | 0.813 | 0.391 |
| Claude-3 | 42.825 | 62.428 | 0.796 | 0.912 | 0.946 | 0.890 | 0.392 | 0.526 | 0.807 | 0.312 |
| CommandR-plus | 20.217 | 43.259 | 0.269 | 0.703 | 0.552 | 0.737 | 0.060 | 0.925 | 0.529 | 0.566 |
| GPT-4 | 64.419 | 78.732 | **0.999** | **1.000** | **1.000** | **1.000** | 0.386 | 0.606 | **0.912** | 0.292 |
| Llama3-70B | 39.510 | 57.523 | 0.607 | 0.807 | 0.756 | 0.854 | 0.021 | 0.927 | 0.667 | 0.464 |
| NVIDIA-NeMo | 57.475 | 72.389 | 0.996 | **1.000** | 0.996 | **1.000** | 0.300 | 0.695 | 0.899 | 0.316 |
| IKUN | 79.778 | 82.118 | 0.991 | **1.000** | **1.000** | **1.000** | 0.267 | 0.665 | 0.894 | 0.280 |
| IKUN-C | 35.101 | 57.300 | **0.999** | 0.925 | 0.812 | **1.000** | 0.248 | 0.747 | 0.844 | 0.350 |
| Occiglot | 22.173 | 37.527 | 0.679 | 0.469 | 0.299 | 0.966 | 0.004 | 0.983 | 0.525 | 0.569 |
| Unbabel-Tower70B | 41.489 | 62.630 | 0.990 | 0.990 | 0.913 | 0.999 | **0.394** | 0.603 | 0.898 | 0.315 |
| CycleL | 30.751 | 58.344 | **0.999** | **1.000** | 0.999 | **1.000** | 0.000 | 0.636 | 0.852 | **0.223** |
| Dubformer | 31.864 | 45.799 | 0.952 | 0.793 | 0.519 | 0.468 | 0.088 | **0.386** | 0.670 | 0.450 |
| IOL_Research | **87.578** | **92.645** | 0.995 | **1.000** | **1.000** | **1.000** | 0.318 | 0.662 | 0.902 | 0.294 |
| MSLC | 50.773 | 74.280 | 0.996 | **1.000** | **1.000** | **1.000** | 0.022 | 0.971 | 0.860 | 0.317 |
| ONLINE-A | 61.528 | 80.215 | **0.999** | **1.000** | **1.000** | **1.000** | 0.274 | 0.681 | 0.896 | 0.296 |
| ONLINE-B | 83.570 | 91.401 | 0.996 | **1.000** | **1.000** | **1.000** | 0.272 | 0.716 | 0.895 | 0.305 |
| ONLINE-G | 80.674 | 91.066 | **0.999** | **1.000** | **1.000** | **1.000** | 0.275 | 0.721 | 0.896 | 0.303 |
| ONLINE-W | 64.504 | 86.163 | **0.999** | **1.000** | **1.000** | **1.000** | 0.089 | 0.903 | 0.870 | 0.304 |
| TSU-HITs | 22.980 | 48.183 | 0.138 | 0.950 | 0.916 | **1.000** | 0.001 | 0.960 | 0.686 | 0.433 |
| TranssionMT | 61.642 | 80.314 | **0.999** | **1.000** | **1.000** | **1.000** | 0.277 | 0.678 | 0.896 | 0.296 |

Table 29: English→Spanish, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 58.461 | 72.866 | 0.994 | **1.000** | **1.000** | **1.000** | 0.040 | 0.911 | 0.862 | 0.287 |
| Claude-3 | 0.104 | 11.418 | 0.032 | 0.253 | 0.233 | 0.967 | 0.027 | 0.086 | 0.225 | 0.443 |
| CommandR-plus | 36.917 | 59.632 | 0.453 | 0.957 | 0.955 | 0.917 | 0.222 | 0.649 | 0.757 | 0.334 |
| GPT-4 | 74.662 | 87.103 | **0.999** | **1.000** | **1.000** | **1.000** | 0.862 | 0.067 | 0.980 | 0.157 |
| Llama3-70B | 68.395 | 84.329 | 0.852 | **1.000** | **1.000** | **1.000** | 0.594 | 0.225 | 0.921 | 0.196 |
| NVIDIA-NeMo | 76.887 | 86.448 | 0.969 | **1.000** | **1.000** | **1.000** | 0.824 | 0.023 | 0.970 | 0.157 |
| IKUN | 90.086 | 91.617 | 0.996 | **1.000** | 0.999 | **1.000** | 0.638 | 0.013 | 0.948 | 0.121 |
| IKUN-C | 49.512 | 71.442 | 0.989 | **1.000** | **1.000** | **1.000** | 0.865 | 0.094 | 0.979 | 0.149 |
| Occiglot | 40.811 | 58.548 | 0.684 | 0.859 | 0.860 | 0.945 | 0.258 | 0.382 | 0.735 | 0.258 |
| Unbabel-Tower70B | 60.752 | 77.323 | 0.989 | **1.000** | **1.000** | **1.000** | **0.974** | 0.011 | **0.995** | 0.147 |
| CycleL | 27.475 | 58.624 | 0.985 | **1.000** | **1.000** | **1.000** | 0.001 | 0.174 | 0.855 | **0.102** |
| Dubformer | 9.617 | 25.795 | 0.985 | 0.748 | 0.360 | 0.013 | 0.001 | 0.173 | 0.417 | 0.443 |
| IOL_Research | **94.731** | **96.852** | 0.996 | **1.000** | **1.000** | **1.000** | 0.729 | **0.002** | 0.961 | 0.138 |
| MSLC | 55.780 | 78.932 | 0.996 | **1.000** | **1.000** | **1.000** | 0.436 | 0.051 | 0.919 | 0.128 |
| ONLINE-A | 69.092 | 85.663 | **0.999** | **1.000** | **1.000** | **1.000** | 0.644 | 0.070 | 0.949 | 0.151 |
| ONLINE-B | 93.077 | 96.375 | 0.996 | **1.000** | **1.000** | **1.000** | 0.737 | 0.009 | 0.962 | 0.141 |
| ONLINE-G | 88.289 | 95.355 | **0.999** | **1.000** | **1.000** | **1.000** | 0.689 | 0.010 | 0.955 | 0.145 |
| ONLINE-W | 69.406 | 88.122 | 0.998 | **1.000** | **1.000** | **1.000** | 0.470 | 0.017 | 0.924 | 0.123 |
| TSU-HITs | 28.341 | 49.497 | 0.093 | **1.000** | 0.991 | **1.000** | 0.012 | 0.706 | 0.728 | 0.326 |
| TranssionMT | 69.177 | 85.712 | **0.999** | **1.000** | **1.000** | **1.000** | 0.644 | 0.069 | 0.949 | 0.151 |

Table 30: English→Spanish, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 44.372 | 45.246 | 0.919 | 0.335 | 0.341 | 0.346 | 0.326 | 0.061 | 0.429 | 0.544 |
| Claude-3 | 51.127 | 53.605 | 0.994 | 0.947 | **0.977** | 0.990 | 0.993 | 0.005 | 0.982 | 0.273 |
| CommandR-plus | 40.586 | 41.471 | 0.939 | 0.146 | 0.149 | 0.160 | 0.135 | 0.087 | 0.277 | 0.641 |
| GPT-4 | 60.036 | 63.499 | **0.999** | 0.838 | 0.854 | 0.860 | 0.863 | 0.001 | 0.880 | 0.322 |
| Llama3-70B | 53.457 | 67.543 | 0.996 | 0.938 | 0.971 | 0.995 | 0.999 | 0.001 | 0.982 | 0.272 |
| NVIDIA-NeMo | 49.910 | 58.849 | 0.984 | 0.925 | 0.968 | 0.987 | 0.991 | **0.000** | 0.972 | 0.276 |
| IKUN | 63.955 | **74.240** | 0.979 | 0.879 | 0.903 | **0.996** | 0.980 | **0.000** | 0.949 | 0.280 |
| IKUN-C | 51.003 | 64.397 | 0.966 | 0.864 | 0.905 | 0.995 | 0.955 | 0.010 | 0.929 | 0.278 |
| Occiglot | 46.531 | 48.071 | 0.808 | 0.277 | 0.241 | 0.386 | 0.152 | 0.040 | 0.325 | 0.518 |
| Unbabel-Tower70B | 51.763 | 57.905 | 0.989 | 0.798 | 0.829 | 0.854 | 0.847 | 0.001 | 0.859 | 0.327 |
| CycleL | 34.687 | 45.926 | 0.000 | 0.058 | 0.127 | 0.649 | 0.001 | 0.001 | 0.119 | 0.521 |
| Dubformer | 19.544 | 27.118 | 0.574 | 0.056 | 0.064 | 0.196 | 0.001 | 0.152 | 0.147 | 0.656 |
| IOL_Research | 57.457 | 67.980 | 0.994 | 0.929 | 0.960 | 0.994 | 0.993 | 0.005 | 0.976 | 0.272 |
| MSLC | 37.791 | 48.802 | 0.887 | 0.509 | 0.498 | 0.621 | 0.534 | 0.010 | 0.590 | 0.450 |
| ONLINE-A | **64.014** | 72.441 | 0.995 | **0.951** | 0.971 | **0.996** | 1.000 | **0.000** | **0.986** | 0.269 |
| ONLINE-B | 63.810 | 44.291 | 0.995 | 0.938 | 0.971 | 0.995 | 1.000 | **0.000** | 0.983 | 0.271 |
| ONLINE-G | 59.288 | 68.441 | **0.999** | 0.942 | **0.977** | 0.991 | 0.996 | **0.000** | 0.984 | 0.269 |
| ONLINE-W | 61.734 | 69.495 | 0.998 | 0.933 | 0.971 | 0.995 | 0.994 | 0.001 | 0.979 | **0.267** |
| TSU-HITs | 0.001 | 10.778 | 0.009 | 0.091 | 0.102 | 0.344 | 0.005 | 0.020 | 0.094 | 0.511 |
| TranssionMT | 63.942 | 70.704 | 0.994 | 0.941 | 0.972 | **0.996** | 1.000 | **0.000** | 0.984 | 0.271 |

Table 31: English→Spanish, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 55.523 | 64.689 | 0.960 | 0.800 | 0.821 | 0.880 | 0.854 | 0.034 | 0.864 | 0.333 |
| Claude-3 | 49.779 | 53.243 | 0.889 | 0.371 | 0.377 | 0.414 | 0.348 | 0.106 | 0.465 | 0.552 |
| CommandR-plus | 36.233 | 45.131 | 0.900 | 0.135 | 0.126 | 0.170 | 0.078 | 0.087 | 0.250 | 0.641 |
| GPT-4 | 54.713 | 59.924 | **0.999** | 0.187 | 0.173 | 0.127 | 0.130 | 0.002 | 0.287 | 0.620 |
| Llama3-70B | 59.554 | 73.115 | 0.996 | 0.931 | 0.971 | 0.995 | 0.996 | 0.004 | 0.980 | 0.273 |
| NVIDIA-NeMo | 60.284 | 67.454 | 0.903 | 0.862 | 0.905 | 0.994 | 0.906 | 0.007 | 0.907 | 0.289 |
| IKUN | 56.665 | 71.158 | 0.980 | 0.874 | 0.913 | **0.998** | 0.979 | 0.007 | 0.949 | 0.278 |
| IKUN-C | 49.748 | 62.920 | 0.980 | 0.873 | 0.912 | 0.993 | 0.972 | 0.011 | 0.944 | 0.276 |
| Occiglot | 17.068 | 27.726 | 0.458 | 0.181 | 0.204 | 0.169 | 0.099 | 0.021 | 0.194 | 0.566 |
| Unbabel-Tower70B | 57.020 | 65.268 | 0.988 | 0.886 | 0.919 | 0.961 | 0.951 | 0.002 | 0.940 | 0.287 |
| CycleL | 14.341 | 28.581 | 0.000 | 0.054 | 0.116 | 0.712 | 0.000 | **0.000** | 0.126 | 0.512 |
| Dubformer | 11.080 | 21.565 | 0.237 | 0.078 | 0.078 | 0.295 | 0.012 | 0.520 | 0.114 | 0.721 |
| IOL_Research | **67.324** | **78.253** | 0.995 | 0.927 | 0.949 | 0.995 | 0.993 | 0.004 | 0.973 | 0.273 |
| MSLC | 41.545 | 49.482 | 0.912 | 0.731 | 0.736 | 0.782 | 0.739 | 0.009 | 0.767 | 0.354 |
| ONLINE-A | 64.470 | 73.336 | 0.995 | **0.951** | 0.971 | 0.996 | **1.000** | **0.000** | **0.986** | 0.269 |
| ONLINE-B | 66.391 | 51.738 | 0.927 | 0.931 | 0.967 | 0.996 | 0.939 | 0.015 | 0.960 | 0.275 |
| ONLINE-G | 64.605 | 73.968 | **0.999** | 0.944 | 0.974 | 0.995 | 0.999 | **0.000** | 0.986 | 0.269 |
| ONLINE-W | 63.707 | 73.669 | 0.998 | 0.929 | **0.978** | 0.994 | 0.994 | 0.001 | 0.979 | **0.266** |
| TSU-HITs | 0.436 | 22.353 | 0.020 | 0.252 | 0.406 | 0.859 | 0.013 | 0.045 | 0.275 | 0.422 |
| TranssionMT | 65.974 | 73.350 | 0.930 | 0.931 | 0.969 | 0.996 | 0.955 | 0.007 | 0.966 | 0.277 |

Table 32: English→Spanish, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 57.243 | 74.550 | 0.999 | 0.944 | 0.955 | 0.996 | 0.994 | 0.005 | 0.969 | 0.241 |
| Claude-3 | **66.823** | **81.945** | 0.998 | **0.969** | **0.982** | 0.994 | 0.996 | 0.004 | **0.983** | 0.246 |
| CommandR-plus | 54.377 | 73.408 | 0.988 | 0.947 | 0.958 | 0.996 | 0.994 | 0.006 | 0.967 | 0.241 |
| GPT-4 | 64.985 | 79.784 | **1.000** | 0.966 | 0.969 | 0.996 | 0.999 | 0.001 | 0.979 | 0.242 |
| Llama3-70B | 61.753 | 77.069 | 0.999 | 0.961 | 0.967 | 0.998 | 0.994 | 0.004 | 0.979 | 0.243 |
| NVIDIA-NeMo | 55.940 | 72.507 | 0.979 | 0.914 | 0.955 | **0.999** | 0.994 | **0.000** | 0.966 | 0.251 |
| CUNI-MH | 57.511 | 75.301 | 0.998 | 0.966 | 0.971 | 0.995 | 0.988 | 0.010 | 0.980 | 0.237 |
| IKUN | 45.469 | 65.478 | **1.000** | 0.898 | 0.914 | 0.995 | 0.985 | 0.005 | 0.939 | 0.241 |
| IKUN-C | 37.968 | 58.621 | 0.996 | 0.848 | 0.901 | 0.995 | 0.971 | 0.009 | 0.908 | 0.237 |
| SCIR-MT | 63.339 | 78.457 | 0.987 | 0.942 | 0.966 | 0.995 | 0.995 | 0.002 | 0.976 | 0.253 |
| Unbabel-Tower70B | 51.206 | 71.180 | 0.990 | 0.936 | 0.957 | 0.996 | 0.990 | 0.010 | 0.961 | 0.238 |
| CUNI-DocTransformer | 58.378 | 75.431 | 0.998 | 0.935 | 0.972 | 0.995 | 0.996 | 0.001 | 0.973 | 0.244 |
| CUNI-GA | 56.400 | 74.149 | 0.998 | 0.931 | 0.966 | 0.994 | 0.999 | **0.000** | 0.972 | 0.243 |
| CUNI-Transformer | 56.400 | 74.149 | 0.998 | 0.931 | 0.966 | 0.994 | 0.999 | **0.000** | 0.972 | 0.243 |
| CycleL | 1.469 | 17.798 | 0.987 | 0.800 | 0.805 | 0.993 | 0.015 | 0.002 | 0.537 | **0.082** |
| CycleL2 | 5.734 | 24.422 | 0.988 | 0.785 | 0.826 | 0.994 | 0.122 | 0.006 | 0.602 | 0.120 |
| IOL_Research | 64.617 | 78.908 | 0.988 | 0.950 | 0.965 | 0.995 | 0.990 | 0.007 | 0.976 | 0.250 |
| ONLINE-A | 63.853 | 79.054 | 0.999 | 0.946 | 0.968 | 0.995 | **1.000** | **0.000** | 0.980 | 0.248 |
| ONLINE-B | 59.851 | 76.425 | 0.998 | 0.936 | 0.963 | 0.995 | 0.998 | 0.002 | 0.974 | 0.247 |
| ONLINE-G | 63.404 | 78.063 | 0.999 | 0.950 | 0.967 | 0.995 | **1.000** | **0.000** | 0.981 | 0.248 |
| ONLINE-W | 55.114 | 73.094 | 0.999 | 0.941 | 0.963 | 0.996 | 0.995 | 0.001 | 0.970 | 0.242 |
| TSU-HITs | 16.169 | 34.946 | 0.081 | 0.545 | 0.725 | 0.977 | 0.596 | 0.047 | 0.565 | 0.385 |
| TranssionMT | 62.123 | 78.598 | 0.999 | 0.949 | 0.971 | 0.995 | 0.999 | 0.001 | 0.979 | 0.246 |

Table 33: English→Czech, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 43.235 | 64.720 | 0.988 | 0.931 | 0.891 | 0.999 | 0.889 | 0.102 | 0.941 | 0.227 |
| Claude-3 | 0.221 | 9.467 | 0.006 | 0.032 | 0.040 | 0.006 | 0.000 | 1.000 | 0.013 | 0.834 |
| CommandR-plus | 13.247 | 31.471 | 0.729 | 0.296 | 0.267 | 0.458 | 0.286 | 0.425 | 0.382 | 0.594 |
| GPT-4 | 19.672 | 40.563 | 0.480 | 0.428 | 0.348 | 0.971 | 0.345 | 0.654 | 0.547 | 0.555 |
| Llama3-70B | 17.102 | 48.921 | 0.778 | 0.777 | 0.765 | 0.800 | 0.783 | 0.217 | 0.777 | 0.335 |
| NVIDIA-NeMo | 48.364 | 70.655 | 0.987 | 0.994 | 0.979 | **1.000** | 0.993 | 0.005 | 0.989 | 0.200 |
| CUNI-MH | 56.704 | 77.481 | 0.998 | **1.000** | **1.000** | **1.000** | 0.988 | 0.012 | 0.998 | 0.190 |
| IKUN | 36.215 | 60.755 | 0.864 | 0.945 | 0.897 | **1.000** | 0.884 | 0.108 | 0.921 | 0.233 |
| IKUN-C | 28.458 | 53.178 | 0.974 | 0.965 | 0.906 | 0.999 | 0.919 | 0.061 | 0.931 | 0.193 |
| SCIR-MT | **76.711** | **86.823** | 0.987 | **1.000** | 0.999 | **1.000** | 0.989 | 0.009 | 0.996 | 0.198 |
| Unbabel-Tower70B | 47.569 | 72.172 | 0.969 | 0.993 | 0.988 | **1.000** | 0.979 | 0.017 | 0.982 | 0.191 |
| CUNI-DocTransformer | 60.086 | 79.628 | 0.998 | 0.996 | 0.996 | **1.000** | 0.991 | 0.007 | 0.997 | 0.191 |
| CUNI-GA | 36.980 | 62.546 | 0.987 | 0.968 | 0.936 | **1.000** | 0.968 | 0.031 | 0.952 | 0.201 |
| CUNI-Transformer | 36.980 | 62.546 | 0.987 | 0.968 | 0.936 | **1.000** | 0.968 | 0.031 | 0.952 | 0.201 |
| CycleL | 1.019 | 17.922 | 0.982 | 0.765 | 0.756 | 0.999 | 0.001 | **0.000** | 0.507 | **0.089** |
| CycleL2 | 3.482 | 21.004 | 0.995 | 0.878 | 0.692 | **1.000** | 0.004 | 0.037 | 0.542 | 0.112 |
| IOL_Research | 34.952 | 55.208 | 0.879 | 0.764 | 0.643 | **1.000** | 0.816 | 0.176 | 0.838 | 0.341 |
| ONLINE-A | 59.187 | 79.085 | **0.999** | **1.000** | 0.998 | **1.000** | 0.971 | 0.029 | 0.994 | 0.199 |
| ONLINE-B | 58.821 | 78.935 | 0.998 | **1.000** | 0.999 | **1.000** | **0.998** | 0.002 | 0.999 | 0.193 |
| ONLINE-G | 58.898 | 79.300 | **0.999** | 1.000 | 1.000 | **1.000** | 0.995 | 0.002 | **0.999** | 0.197 |
| ONLINE-W | 57.748 | 77.987 | 0.998 | **1.000** | **1.000** | **1.000** | **0.998** | 0.002 | 0.999 | 0.195 |
| TSU-HITs | 16.823 | 37.143 | 0.029 | 0.749 | 0.843 | 0.978 | 0.449 | 0.177 | 0.600 | 0.267 |
| TranssionMT | 58.756 | 79.292 | **0.999** | **1.000** | 0.998 | **1.000** | 0.972 | 0.028 | 0.995 | 0.199 |

Table 34: English→Czech, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 43.861 | 66.893 | 0.999 | 0.999 | 0.996 | **1.000** | 0.295 | 0.608 | 0.894 | 0.248 |
| Claude-3 | 3.979 | 18.708 | 0.127 | 0.481 | 0.518 | 0.444 | 0.229 | 0.688 | 0.307 | 0.512 |
| CommandR-plus | 32.040 | 48.780 | 0.627 | 0.913 | 0.831 | **1.000** | 0.081 | 0.858 | 0.700 | 0.365 |
| GPT-4 | 50.035 | 68.241 | **1.000** | 0.998 | 0.995 | **1.000** | **0.775** | 0.209 | **0.966** | 0.201 |
| Llama3-70B | 43.095 | 59.015 | 0.780 | 0.783 | 0.732 | 0.944 | 0.192 | 0.782 | 0.726 | 0.403 |
| NVIDIA-NeMo | **66.822** | 75.097 | 0.989 | 0.998 | 0.996 | 0.996 | 0.038 | 0.956 | 0.859 | 0.311 |
| CUNI-MH | 39.866 | 58.411 | 0.996 | 0.999 | 0.993 | **1.000** | 0.273 | 0.714 | 0.893 | 0.259 |
| IKUN | 51.933 | 71.198 | 0.998 | **1.000** | **1.000** | **1.000** | 0.600 | 0.372 | 0.941 | 0.204 |
| IKUN-C | 26.890 | 46.890 | 0.995 | 0.919 | 0.733 | **1.000** | 0.062 | 0.936 | 0.777 | 0.345 |
| SCIR-MT | 61.409 | 69.918 | 0.987 | **1.000** | 0.999 | **1.000** | 0.073 | 0.907 | 0.866 | 0.302 |
| Unbabel-Tower70B | 39.160 | 61.603 | **1.000** | 0.996 | 0.990 | **1.000** | 0.449 | 0.481 | 0.910 | 0.238 |
| CUNI-DocTransformer | 53.662 | 74.090 | 0.996 | **1.000** | **1.000** | **1.000** | 0.187 | 0.802 | 0.883 | 0.279 |
| CUNI-GA | 58.498 | **79.747** | 0.998 | **1.000** | **1.000** | **1.000** | 0.624 | 0.360 | 0.946 | 0.217 |
| CUNI-Transformer | 58.498 | **79.747** | 0.998 | **1.000** | **1.000** | **1.000** | 0.624 | 0.360 | 0.946 | 0.217 |
| CycleL | 1.240 | 17.016 | 0.984 | 0.978 | 0.824 | 0.995 | 0.000 | **0.098** | 0.541 | **0.058** |
| CycleL2 | 9.402 | 29.260 | 0.994 | 0.998 | 0.977 | **1.000** | 0.000 | 0.317 | 0.703 | 0.087 |
| IOL_Research | 46.391 | 65.294 | 0.999 | 0.999 | **1.000** | **1.000** | 0.588 | 0.375 | 0.940 | 0.216 |
| ONLINE-A | 56.372 | 77.344 | 0.999 | **1.000** | **1.000** | **1.000** | 0.359 | 0.592 | 0.908 | 0.259 |
| ONLINE-B | 47.934 | 64.659 | 0.998 | **1.000** | 0.991 | **1.000** | 0.301 | 0.667 | 0.898 | 0.262 |
| ONLINE-G | 37.706 | 66.368 | 0.999 | 0.991 | 0.999 | **1.000** | 0.760 | 0.131 | 0.948 | 0.188 |
| ONLINE-W | 56.533 | 78.862 | 0.999 | **1.000** | **1.000** | **1.000** | 0.078 | 0.920 | 0.868 | 0.293 |
| TSU-HITs | 13.616 | 36.001 | 0.061 | 0.873 | 0.940 | 0.998 | 0.002 | 0.771 | 0.564 | 0.341 |
| TranssionMT | 49.884 | 67.486 | 0.999 | **1.000** | 0.990 | 0.993 | 0.307 | 0.659 | 0.897 | 0.265 |

Table 35: English→Czech, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 55.072 | 75.740 | 0.996 | **1.000** | **1.000** | **1.000** | 0.125 | 0.542 | 0.874 | 0.180 |
| Claude-3 | 0.767 | 9.381 | 0.013 | 0.834 | 0.682 | 0.881 | 0.018 | 0.346 | 0.349 | 0.286 |
| CommandR-plus | 49.165 | 68.286 | 0.868 | **1.000** | 0.996 | **1.000** | 0.419 | 0.428 | 0.895 | 0.201 |
| GPT-4 | 54.286 | 75.679 | 0.999 | **1.000** | **1.000** | **1.000** | 0.460 | 0.229 | 0.923 | 0.143 |
| Llama3-70B | 42.181 | 65.074 | 0.971 | **1.000** | **1.000** | **1.000** | **0.661** | 0.093 | 0.947 | 0.117 |
| NVIDIA-NeMo | 63.608 | 74.237 | 0.907 | 0.998 | 0.995 | 0.950 | 0.132 | 0.696 | 0.840 | 0.236 |
| CUNI-MH | 50.682 | 71.752 | 0.998 | **1.000** | **1.000** | **1.000** | 0.652 | 0.182 | **0.950** | 0.126 |
| IKUN | 73.046 | 82.347 | 0.999 | **1.000** | 0.985 | **1.000** | 0.067 | 0.447 | 0.864 | 0.149 |
| IKUN-C | 41.076 | 65.308 | 0.989 | **1.000** | **1.000** | **1.000** | 0.542 | 0.259 | 0.933 | 0.137 |
| SCIR-MT | **82.937** | **86.289** | 0.987 | **1.000** | **1.000** | **1.000** | 0.162 | 0.364 | 0.878 | 0.172 |
| Unbabel-Tower70B | 56.279 | 74.296 | 0.998 | **1.000** | **1.000** | **1.000** | 0.506 | 0.248 | 0.929 | 0.144 |
| CUNI-DocTransformer | 61.577 | 81.949 | 0.998 | **1.000** | **1.000** | **1.000** | 0.209 | 0.524 | 0.887 | 0.181 |
| CUNI-GA | 70.410 | 85.330 | 0.998 | **1.000** | **1.000** | **1.000** | 0.059 | 0.443 | 0.865 | 0.170 |
| CUNI-Transformer | 70.410 | 85.330 | 0.998 | **1.000** | **1.000** | **1.000** | 0.059 | 0.443 | 0.865 | 0.170 |
| CycleL | 0.446 | 18.362 | 0.859 | 0.887 | 0.961 | 0.998 | 0.000 | **0.021** | 0.529 | 0.047 |
| CycleL2 | 6.341 | 28.958 | 0.974 | 0.996 | 0.996 | **1.000** | 0.000 | 0.181 | 0.704 | **0.042** |
| IOL_Research | 50.768 | 73.431 | **1.000** | **1.000** | **1.000** | **1.000** | 0.195 | 0.321 | 0.885 | 0.144 |
| ONLINE-A | 60.436 | 80.049 | 0.999 | **1.000** | **1.000** | **1.000** | 0.192 | 0.397 | 0.884 | 0.171 |
| ONLINE-B | 77.459 | 83.866 | 0.998 | **1.000** | **1.000** | **1.000** | 0.218 | 0.446 | 0.888 | 0.176 |
| ONLINE-G | 57.235 | 77.528 | 0.999 | **1.000** | **1.000** | **1.000** | 0.186 | 0.460 | 0.884 | 0.179 |
| ONLINE-W | 59.116 | 81.841 | 0.999 | **1.000** | **1.000** | **1.000** | 0.031 | 0.670 | 0.861 | 0.196 |
| TSU-HITs | 6.089 | 29.588 | 0.028 | 0.849 | 0.951 | 0.999 | 0.001 | 0.618 | 0.531 | 0.291 |
| TranssionMT | 60.824 | 80.180 | 0.999 | **1.000** | **1.000** | **1.000** | 0.195 | 0.397 | 0.885 | 0.171 |

Table 36: English→Czech, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 41.106 | 43.140 | 0.928 | 0.879 | 0.896 | 0.985 | 0.927 | 0.060 | 0.915 | 0.277 |
| Claude-3 | 32.438 | 37.308 | 0.973 | 0.944 | 0.956 | 0.976 | 0.968 | 0.026 | 0.962 | 0.265 |
| CommandR-plus | 29.014 | 31.981 | 0.958 | 0.623 | 0.636 | 0.698 | 0.689 | 0.062 | 0.697 | 0.400 |
| GPT-4 | 40.384 | 46.163 | **1.000** | 0.397 | 0.393 | 0.394 | 0.410 | 0.058 | 0.490 | 0.514 |
| Llama3-70B | 39.552 | 51.856 | 0.999 | **0.968** | **0.974** | 0.995 | 0.998 | 0.001 | **0.981** | 0.241 |
| NVIDIA-NeMo | 41.978 | 44.362 | 0.974 | 0.826 | 0.832 | 0.869 | 0.869 | 0.012 | 0.863 | 0.303 |
| CUNI-MH | 26.635 | 39.740 | **1.000** | 0.690 | 0.627 | **1.000** | 0.000 | **0.000** | 0.474 | **0.098** |
| IKUN | **65.224** | **70.139** | 0.905 | 0.808 | 0.864 | 0.987 | 0.909 | 0.005 | 0.840 | 0.227 |
| IKUN-C | 45.966 | 50.534 | 0.936 | 0.859 | 0.895 | 0.979 | 0.908 | 0.017 | 0.877 | 0.235 |
| SCIR-MT | 31.395 | 47.915 | 0.989 | 0.940 | 0.967 | 0.991 | 0.995 | 0.002 | 0.974 | 0.251 |
| Unbabel-Tower70B | 38.998 | 46.846 | 0.955 | 0.892 | 0.920 | 0.968 | 0.944 | 0.017 | 0.924 | 0.254 |
| CUNI-DocTransformer | 13.558 | 35.236 | 0.998 | 0.444 | 0.449 | 0.441 | 0.492 | 0.062 | 0.540 | 0.496 |
| CUNI-GA | 12.724 | 32.388 | 0.942 | 0.174 | 0.149 | 0.116 | 0.198 | 0.087 | 0.275 | 0.639 |
| CUNI-Transformer | 12.724 | 32.388 | 0.942 | 0.174 | 0.149 | 0.116 | 0.198 | 0.087 | 0.275 | 0.639 |
| CycleL | 0.608 | 11.178 | 0.000 | 0.089 | 0.073 | 0.647 | 0.000 | 0.001 | 0.116 | 0.460 |
| CycleL2 | 7.851 | 18.916 | 0.000 | 0.099 | 0.126 | 0.787 | 0.000 | 0.015 | 0.145 | 0.429 |
| IOL_Research | 29.386 | 38.865 | 0.993 | 0.939 | 0.956 | 0.983 | 0.979 | 0.007 | 0.963 | 0.240 |
| ONLINE-A | 49.782 | 59.772 | 0.996 | 0.936 | 0.967 | 0.995 | 0.998 | **0.000** | 0.975 | 0.246 |
| ONLINE-B | 51.148 | 35.081 | 0.994 | 0.935 | 0.965 | 0.988 | 0.994 | 0.001 | 0.972 | 0.247 |
| ONLINE-G | 42.713 | 48.902 | 0.998 | 0.903 | 0.961 | 0.995 | 0.993 | 0.005 | 0.967 | 0.255 |
| ONLINE-W | 52.745 | 58.313 | 0.965 | 0.923 | 0.939 | 0.994 | 0.925 | 0.037 | 0.936 | 0.238 |
| TSU-HITs | 0.000 | 6.347 | 0.007 | 0.100 | 0.177 | 0.902 | 0.006 | 0.007 | 0.177 | 0.409 |
| TranssionMT | 50.961 | 58.322 | 0.993 | 0.949 | 0.968 | 0.993 | **0.999** | **0.000** | 0.976 | 0.245 |

Table 37: English→Czech, 0-shot JSON format

412

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 29.727 | 38.018 | 0.994 | 0.927 | 0.953 | 0.995 | 0.991 | 0.007 | 0.963 | 0.244 |
| Claude-3 | 5.118 | 16.803 | 0.200 | 0.106 | 0.125 | 0.573 | 0.050 | 0.683 | 0.212 | 0.698 |
| CommandR-plus | 21.301 | 28.347 | 0.974 | 0.465 | 0.476 | 0.513 | 0.518 | 0.088 | 0.561 | 0.478 |
| GPT-4 | 28.891 | 37.329 | **0.999** | 0.073 | 0.043 | 0.004 | 0.054 | 0.087 | 0.186 | 0.684 |
| Llama3-70B | 28.057 | 41.364 | **0.999** | **0.960** | 0.966 | **0.996** | **0.999** | 0.001 | **0.978** | 0.244 |
| NVIDIA-NeMo | 27.813 | 34.402 | 0.890 | 0.448 | 0.453 | 0.494 | 0.457 | 0.049 | 0.518 | 0.486 |
| CUNI-MH | 31.838 | 35.068 | 0.994 | 0.936 | 0.966 | 0.989 | 0.987 | 0.010 | 0.962 | **0.237** |
| IKUN | **45.464** | **55.355** | 0.310 | 0.357 | 0.382 | **0.996** | 0.269 | 0.148 | 0.402 | 0.359 |
| IKUN-C | 31.924 | 42.558 | 0.563 | 0.552 | 0.573 | 0.783 | 0.333 | 0.038 | 0.478 | 0.290 |
| SCIR-MT | 24.259 | 42.627 | 0.989 | 0.940 | 0.967 | 0.991 | 0.995 | 0.002 | 0.974 | 0.251 |
| Unbabel-Tower70B | 28.042 | 38.886 | 0.967 | 0.854 | 0.869 | 0.901 | 0.887 | 0.017 | 0.885 | 0.282 |
| CUNI-DocTransformer | 6.961 | 28.507 | 0.998 | 0.859 | 0.890 | 0.898 | 0.909 | 0.012 | 0.904 | 0.290 |
| CUNI-GA | 3.680 | 20.055 | 0.976 | 0.436 | 0.110 | 0.073 | 0.001 | 0.024 | 0.231 | 0.353 |
| CUNI-Transformer | 3.680 | 20.055 | 0.976 | 0.436 | 0.110 | 0.073 | 0.001 | 0.024 | 0.231 | 0.353 |
| CycleL | 0.122 | 7.126 | 0.000 | 0.078 | 0.078 | 0.621 | 0.000 | 0.002 | 0.111 | 0.464 |
| CycleL2 | 1.307 | 11.054 | 0.000 | 0.097 | 0.113 | 0.815 | 0.000 | 0.017 | 0.146 | 0.428 |
| IOL_Research | 18.249 | 31.706 | 0.984 | 0.922 | 0.949 | 0.974 | 0.967 | 0.012 | 0.951 | 0.245 |
| ONLINE-A | 34.158 | 45.736 | 0.996 | 0.936 | **0.969** | 0.995 | 0.998 | **0.000** | 0.976 | 0.246 |
| ONLINE-B | 34.468 | 28.263 | 0.994 | 0.876 | 0.936 | 0.993 | 0.984 | 0.006 | 0.943 | 0.248 |
| ONLINE-G | 29.851 | 39.063 | 0.998 | 0.903 | 0.961 | 0.995 | 0.993 | 0.005 | 0.967 | 0.255 |
| ONLINE-W | 36.605 | 44.861 | 0.890 | 0.820 | 0.834 | 0.737 | 0.692 | 0.011 | 0.767 | 0.277 |
| TSU-HITs | 0.001 | 7.399 | 0.009 | 0.154 | 0.197 | 0.984 | 0.005 | 0.006 | 0.208 | 0.395 |
| TranssionMT | 34.162 | 45.736 | 0.998 | 0.936 | **0.969** | 0.995 | **0.999** | **0.000** | 0.976 | 0.245 |

Table 38: English→Czech, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | **31.789** | 44.156 | **1.000** | 0.119 | 0.887 | 0.995 | 0.974 | 0.013 | 0.700 | 0.234 |
| Claude-3 | 6.160 | 52.796 | 0.994 | **0.246** | **0.928** | 0.989 | 0.969 | 0.021 | **0.757** | 0.232 |
| CommandR-plus | 12.165 | 44.248 | 0.990 | 0.146 | 0.896 | 0.993 | 0.966 | 0.027 | 0.712 | 0.239 |
| GPT-4 | 11.698 | 49.684 | 0.999 | 0.201 | 0.918 | 0.993 | 0.974 | 0.012 | 0.743 | 0.233 |
| Llama3-70B | 14.886 | 45.139 | 0.998 | 0.127 | 0.901 | 0.991 | 0.972 | 0.016 | 0.709 | 0.235 |
| NVIDIA-NeMo | 1.315 | 35.577 | 0.971 | 0.062 | 0.847 | 0.984 | 0.942 | 0.023 | 0.655 | 0.252 |
| IKUN | 2.722 | 34.863 | 0.995 | 0.069 | 0.827 | 0.991 | 0.947 | 0.038 | 0.652 | 0.243 |
| IKUN-C | 4.261 | 29.898 | 0.989 | 0.039 | 0.832 | 0.989 | 0.931 | 0.045 | 0.627 | 0.236 |
| Unbabel-Tower70B | 2.125 | 40.668 | 0.994 | 0.105 | 0.887 | 0.988 | 0.974 | 0.016 | 0.696 | 0.243 |
| CycleL | 0.057 | 3.676 | 0.837 | 0.007 | 0.319 | 0.930 | 0.075 | 0.040 | 0.311 | 0.289 |
| CycleL2 | 0.000 | 0.779 | 0.032 | 0.000 | 0.073 | 0.635 | 0.004 | **0.009** | 0.106 | 0.475 |
| HW-TSC | 18.593 | 47.754 | 0.999 | 0.195 | 0.916 | 0.993 | 0.965 | 0.024 | 0.736 | **0.230** |
| IOL_Research | 28.529 | **54.058** | 0.999 | 0.230 | 0.919 | 0.991 | 0.965 | 0.027 | 0.752 | 0.232 |
| ONLINE-A | 11.048 | 49.271 | 0.999 | 0.190 | 0.876 | 0.996 | 0.961 | 0.026 | 0.727 | 0.235 |
| ONLINE-B | 2.844 | 45.939 | 0.999 | 0.143 | 0.891 | 0.991 | 0.963 | 0.027 | 0.711 | 0.241 |
| ONLINE-G | 2.939 | 42.534 | 0.998 | 0.129 | 0.907 | 0.996 | 0.961 | 0.028 | 0.706 | 0.238 |
| ONLINE-W | 3.376 | 44.271 | 0.999 | 0.144 | 0.887 | **0.998** | 0.962 | 0.027 | 0.707 | 0.240 |
| UvA-MT | 0.668 | 34.492 | 0.978 | 0.011 | 0.807 | 0.985 | **0.977** | 0.015 | 0.641 | 0.276 |

Table 39: English→Chinese, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 5.186 | 40.417 | 0.996 | 0.004 | 0.976 | **1.000** | 0.968 | 0.027 | 0.698 | 0.226 |
| Claude-3 | 0.019 | 1.919 | 0.006 | 0.001 | 0.077 | 0.262 | 0.000 | 0.998 | 0.051 | 0.772 |
| CommandR-plus | 1.494 | 34.071 | 0.918 | 0.007 | 0.834 | 0.969 | 0.869 | 0.109 | 0.629 | 0.282 |
| GPT-4 | **6.933** | 30.820 | 0.919 | 0.006 | 0.518 | 0.993 | 0.611 | 0.386 | 0.532 | 0.383 |
| Llama3-70B | 0.126 | 27.102 | 0.891 | 0.005 | 0.870 | 0.895 | 0.882 | 0.115 | 0.629 | 0.294 |
| NVIDIA-NeMo | 0.581 | 35.996 | 0.983 | **0.010** | 0.957 | 0.993 | 0.953 | 0.040 | 0.682 | 0.240 |
| IKUN | 2.194 | 34.039 | 0.973 | 0.004 | 0.922 | **1.000** | 0.968 | 0.028 | 0.665 | 0.236 |
| IKUN-C | 1.653 | 32.757 | 0.965 | 0.006 | 0.951 | **1.000** | 0.961 | 0.033 | 0.671 | 0.215 |
| Unbabel-Tower70B | 2.875 | 37.019 | 0.976 | 0.001 | 0.958 | **1.000** | 0.957 | 0.040 | 0.686 | 0.239 |
| CycleL | 0.084 | 13.462 | 0.912 | 0.006 | **0.994** | **1.000** | 0.007 | 0.121 | 0.424 | **0.188** |
| CycleL2 | 0.007 | 0.705 | 0.009 | 0.000 | 0.049 | 0.880 | 0.000 | **0.011** | 0.134 | 0.443 |
| HW-TSC | 6.806 | 38.650 | **0.999** | 0.004 | 0.941 | **1.000** | 0.984 | 0.013 | 0.689 | 0.239 |
| IOL_Research | 4.453 | 40.267 | 0.988 | 0.002 | 0.862 | **1.000** | 0.889 | 0.108 | 0.667 | 0.272 |
| ONLINE-A | 1.079 | 42.426 | **0.999** | 0.005 | 0.968 | **1.000** | 0.968 | 0.029 | 0.701 | 0.240 |
| ONLINE-B | 0.964 | **43.437** | **0.999** | 0.002 | 0.976 | **1.000** | 0.974 | 0.021 | **0.704** | 0.239 |
| ONLINE-G | 1.689 | 38.324 | 0.998 | 0.009 | 0.962 | **1.000** | 0.965 | 0.032 | 0.697 | 0.231 |
| ONLINE-W | 2.615 | 38.154 | **0.999** | 0.009 | 0.919 | **1.000** | 0.982 | 0.015 | 0.684 | 0.245 |
| UvA-MT | 0.602 | 37.625 | 0.991 | 0.009 | 0.962 | **1.000** | **0.985** | 0.013 | 0.695 | 0.244 |

Table 40: English→Chinese, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.761 | 24.987 | 0.929 | 0.005 | 0.468 | 0.983 | 0.006 | 0.984 | 0.412 | 0.458 |
| Claude-3 | 0.045 | 11.015 | 0.395 | **0.115** | 0.488 | 0.529 | 0.081 | 0.834 | 0.303 | 0.557 |
| CommandR-plus | 0.049 | 21.207 | 0.486 | 0.004 | 0.494 | 0.810 | 0.037 | 0.930 | 0.315 | 0.523 |
| GPT-4 | 3.758 | **58.453** | **0.999** | 0.007 | **1.000** | **1.000** | 0.091 | 0.882 | **0.587** | 0.351 |
| Llama3-70B | 0.278 | 43.325 | 0.973 | 0.031 | 0.931 | 0.990 | 0.091 | 0.841 | 0.549 | 0.335 |
| NVIDIA-NeMo | 0.270 | 30.835 | 0.968 | 0.005 | 0.834 | 0.979 | 0.006 | 0.880 | 0.512 | 0.364 |
| IKUN | 0.621 | 47.653 | 0.996 | 0.005 | 0.995 | **1.000** | 0.075 | 0.783 | 0.581 | 0.308 |
| IKUN-C | 0.200 | 20.826 | 0.991 | 0.005 | 0.646 | 0.999 | 0.033 | 0.951 | 0.439 | 0.390 |
| Unbabel-Tower70B | 1.298 | 35.764 | 0.996 | 0.001 | 0.876 | **1.000** | 0.027 | 0.967 | 0.544 | 0.384 |
| CycleL | 0.009 | 2.377 | 0.892 | 0.006 | 0.264 | 0.999 | 0.000 | 0.158 | 0.309 | **0.287** |
| CycleL2 | 0.000 | 1.037 | 0.006 | 0.000 | 0.109 | 0.973 | 0.000 | **0.127** | 0.155 | 0.435 |
| HW-TSC | **5.378** | 37.590 | **0.999** | 0.004 | 0.980 | **1.000** | 0.104 | 0.814 | 0.582 | 0.334 |
| IOL_Research | 0.741 | 57.559 | 0.998 | 0.009 | **1.000** | **1.000** | 0.033 | 0.958 | 0.578 | 0.355 |
| ONLINE-A | 1.040 | 46.701 | **0.999** | 0.005 | 0.996 | **1.000** | 0.034 | 0.862 | 0.577 | 0.339 |
| ONLINE-B | 0.053 | 43.261 | **0.999** | 0.005 | 0.999 | **1.000** | 0.047 | 0.922 | 0.579 | 0.380 |
| ONLINE-G | 0.216 | 28.849 | 0.998 | 0.010 | 0.805 | **1.000** | 0.015 | 0.984 | 0.511 | 0.386 |
| ONLINE-W | 0.402 | 33.369 | 0.996 | 0.015 | 0.907 | **1.000** | **0.152** | 0.814 | 0.560 | 0.345 |
| UvA-MT | 0.135 | 26.099 | 0.982 | 0.010 | 0.818 | 0.990 | 0.026 | 0.880 | 0.496 | 0.356 |

Table 41: English→Chinese, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.265 | 19.780 | 0.998 | 0.010 | 0.777 | 0.995 | 0.010 | 0.782 | 0.435 | 0.325 |
| Claude-3 | 0.836 | 51.981 | 0.879 | 0.007 | 0.976 | 0.984 | 0.012 | 0.421 | 0.536 | 0.270 |
| CommandR-plus | 0.054 | 43.878 | 0.881 | 0.006 | 0.916 | 0.949 | 0.012 | 0.465 | 0.513 | 0.281 |
| GPT-4 | 1.904 | 54.773 | **0.999** | 0.006 | **1.000** | **1.000** | 0.017 | 0.306 | 0.576 | 0.230 |
| Llama3-70B | 0.083 | 46.518 | 0.988 | 0.009 | **1.000** | **1.000** | 0.007 | 0.459 | 0.572 | 0.233 |
| NVIDIA-NeMo | 0.606 | 38.384 | 0.994 | 0.005 | **1.000** | **1.000** | 0.001 | 0.711 | 0.573 | 0.284 |
| IKUN | 0.512 | **61.085** | 0.994 | 0.004 | 0.999 | **1.000** | 0.007 | 0.393 | 0.572 | **0.225** |
| IKUN-C | 0.123 | 25.144 | 0.895 | 0.001 | 0.972 | **1.000** | 0.011 | 0.595 | 0.520 | 0.266 |
| Unbabel-Tower70B | 0.170 | 42.703 | 0.957 | 0.006 | 0.999 | **1.000** | 0.013 | 0.510 | 0.568 | 0.261 |
| CycleL | 0.003 | 2.229 | 0.393 | 0.001 | 0.483 | 0.993 | 0.000 | 0.066 | 0.267 | 0.314 |
| CycleL2 | 0.000 | 0.891 | 0.011 | 0.000 | 0.007 | 0.941 | 0.000 | **0.011** | 0.137 | 0.436 |
| HW-TSC | **5.511** | 50.392 | **0.999** | 0.004 | **1.000** | **1.000** | 0.016 | 0.453 | 0.575 | 0.248 |
| IOL_Research | 0.116 | 51.409 | 0.996 | 0.004 | **1.000** | **1.000** | 0.013 | 0.327 | 0.574 | 0.230 |
| ONLINE-A | 1.040 | 52.306 | **0.999** | 0.005 | **1.000** | **1.000** | 0.009 | 0.728 | 0.575 | 0.286 |
| ONLINE-B | 0.031 | 48.627 | **0.999** | 0.005 | **1.000** | **1.000** | **0.031** | 0.318 | **0.576** | 0.230 |
| ONLINE-G | 0.140 | 37.881 | 0.998 | 0.010 | **1.000** | **1.000** | 0.020 | 0.821 | 0.575 | 0.297 |
| ONLINE-W | 0.129 | 44.726 | 0.998 | **0.012** | **1.000** | **1.000** | 0.006 | 0.472 | 0.574 | 0.235 |
| UvA-MT | 0.047 | 24.717 | 0.958 | 0.009 | 0.987 | **1.000** | 0.001 | 0.535 | 0.543 | 0.247 |

Table 42: English→Chinese, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 26.338 | 30.138 | 0.905 | 0.049 | 0.612 | 0.753 | 0.668 | 0.064 | 0.514 | 0.366 |
| Claude-3 | 29.992 | 30.134 | 0.909 | **0.201** | 0.829 | 0.902 | 0.875 | 0.075 | 0.678 | 0.289 |
| CommandR-plus | 22.717 | 25.593 | 0.949 | 0.053 | 0.552 | 0.672 | 0.614 | 0.086 | 0.490 | 0.415 |
| GPT-4 | **37.792** | 42.756 | **0.999** | 0.018 | 0.173 | 0.184 | 0.186 | 0.033 | 0.252 | 0.616 |
| Llama3-70B | 34.428 | 43.116 | 0.947 | 0.105 | 0.903 | 0.979 | 0.930 | 0.028 | 0.679 | 0.241 |
| NVIDIA-NeMo | 30.861 | 28.789 | 0.962 | 0.033 | 0.670 | 0.761 | 0.742 | 0.035 | 0.531 | 0.348 |
| IKUN | 34.710 | **43.744** | 0.987 | 0.051 | 0.841 | **0.989** | 0.947 | 0.028 | 0.646 | 0.248 |
| IKUN-C | 34.992 | 39.814 | 0.974 | 0.037 | 0.732 | 0.892 | 0.834 | 0.047 | 0.573 | 0.287 |
| Unbabel-Tower70B | 34.411 | 41.251 | 0.994 | 0.050 | 0.568 | 0.645 | 0.632 | 0.027 | 0.499 | 0.407 |
| CycleL | 0.000 | 3.706 | 0.001 | 0.000 | 0.002 | 0.330 | 0.000 | **0.000** | 0.048 | 0.524 |
| CycleL2 | 0.000 | 0.584 | 0.004 | 0.000 | 0.044 | 0.632 | 0.000 | 0.004 | 0.097 | 0.480 |
| HW-TSC | 0.000 | 8.198 | 0.321 | 0.015 | 0.170 | 0.776 | 0.130 | 0.076 | 0.216 | 0.444 |
| IOL_Research | 9.759 | 15.764 | 0.903 | 0.098 | 0.720 | 0.827 | 0.750 | 0.077 | 0.577 | 0.306 |
| ONLINE-A | 1.980 | 17.699 | 0.869 | 0.106 | 0.820 | 0.853 | 0.815 | 0.033 | 0.607 | 0.276 |
| ONLINE-B | 35.639 | 22.384 | 0.996 | 0.104 | **0.916** | 0.976 | 0.947 | 0.033 | 0.691 | **0.239** |
| ONLINE-G | 34.841 | 29.900 | 0.996 | 0.130 | 0.870 | 0.983 | **0.949** | 0.033 | **0.692** | 0.244 |
| ONLINE-W | 23.692 | 23.809 | 0.847 | 0.018 | 0.233 | 0.370 | 0.241 | 0.045 | 0.274 | 0.550 |
| UvA-MT | 37.581 | 39.663 | 0.931 | 0.004 | 0.272 | 0.359 | 0.344 | 0.206 | 0.304 | 0.585 |

Table 43: English→Chinese, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 12.503 | 28.651 | 0.999 | 0.106 | 0.887 | 0.982 | 0.955 | 0.027 | 0.687 | 0.236 |
| Claude-3 | 0.305 | 8.270 | 0.015 | 0.001 | 0.111 | 0.670 | 0.021 | 0.945 | 0.127 | 0.689 |
| CommandR-plus | 9.258 | 20.579 | 0.858 | 0.078 | 0.679 | 0.862 | 0.722 | 0.129 | 0.554 | 0.341 |
| GPT-4 | 20.497 | 31.347 | 0.999 | 0.002 | 0.006 | 0.002 | 0.021 | 0.032 | 0.150 | 0.691 |
| Llama3-70B | 14.202 | 31.556 | 0.994 | 0.126 | 0.909 | 0.987 | **0.961** | 0.017 | 0.706 | 0.240 |
| NVIDIA-NeMo | 20.277 | 20.133 | 0.035 | 0.002 | 0.009 | 0.015 | 0.001 | 0.207 | 0.010 | 0.621 |
| IKUN | 15.178 | **33.436** | 0.989 | 0.038 | 0.802 | 0.976 | 0.957 | 0.023 | 0.638 | 0.259 |
| IKUN-C | 13.816 | 29.137 | 0.987 | 0.035 | 0.796 | 0.955 | 0.897 | 0.037 | 0.602 | 0.247 |
| Unbabel-Tower70B | 18.475 | 31.454 | 0.998 | 0.054 | 0.493 | 0.551 | 0.545 | 0.027 | 0.453 | 0.445 |
| CycleL | 0.000 | 2.247 | 0.000 | 0.000 | 0.004 | 0.345 | 0.000 | **0.001** | 0.050 | 0.523 |
| CycleL2 | 0.000 | 0.327 | 0.005 | 0.000 | 0.054 | 0.646 | 0.000 | 0.004 | 0.101 | 0.475 |
| HW-TSC | 0.000 | 3.635 | 0.000 | 0.000 | 0.109 | **0.999** | 0.000 | 0.017 | 0.158 | 0.416 |
| IOL_Research | 6.086 | 15.506 | 0.960 | 0.105 | 0.780 | 0.842 | 0.825 | 0.064 | 0.620 | 0.284 |
| ONLINE-A | 0.241 | 9.912 | **1.000** | 0.106 | 0.841 | 0.984 | 0.815 | 0.042 | 0.647 | 0.237 |
| ONLINE-B | 16.305 | 16.028 | 0.999 | 0.116 | **0.914** | 0.976 | 0.956 | 0.021 | 0.697 | **0.236** |
| ONLINE-G | 15.830 | 18.245 | 0.999 | **0.148** | 0.896 | 0.994 | 0.953 | 0.031 | **0.709** | 0.237 |
| ONLINE-W | 12.605 | 17.025 | 0.111 | 0.005 | 0.188 | 0.974 | 0.094 | 0.300 | 0.206 | 0.441 |
| UvA-MT | **21.160** | 29.706 | 0.015 | 0.000 | 0.000 | 0.000 | 0.000 | 0.136 | 0.002 | 0.607 |

Table 44: English→Chinese, 1-shot JSON format

417

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 49.791 | 70.611 | **1.000** | 0.905 | 0.914 | 0.999 | 0.972 | 0.023 | 0.947 | 0.251 |
| Claude-3 | **62.653** | **79.565** | 0.998 | **0.945** | 0.968 | **1.000** | 0.966 | 0.034 | **0.973** | 0.253 |
| CommandR-plus | 52.187 | 72.206 | 0.995 | 0.927 | 0.944 | 0.999 | 0.971 | 0.029 | 0.958 | 0.249 |
| GPT-4 | 54.848 | 74.440 | 0.999 | 0.931 | **0.969** | 0.999 | 0.966 | 0.034 | 0.968 | 0.246 |
| Llama3-70B | 49.780 | 70.822 | 0.999 | 0.902 | 0.938 | 0.999 | 0.971 | 0.027 | 0.949 | 0.251 |
| NVIDIA-NeMo | 47.690 | 67.971 | 0.969 | 0.905 | 0.936 | **1.000** | 0.967 | **0.010** | 0.940 | 0.247 |
| IKUN | 34.437 | 58.673 | 0.988 | 0.868 | 0.914 | 0.999 | 0.949 | 0.044 | 0.908 | 0.241 |
| IKUN-C | 36.359 | 59.479 | 0.993 | 0.881 | 0.929 | **1.000** | 0.965 | 0.027 | 0.918 | **0.234** |
| Unbabel-Tower70B | 49.401 | 71.358 | 0.991 | 0.911 | 0.947 | 0.999 | 0.962 | 0.034 | 0.954 | 0.251 |
| CycleL | 0.928 | 14.750 | 0.000 | 0.290 | 0.446 | 0.999 | 0.040 | 0.047 | 0.270 | 0.359 |
| Dubformer | 49.968 | 71.853 | 0.987 | 0.918 | 0.946 | 0.998 | 0.971 | 0.029 | 0.956 | 0.250 |
| IOL_Research | 59.265 | 76.482 | 0.979 | 0.919 | 0.944 | **1.000** | **0.984** | 0.016 | 0.960 | 0.256 |
| ONLINE-A | 53.505 | 72.787 | 0.996 | 0.930 | 0.949 | 0.999 | 0.966 | 0.028 | 0.956 | 0.248 |
| ONLINE-B | 52.012 | 72.139 | 0.998 | 0.917 | 0.946 | **1.000** | 0.972 | 0.026 | 0.956 | 0.250 |
| ONLINE-G | 47.843 | 70.719 | 0.996 | 0.917 | 0.938 | **1.000** | 0.962 | 0.034 | 0.951 | 0.251 |
| ONLINE-W | 56.473 | 74.051 | 0.999 | 0.928 | 0.944 | **1.000** | 0.965 | 0.033 | 0.959 | 0.252 |
| TranssionMT | 54.465 | 74.167 | 0.995 | 0.935 | 0.951 | **1.000** | 0.969 | 0.027 | 0.961 | 0.251 |

Table 45: English→Ukrainian, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 43.208 | 62.841 | 0.998 | 0.767 | 0.698 | 0.999 | 0.635 | 0.365 | 0.832 | 0.331 |
| Claude-3 | 0.192 | 0.963 | 0.024 | 0.007 | 0.004 | 0.012 | 0.009 | 0.984 | 0.009 | 0.838 |
| CommandR-plus | 43.779 | 70.052 | 0.884 | 0.869 | 0.863 | 0.958 | 0.756 | 0.225 | 0.872 | 0.280 |
| GPT-4 | 22.972 | 42.060 | 0.621 | 0.430 | 0.335 | 0.961 | 0.305 | 0.694 | 0.566 | 0.547 |
| Llama3-70B | 2.941 | 11.270 | 0.251 | 0.166 | 0.143 | 0.211 | 0.127 | 0.873 | 0.181 | 0.738 |
| NVIDIA-NeMo | **54.317** | **77.552** | 0.983 | 0.999 | **0.999** | **1.000** | 0.875 | 0.125 | 0.979 | 0.210 |
| IKUN | 27.427 | 61.364 | 0.928 | 0.972 | 0.987 | **1.000** | **0.924** | **0.076** | 0.933 | 0.193 |
| IKUN-C | 24.366 | 57.665 | 0.995 | 0.960 | 0.965 | **1.000** | 0.916 | 0.084 | 0.916 | **0.187** |
| Unbabel-Tower70B | 38.592 | 73.353 | 0.991 | 0.995 | **0.999** | **1.000** | 0.912 | 0.088 | 0.983 | 0.198 |
| CycleL | 0.493 | 15.506 | 0.000 | 0.603 | 0.567 | **1.000** | 0.001 | 0.082 | 0.313 | 0.293 |
| Dubformer | 15.405 | 34.623 | 0.523 | 0.454 | 0.466 | 0.700 | 0.280 | 0.610 | 0.482 | 0.545 |
| IOL_Research | 36.206 | 54.753 | 0.973 | 0.638 | 0.493 | **1.000** | 0.499 | 0.498 | 0.753 | 0.412 |
| ONLINE-A | 47.835 | 76.254 | 0.995 | **1.000** | 0.999 | **1.000** | 0.764 | 0.234 | 0.965 | 0.225 |
| ONLINE-B | 50.403 | 77.296 | 0.998 | 0.998 | **0.999** | **1.000** | 0.923 | 0.077 | **0.988** | 0.198 |
| ONLINE-G | 50.344 | 75.798 | **0.999** | 0.999 | **0.999** | **1.000** | 0.880 | 0.120 | 0.979 | 0.202 |
| ONLINE-W | 48.888 | 75.292 | 0.906 | 0.999 | 0.998 | **1.000** | 0.882 | 0.118 | 0.969 | 0.227 |
| TranssionMT | 47.024 | 76.579 | 0.995 | 0.999 | 0.998 | **1.000** | 0.802 | 0.198 | 0.970 | 0.218 |

Table 46: English→Ukrainian, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 38.170 | 66.310 | 0.969 | 0.990 | 0.991 | **1.000** | 0.004 | 0.965 | 0.842 | 0.315 |
| Claude-3 | 34.875 | 63.514 | 0.918 | 0.939 | 0.940 | 0.955 | **0.202** | 0.558 | 0.823 | 0.280 |
| CommandR-plus | 18.276 | 41.153 | 0.395 | 0.849 | 0.668 | 0.913 | 0.009 | 0.979 | 0.592 | 0.462 |
| GPT-4 | 40.735 | 66.885 | 0.996 | **1.000** | **1.000** | **1.000** | 0.011 | 0.962 | 0.857 | 0.301 |
| Llama3-70B | 22.647 | 41.801 | 0.770 | 0.499 | 0.395 | 0.973 | 0.011 | 0.976 | 0.547 | 0.532 |
| NVIDIA-NeMo | **44.936** | 66.430 | 0.994 | 0.996 | 0.991 | **1.000** | 0.002 | 0.994 | 0.854 | 0.312 |
| IKUN | 31.820 | 64.279 | 0.996 | **1.000** | 0.999 | **1.000** | 0.040 | 0.521 | 0.861 | **0.217** |
| IKUN-C | 20.981 | 50.000 | 0.996 | 0.925 | 0.860 | **1.000** | 0.012 | 0.983 | 0.759 | 0.338 |
| Unbabel-Tower70B | 26.806 | 54.532 | 0.982 | 0.941 | 0.873 | 0.999 | 0.023 | 0.956 | 0.792 | 0.348 |
| CycleL | 0.190 | 11.645 | 0.000 | 0.953 | 0.843 | 0.857 | 0.000 | 0.346 | 0.380 | 0.247 |
| Dubformer | 16.325 | 24.645 | 0.984 | 0.568 | 0.231 | 0.245 | 0.010 | **0.246** | 0.359 | 0.528 |
| IOL_Research | 40.166 | **73.445** | 0.990 | 0.995 | 0.995 | **1.000** | 0.033 | 0.646 | 0.858 | 0.258 |
| ONLINE-A | 40.582 | 72.313 | 0.996 | **1.000** | **1.000** | **1.000** | 0.020 | 0.645 | 0.859 | 0.250 |
| ONLINE-B | 37.933 | 72.460 | 0.998 | **1.000** | **1.000** | **1.000** | 0.035 | 0.635 | **0.862** | 0.252 |
| ONLINE-G | 26.133 | 59.435 | **0.999** | 0.983 | 0.995 | **1.000** | 0.005 | 0.398 | 0.808 | 0.217 |
| ONLINE-W | 37.025 | 71.965 | **0.999** | 0.996 | **1.000** | **1.000** | 0.033 | 0.394 | 0.856 | 0.217 |
| TranssionMT | 40.563 | 72.472 | 0.996 | **1.000** | **1.000** | **1.000** | 0.018 | 0.644 | 0.859 | 0.251 |

Table 47: English→Ukrainian, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 46.559 | 73.107 | **1.000** | **1.000** | **1.000** | **1.000** | 0.029 | 0.787 | 0.861 | 0.221 |
| Claude-3 | 2.037 | 20.426 | 0.198 | 0.889 | 0.496 | 0.983 | 0.027 | 0.244 | 0.412 | 0.272 |
| CommandR-plus | 50.447 | 71.338 | 0.827 | 0.998 | 0.998 | **1.000** | 0.012 | 0.870 | 0.830 | 0.258 |
| GPT-4 | 49.840 | 73.623 | 0.994 | **1.000** | **1.000** | **1.000** | 0.032 | 0.792 | 0.861 | 0.222 |
| Llama3-70B | 36.700 | 60.591 | 0.969 | **1.000** | **1.000** | **1.000** | **0.173** | 0.622 | **0.877** | 0.215 |
| NVIDIA-NeMo | 49.107 | 70.904 | 0.958 | **1.000** | **1.000** | **1.000** | 0.013 | 0.909 | 0.853 | 0.248 |
| IKUN | 70.906 | 82.849 | **1.000** | 0.996 | 0.972 | **1.000** | 0.024 | 0.542 | 0.853 | **0.174** |
| IKUN-C | 30.631 | 60.900 | 0.994 | 0.999 | 0.998 | **1.000** | 0.021 | 0.916 | 0.854 | 0.239 |
| Unbabel-Tower70B | 47.335 | 72.274 | 0.996 | **1.000** | **1.000** | **1.000** | 0.035 | 0.892 | 0.862 | 0.246 |
| CycleL | 0.255 | 17.188 | 0.000 | 0.841 | 0.977 | 0.944 | 0.000 | 0.099 | 0.395 | 0.192 |
| Dubformer | 5.538 | 8.974 | 0.998 | 0.491 | 0.049 | 0.037 | 0.006 | **0.043** | 0.236 | 0.496 |
| IOL_Research | 65.240 | 84.223 | 0.996 | **1.000** | **1.000** | **1.000** | 0.015 | 0.672 | 0.859 | 0.198 |
| ONLINE-A | 64.032 | 84.119 | 0.996 | **1.000** | **1.000** | **1.000** | 0.029 | 0.546 | 0.861 | 0.178 |
| ONLINE-B | **74.871** | **88.689** | 0.998 | **1.000** | **1.000** | **1.000** | 0.022 | 0.550 | 0.860 | 0.177 |
| ONLINE-G | 59.770 | 82.970 | 0.994 | **1.000** | **1.000** | **1.000** | 0.024 | 0.570 | 0.860 | 0.179 |
| ONLINE-W | 53.733 | 80.823 | 0.999 | **1.000** | **1.000** | **1.000** | 0.026 | 0.640 | 0.861 | 0.194 |
| TranssionMT | 65.081 | 84.691 | 0.996 | **1.000** | **1.000** | **1.000** | 0.033 | 0.537 | 0.861 | 0.179 |

Table 48: English→Ukrainian, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 4.388 | 22.288 | 0.995 | 0.913 | 0.919 | 0.993 | 0.958 | 0.026 | 0.942 | **0.248** |
| Claude-3 | **22.891** | **39.389** | 0.998 | **0.950** | **0.968** | 0.998 | 0.966 | 0.029 | **0.976** | 0.255 |
| CommandR-plus | 0.126 | 9.897 | 0.748 | 0.689 | 0.722 | 0.979 | 0.733 | 0.160 | 0.759 | 0.340 |
| GPT-4 | 13.402 | 32.382 | **0.999** | 0.297 | 0.279 | 0.285 | 0.329 | 0.213 | 0.396 | 0.588 |
| Llama3-70B | 11.288 | 29.999 | **0.999** | 0.903 | 0.942 | **1.000** | **0.968** | 0.031 | 0.951 | 0.252 |
| NVIDIA-NeMo | 12.668 | 27.319 | 0.428 | 0.492 | 0.519 | 0.974 | 0.493 | 0.166 | 0.543 | 0.368 |
| IKUN | 6.190 | 25.561 | 0.965 | 0.875 | 0.890 | 0.971 | 0.894 | 0.070 | 0.885 | 0.251 |
| IKUN-C | 4.751 | 23.331 | 0.951 | 0.734 | 0.761 | 0.842 | 0.818 | 0.089 | 0.783 | 0.310 |
| Unbabel-Tower70B | 5.221 | 23.735 | 0.977 | 0.903 | 0.930 | 0.989 | 0.939 | 0.050 | 0.938 | 0.257 |
| CycleL | 0.000 | 3.299 | 0.000 | 0.118 | 0.168 | 0.982 | 0.000 | **0.000** | 0.181 | 0.390 |
| Dubformer | 0.707 | 12.279 | 0.384 | 0.037 | 0.027 | 0.208 | 0.039 | 0.655 | 0.107 | 0.761 |
| IOL_Research | 3.439 | 15.528 | 0.769 | 0.639 | 0.643 | 0.873 | 0.694 | 0.031 | 0.684 | 0.329 |
| ONLINE-A | 13.719 | 35.538 | 0.966 | 0.906 | 0.930 | 0.998 | 0.941 | 0.056 | 0.942 | 0.263 |
| ONLINE-B | 13.641 | 22.733 | 0.987 | 0.905 | 0.935 | 0.991 | 0.955 | 0.038 | 0.945 | 0.255 |
| ONLINE-G | 11.368 | 26.091 | 0.148 | 0.108 | 0.211 | 0.753 | 0.119 | 0.103 | 0.222 | 0.446 |
| ONLINE-W | 13.269 | 31.455 | 0.996 | 0.933 | 0.938 | 0.994 | 0.950 | 0.043 | 0.956 | 0.255 |
| TranssionMT | 13.692 | 35.536 | 0.991 | 0.922 | 0.942 | 0.998 | 0.961 | 0.033 | 0.957 | 0.255 |

Table 49: English→Ukrainian, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 1.063 | 14.715 | 0.998 | 0.917 | 0.945 | **0.999** | **0.973** | 0.022 | 0.958 | 0.246 |
| Claude-3 | **4.693** | 19.766 | 0.823 | 0.092 | 0.092 | 0.182 | 0.137 | 0.388 | 0.229 | 0.715 |
| CommandR-plus | 0.040 | 7.585 | 0.288 | 0.360 | 0.400 | 0.972 | 0.351 | 0.628 | 0.492 | 0.580 |
| GPT-4 | 3.118 | 20.202 | **0.999** | 0.043 | 0.009 | 0.006 | 0.086 | 0.278 | 0.169 | 0.720 |
| Llama3-70B | 2.764 | 19.104 | 0.962 | 0.903 | 0.934 | 0.994 | 0.950 | 0.047 | 0.936 | 0.258 |
| NVIDIA-NeMo | 1.287 | 19.588 | 0.987 | **0.985** | 0.632 | **0.999** | 0.001 | 0.005 | 0.516 | **0.062** |
| IKUN | 1.227 | 16.163 | 0.940 | 0.860 | 0.882 | 0.968 | 0.905 | 0.048 | 0.877 | 0.254 |
| IKUN-C | 0.886 | 14.892 | 0.923 | 0.849 | 0.852 | 0.990 | 0.849 | 0.072 | 0.860 | 0.256 |
| Unbabel-Tower70B | 1.202 | 15.677 | 0.991 | 0.918 | **0.962** | 0.998 | 0.965 | 0.029 | **0.960** | 0.249 |
| CycleL | 0.000 | 1.701 | 0.000 | 0.122 | 0.168 | 0.988 | 0.000 | **0.001** | 0.183 | 0.389 |
| Dubformer | 0.695 | 12.893 | 0.600 | 0.043 | 0.011 | 0.091 | 0.049 | 0.474 | 0.116 | 0.735 |
| IOL_Research | 0.689 | 9.081 | 0.130 | 0.106 | 0.228 | 0.816 | 0.105 | 0.070 | 0.227 | 0.429 |
| ONLINE-A | 3.114 | **21.522** | 0.966 | 0.908 | 0.931 | 0.998 | 0.941 | 0.049 | 0.943 | 0.262 |
| ONLINE-B | 3.164 | 12.686 | 0.991 | 0.909 | 0.942 | 0.996 | 0.965 | 0.033 | 0.950 | 0.253 |
| ONLINE-G | 2.564 | 13.911 | 0.094 | 0.076 | 0.207 | 0.837 | 0.077 | 0.043 | 0.205 | 0.428 |
| ONLINE-W | 3.092 | 18.793 | 0.996 | 0.917 | 0.929 | 0.991 | 0.945 | 0.049 | 0.948 | 0.261 |
| TranssionMT | 3.155 | 21.520 | 0.993 | 0.913 | 0.944 | **0.999** | 0.967 | 0.029 | 0.956 | 0.256 |

Table 50: English→Ukrainian, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 13.449 | 35.122 | 0.996 | 0.820 | 0.897 | 0.972 | 0.933 | 0.005 | 0.817 | 0.167 |
| Claude-3 | 55.420 | **75.544** | 0.999 | **0.945** | **0.974** | 0.988 | 0.996 | 0.004 | 0.971 | 0.245 |
| CommandR-plus | 20.222 | 44.344 | 0.509 | 0.847 | 0.887 | 0.990 | 0.962 | 0.023 | 0.798 | 0.273 |
| GPT-4 | 42.953 | 65.458 | **1.000** | 0.909 | 0.963 | 0.989 | 0.996 | 0.004 | 0.951 | 0.233 |
| Llama3-70B | 38.608 | 60.739 | 0.991 | 0.896 | 0.946 | 0.989 | 0.996 | 0.001 | 0.936 | 0.229 |
| IKUN | 31.698 | 55.417 | 0.967 | 0.749 | 0.832 | 0.990 | 0.950 | 0.018 | 0.865 | 0.273 |
| IKUN-C | 25.692 | 49.700 | 0.983 | 0.733 | 0.824 | 0.990 | 0.945 | 0.029 | 0.839 | 0.251 |
| Unbabel-Tower70B | 44.358 | 67.090 | 0.999 | 0.897 | 0.949 | **0.995** | 0.991 | 0.009 | 0.949 | 0.244 |
| AMI | 52.729 | 72.148 | 0.998 | 0.940 | 0.953 | **0.995** | 1.000 | **0.000** | 0.970 | 0.245 |
| CycleL | 10.383 | 29.998 | 0.929 | 0.786 | 0.875 | 0.994 | 0.460 | 0.017 | 0.699 | **0.150** |
| Dubformer | 41.037 | 61.391 | 0.978 | 0.874 | 0.912 | 0.953 | 0.955 | 0.022 | 0.914 | 0.260 |
| IOL_Research | 45.690 | 64.846 | 0.988 | 0.879 | 0.929 | 0.989 | 0.993 | 0.005 | 0.941 | 0.253 |
| ONLINE-A | 55.587 | 73.600 | 0.999 | 0.930 | 0.957 | 0.990 | 0.998 | 0.001 | 0.968 | 0.249 |
| ONLINE-B | 57.116 | 73.904 | 0.998 | 0.942 | 0.963 | 0.991 | **1.000** | **0.000** | **0.974** | 0.248 |
| ONLINE-G | 47.642 | 67.534 | 0.998 | 0.906 | 0.938 | 0.991 | 0.989 | 0.001 | 0.951 | 0.246 |
| ONLINE-W | | | | | | NA | | | | |
| TSU-HITs | 8.553 | 28.192 | 0.317 | 0.493 | 0.732 | 0.979 | 0.676 | 0.023 | 0.570 | 0.337 |
| TranssionMT | **57.314** | 74.708 | 0.999 | 0.940 | 0.965 | 0.990 | **1.000** | **0.000** | 0.973 | 0.249 |

Table 51: English→Icelandic, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 7.573 | 32.205 | 0.897 | 0.889 | 0.827 | 0.990 | 0.750 | 0.160 | 0.752 | 0.175 |
| Claude-3 | 0.417 | 10.167 | 0.010 | 0.037 | 0.047 | 0.076 | 0.002 | 0.996 | 0.033 | 0.822 |
| CommandR-plus | 2.692 | 17.322 | 0.337 | 0.196 | 0.180 | 0.370 | 0.195 | 0.605 | 0.227 | 0.667 |
| GPT-4 | 15.660 | 33.955 | 0.480 | 0.293 | 0.223 | 0.971 | 0.411 | 0.588 | 0.482 | 0.582 |
| Llama3-70B | 0.709 | 11.013 | 0.076 | 0.069 | 0.067 | 0.078 | 0.042 | 0.958 | 0.064 | 0.800 |
| IKUN | 42.666 | 67.063 | 0.854 | 0.998 | 0.990 | **1.000** | 0.996 | 0.002 | 0.972 | 0.212 |
| IKUN-C | 38.746 | 63.561 | 0.983 | 0.996 | 0.987 | 0.999 | 0.990 | 0.009 | 0.988 | 0.179 |
| Unbabel-Tower70B | 39.320 | 65.432 | 0.917 | 0.988 | 0.982 | **1.000** | 0.963 | 0.037 | 0.972 | 0.205 |
| AMI | 54.415 | 74.927 | 0.998 | 0.999 | 0.998 | **1.000** | **0.999** | 0.001 | 0.997 | 0.192 |
| CycleL | 8.093 | 30.233 | 0.958 | 0.933 | 0.929 | **1.000** | 0.048 | 0.332 | 0.674 | **0.142** |
| Dubformer | 11.780 | 31.937 | 0.433 | 0.452 | 0.438 | 0.644 | 0.356 | 0.576 | 0.456 | 0.538 |
| IOL_Research | 21.003 | 41.958 | 0.996 | 0.465 | 0.409 | 0.994 | 0.815 | 0.181 | 0.720 | 0.398 |
| ONLINE-A | 58.224 | 76.680 | **0.999** | 1.000 | 1.000 | **1.000** | 0.994 | 0.005 | 0.997 | 0.195 |
| ONLINE-B | **61.327** | 78.012 | 0.996 | 0.999 | 0.999 | **1.000** | **0.999** | 0.000 | **0.999** | 0.194 |
| ONLINE-G | 48.915 | 70.410 | 0.998 | 0.996 | 1.000 | **1.000** | 0.993 | 0.002 | 0.998 | 0.185 |
| ONLINE-W | | | | | | NA | | | | |
| TSU-HITs | 3.107 | 17.806 | 0.394 | 0.416 | 0.379 | 0.996 | 0.124 | 0.403 | 0.360 | 0.370 |
| TranssionMT | 61.273 | **78.416** | 0.998 | 0.999 | 0.999 | **1.000** | 0.996 | 0.001 | 0.999 | 0.194 |

Table 52: English→Icelandic, direct

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 22.589 | 40.223 | 0.979 | 0.958 | 0.890 | 0.919 | 0.044 | 0.789 | 0.767 | 0.244 |
| Claude-3 | 2.933 | 15.668 | 0.006 | 0.364 | 0.428 | 0.299 | 0.125 | 0.792 | 0.187 | 0.603 |
| CommandR-plus | 3.358 | 19.349 | 0.042 | 0.280 | 0.263 | 0.935 | 0.011 | 0.983 | 0.252 | 0.625 |
| GPT-4 | 51.697 | 66.383 | 0.998 | 0.965 | 0.938 | 0.999 | 0.482 | 0.493 | 0.900 | 0.247 |
| Llama3-70B | 16.466 | 35.615 | 0.535 | 0.517 | 0.472 | 0.931 | 0.115 | 0.864 | 0.513 | 0.511 |
| IKUN | **71.668** | **78.373** | 0.996 | **1.000** | **1.000** | **1.000** | 0.186 | 0.802 | 0.882 | 0.266 |
| IKUN-C | 22.812 | 48.582 | 0.998 | 0.935 | 0.874 | **1.000** | 0.094 | 0.903 | 0.793 | 0.310 |
| Unbabel-Tower70B | 33.879 | 55.577 | 0.996 | 0.938 | 0.843 | **1.000** | 0.212 | 0.760 | 0.829 | 0.318 |
| AMI | 55.611 | 71.232 | 0.998 | 0.994 | 0.976 | **1.000** | 0.550 | 0.411 | 0.930 | 0.233 |
| CycleL | 9.513 | 29.310 | 0.957 | 0.995 | 0.961 | **1.000** | 0.004 | 0.453 | 0.696 | **0.133** |
| Dubformer | 17.515 | 27.491 | 0.879 | 0.627 | 0.293 | 0.202 | 0.022 | 0.289 | 0.455 | 0.521 |
| IOL_Research | 58.323 | 71.321 | 0.995 | **1.000** | **1.000** | **1.000** | 0.772 | 0.196 | 0.967 | 0.182 |
| ONLINE-A | 64.175 | 76.302 | **0.999** | **1.000** | 0.999 | **1.000** | 0.599 | 0.364 | 0.942 | 0.226 |
| ONLINE-B | 64.864 | 75.933 | 0.998 | **1.000** | **1.000** | **1.000** | **0.814** | **0.170** | **0.973** | 0.201 |
| ONLINE-G | 36.759 | 58.436 | **0.999** | 0.998 | 0.989 | **1.000** | 0.165 | 0.603 | 0.870 | 0.252 |
| ONLINE-W | | | | | | NA | | | | |
| TSU-HITs | 2.741 | 16.834 | 0.656 | 0.319 | 0.257 | 0.998 | 0.002 | 0.858 | 0.344 | 0.449 |
| TranssionMT | 64.665 | 76.137 | **0.999** | **1.000** | **1.000** | **1.000** | 0.807 | 0.177 | 0.972 | 0.202 |

Table 53: English→Icelandic, 0-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 16.281 | 38.598 | 0.977 | 0.993 | 0.989 | 0.999 | 0.069 | 0.390 | 0.805 | **0.093** |
| Claude-3 | 0.489 | 7.905 | 0.005 | 0.475 | 0.517 | 0.911 | 0.013 | 0.130 | 0.274 | 0.357 |
| CommandR-plus | 6.805 | 29.615 | 0.129 | 0.834 | 0.814 | **1.000** | 0.053 | 0.846 | 0.523 | 0.373 |
| GPT-4 | 42.211 | 69.479 | **1.000** | **1.000** | **1.000** | **1.000** | 0.455 | 0.083 | 0.922 | 0.112 |
| Llama3-70B | 36.874 | 64.809 | 0.944 | **1.000** | **1.000** | **1.000** | 0.543 | 0.088 | 0.927 | 0.111 |
| IKUN | **86.193** | 88.623 | 0.998 | **1.000** | 0.990 | **1.000** | 0.341 | 0.098 | 0.904 | 0.093 |
| IKUN-C | 32.649 | 58.408 | 0.996 | **1.000** | **1.000** | **1.000** | 0.401 | 0.250 | 0.912 | 0.133 |
| Unbabel-Tower70B | 43.206 | 67.960 | 0.994 | **1.000** | **1.000** | **1.000** | 0.461 | 0.132 | 0.921 | 0.138 |
| AMI | 58.994 | 75.683 | 0.998 | **1.000** | **1.000** | **1.000** | 0.378 | 0.038 | 0.911 | 0.128 |
| CycleL | 4.677 | 28.065 | 0.239 | 0.996 | 0.996 | **1.000** | 0.043 | 0.274 | 0.595 | 0.165 |
| Dubformer | 6.587 | 19.173 | 0.996 | 0.640 | 0.246 | 0.002 | 0.001 | 0.344 | 0.271 | 0.483 |
| IOL_Research | 59.652 | 75.700 | 0.999 | **1.000** | **1.000** | **1.000** | **0.589** | 0.047 | **0.941** | 0.101 |
| ONLINE-A | 85.107 | **90.140** | 0.999 | **1.000** | **1.000** | **1.000** | 0.443 | **0.017** | 0.920 | 0.116 |
| ONLINE-B | 85.157 | 89.894 | 0.998 | **1.000** | **1.000** | **1.000** | 0.405 | 0.034 | 0.915 | 0.116 |
| ONLINE-G | 53.725 | 74.475 | 0.999 | **1.000** | **1.000** | **1.000** | 0.343 | 0.168 | 0.906 | 0.124 |
| ONLINE-W | | | | | NA | | | | | |
| TSU-HITs | 2.834 | 23.385 | 0.089 | 0.867 | 0.864 | 0.999 | 0.001 | 0.318 | 0.467 | 0.227 |
| TranssionMT | 85.075 | 90.014 | 0.999 | **1.000** | **1.000** | **1.000** | 0.411 | 0.033 | 0.916 | 0.116 |

Table 54: English→Icelandic, 1-shot

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 37.596 | 36.308 | 0.923 | 0.591 | 0.634 | 0.734 | 0.662 | 0.049 | 0.617 | 0.307 |
| Claude-3 | 43.969 | 45.621 | 0.980 | 0.925 | 0.956 | 0.973 | 0.979 | 0.020 | 0.956 | 0.260 |
| CommandR-plus | 36.022 | 34.809 | 0.799 | 0.520 | 0.575 | 0.720 | 0.621 | 0.182 | 0.591 | 0.400 |
| GPT-4 | 44.196 | 47.459 | **1.000** | 0.727 | 0.763 | 0.771 | 0.780 | 0.002 | 0.784 | 0.332 |
| Llama3-70B | 38.483 | 43.123 | 0.996 | 0.897 | 0.929 | 0.985 | 0.988 | 0.004 | 0.928 | **0.231** |
| IKUN | 42.347 | 49.631 | 0.778 | 0.718 | 0.767 | 0.942 | 0.797 | 0.055 | 0.759 | 0.275 |
| IKUN-C | 28.758 | 39.915 | 0.892 | 0.704 | 0.760 | 0.858 | 0.837 | 0.050 | 0.767 | 0.295 |
| Unbabel-Tower70B | 43.569 | 47.304 | 0.994 | 0.922 | **0.957** | 0.990 | 0.991 | 0.002 | **0.961** | 0.248 |
| AMI | 41.397 | 48.607 | 0.829 | 0.813 | 0.848 | 0.994 | 0.868 | 0.002 | 0.860 | 0.284 |
| CycleL | 11.962 | 22.576 | 0.000 | 0.072 | 0.120 | 0.428 | 0.000 | **0.000** | 0.089 | 0.485 |
| Dubformer | 12.767 | 21.934 | 0.233 | 0.061 | 0.059 | 0.283 | 0.004 | 0.146 | 0.094 | 0.576 |
| IOL_Research | 17.865 | 31.474 | 0.995 | 0.882 | 0.940 | 0.988 | 0.989 | 0.005 | 0.938 | 0.244 |
| ONLINE-A | 39.385 | 45.191 | 0.994 | 0.920 | 0.944 | **0.995** | 0.991 | 0.001 | 0.961 | 0.251 |
| ONLINE-B | **52.890** | 35.070 | 0.980 | 0.923 | 0.945 | 0.978 | 0.985 | 0.002 | 0.955 | 0.254 |
| ONLINE-G | 37.199 | 45.364 | 0.998 | 0.906 | 0.951 | 0.993 | **0.998** | 0.001 | 0.959 | 0.244 |
| ONLINE-W | | | | | NA | | | | | |
| TSU-HITs | 0.000 | 1.387 | 0.004 | 0.054 | 0.175 | 0.947 | 0.000 | 0.821 | 0.169 | 0.567 |
| TranssionMT | 52.887 | **57.312** | 0.977 | **0.931** | 0.945 | 0.991 | 0.987 | 0.005 | 0.960 | 0.255 |

Table 55: English→Icelandic, 0-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 25.843 | 29.791 | 0.805 | 0.448 | 0.460 | 0.616 | 0.409 | 0.109 | 0.466 | 0.400 |
| Claude-3 | 3.269 | 12.922 | 0.021 | 0.035 | 0.044 | 0.280 | 0.001 | 0.976 | 0.075 | 0.838 |
| CommandR-plus | 29.694 | 30.461 | 0.831 | 0.252 | 0.247 | 0.376 | 0.239 | 0.086 | 0.337 | 0.540 |
| GPT-4 | 36.402 | 41.005 | 0.991 | 0.078 | 0.035 | 0.009 | 0.005 | 0.032 | 0.177 | 0.677 |
| Llama3-70B | 36.940 | 45.403 | 0.971 | 0.874 | 0.907 | 0.979 | 0.972 | 0.023 | 0.915 | 0.250 |
| IKUN | 31.212 | 43.813 | 0.936 | 0.871 | 0.920 | 0.991 | 0.940 | 0.006 | 0.905 | 0.254 |
| IKUN-C | 12.917 | 29.192 | 0.455 | 0.420 | 0.480 | 0.836 | 0.426 | 0.111 | 0.476 | 0.372 |
| Unbabel-Tower70B | 36.851 | 44.929 | 0.996 | **0.935** | **0.967** | **0.995** | 0.995 | 0.005 | **0.970** | 0.246 |
| AMI | 35.669 | 48.173 | 0.837 | 0.812 | 0.853 | 0.994 | 0.863 | 0.048 | 0.860 | 0.285 |
| CycleL | 3.474 | 14.663 | 0.000 | 0.078 | 0.127 | 0.435 | 0.000 | **0.000** | 0.091 | 0.482 |
| Dubformer | 26.980 | 28.818 | 0.732 | 0.070 | 0.034 | 0.045 | 0.004 | 0.255 | 0.139 | 0.723 |
| IOL_Research | 18.067 | 33.650 | 0.996 | 0.857 | 0.931 | 0.989 | 0.990 | 0.004 | 0.931 | 0.248 |
| ONLINE-A | 30.906 | 42.166 | 0.994 | 0.918 | 0.946 | **0.995** | 0.991 | 0.001 | 0.961 | 0.251 |
| ONLINE-B | 43.584 | 35.063 | 0.891 | 0.892 | 0.934 | 0.993 | 0.946 | 0.005 | 0.918 | 0.261 |
| ONLINE-G | 30.524 | 42.300 | **0.998** | 0.906 | 0.951 | 0.993 | **0.998** | 0.001 | 0.959 | **0.244** |
| ONLINE-W | | | | | | NA | | | | |
| TSU-HITs | 0.000 | 3.586 | 0.021 | 0.118 | 0.318 | 0.900 | 0.031 | 0.082 | 0.212 | 0.413 |
| TranssionMT | **43.597** | **53.077** | 0.881 | 0.897 | 0.934 | 0.993 | 0.945 | 0.005 | 0.920 | 0.264 |

Table 56: English→Icelandic, 1-shot JSON format

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 16.547 | 35.168 | 0.995 | 0.072 | 0.841 | 0.990 | 0.542 | 0.328 | 0.597 | **0.277** |
| Claude-3 | 3.943 | 38.065 | 0.993 | 0.111 | 0.869 | 0.994 | 0.561 | 0.311 | **0.628** | 0.280 |
| CommandR-plus | 7.728 | 35.127 | 0.993 | 0.084 | 0.816 | 0.980 | 0.541 | 0.322 | 0.599 | 0.286 |
| GPT-4 | 15.472 | 39.233 | 0.999 | 0.082 | 0.853 | 0.995 | 0.558 | 0.341 | 0.617 | 0.288 |
| Llama3-70B | **18.386** | 32.080 | 0.998 | 0.059 | 0.845 | 0.993 | 0.569 | 0.339 | 0.596 | 0.280 |
| IKUN | 1.519 | 28.192 | 0.996 | 0.039 | 0.796 | **0.996** | 0.463 | 0.411 | 0.556 | 0.292 |
| IKUN-C | 5.156 | 23.669 | 0.988 | 0.021 | 0.761 | **0.996** | 0.390 | 0.420 | 0.512 | 0.289 |
| Unbabel-Tower70B | 6.585 | 36.271 | 0.996 | 0.076 | 0.830 | 0.987 | 0.550 | 0.317 | 0.602 | 0.284 |
| CycleL | 0.013 | 2.344 | 0.406 | 0.004 | 0.257 | 0.869 | 0.022 | **0.065** | 0.224 | 0.371 |
| DLUT_GTCOM | 0.735 | 30.945 | 0.830 | 0.006 | 0.789 | 0.969 | 0.556 | 0.323 | 0.544 | 0.343 |
| IOL_Research | 16.514 | 39.294 | 0.998 | 0.104 | 0.847 | **0.996** | **0.590** | 0.304 | 0.623 | 0.279 |
| MSLC | 9.124 | 29.066 | 0.995 | 0.071 | 0.815 | 0.940 | 0.542 | 0.335 | 0.571 | 0.282 |
| NTTSU | 0.456 | 32.324 | 0.999 | 0.005 | 0.792 | 0.976 | 0.580 | 0.266 | 0.574 | 0.297 |
| ONLINE-A | 4.688 | **39.838** | **1.000** | **0.125** | 0.853 | 0.993 | 0.449 | 0.398 | 0.618 | 0.292 |
| ONLINE-B | 1.534 | 38.803 | 0.998 | 0.120 | 0.864 | 0.989 | 0.466 | 0.360 | 0.619 | 0.287 |
| ONLINE-G | 2.440 | 33.098 | 0.998 | 0.087 | 0.841 | 0.990 | 0.482 | 0.360 | 0.598 | 0.290 |
| ONLINE-W | 2.803 | 38.856 | 0.990 | 0.111 | **0.871** | 0.995 | 0.463 | 0.344 | 0.611 | 0.285 |
| Team-J | 0.573 | 28.582 | 0.999 | 0.007 | 0.788 | 0.988 | 0.550 | 0.294 | 0.566 | 0.307 |
| UvA-MT | 0.413 | 32.523 | 0.996 | 0.007 | 0.776 | 0.961 | 0.579 | 0.268 | 0.566 | 0.298 |

Table 57: Japanese→Chinese, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 2.588 | 30.164 | 0.996 | 0.016 | **1.000** | 0.807 | 0.474 | 0.504 | 0.555 | 0.241 |
| Claude-3 | 0.180 | 21.355 | 0.253 | 0.021 | 0.977 | 0.955 | 0.393 | 0.526 | 0.425 | 0.334 |
| CommandR-plus | 2.022 | 32.787 | 0.907 | 0.015 | 0.999 | 0.960 | 0.465 | 0.493 | 0.579 | 0.231 |
| GPT-4 | 7.246 | 42.392 | 0.996 | 0.020 | **1.000** | 1.000 | 0.520 | 0.447 | 0.640 | 0.205 |
| Llama3-70B | 1.227 | 27.591 | 0.994 | 0.012 | 0.999 | 0.871 | 0.504 | 0.470 | 0.578 | 0.228 |
| IKUN | 0.186 | 34.691 | 0.993 | 0.018 | **1.000** | 1.000 | 0.461 | 0.481 | 0.597 | 0.211 |
| IKUN-C | 0.217 | 14.227 | 0.802 | 0.010 | 0.979 | 0.998 | 0.236 | 0.589 | 0.459 | 0.257 |
| Unbabel-Tower70B | 2.887 | 37.396 | 0.991 | 0.011 | 0.999 | 0.994 | 0.509 | 0.463 | 0.623 | 0.210 |
| CycleL | 0.002 | 1.407 | 0.324 | 0.007 | 0.529 | 0.732 | 0.006 | **0.022** | 0.228 | 0.352 |
| DLUT_GTCOM | 0.205 | 26.707 | 0.722 | 0.017 | **1.000** | 1.000 | 0.471 | 0.468 | 0.552 | 0.268 |
| IOL_Research | **10.974** | **47.947** | 0.994 | **0.082** | **1.000** | 1.000 | 0.458 | 0.494 | **0.642** | **0.203** |
| MSLC | 4.476 | 32.966 | 0.999 | 0.015 | 0.994 | 0.772 | **0.530** | 0.448 | 0.575 | 0.239 |
| NTTSU | 0.026 | 29.887 | **1.000** | 0.011 | **1.000** | 1.000 | 0.490 | 0.393 | 0.611 | 0.208 |
| ONLINE-A | 0.108 | 42.229 | **1.000** | 0.011 | 0.999 | 0.999 | 0.490 | 0.446 | 0.636 | 0.205 |
| ONLINE-B | 0.600 | 40.376 | 0.999 | 0.010 | **1.000** | 1.000 | 0.519 | 0.439 | 0.636 | 0.205 |
| ONLINE-G | 0.182 | 26.305 | 0.996 | 0.012 | **1.000** | 1.000 | 0.315 | 0.514 | 0.567 | 0.216 |
| ONLINE-W | 1.180 | 44.101 | 0.995 | 0.022 | **1.000** | 1.000 | 0.493 | 0.453 | 0.635 | 0.206 |
| Team-J | 0.041 | 28.167 | 0.999 | 0.011 | **1.000** | 1.000 | 0.448 | 0.436 | 0.607 | 0.218 |
| UvA-MT | 0.057 | 26.219 | 0.996 | 0.016 | 0.999 | 0.878 | 0.242 | 0.613 | 0.530 | 0.246 |

Table 58: Japanese→Chinese, direct (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.328 | 26.334 | 0.889 | 0.009 | 0.927 | 0.934 | 0.441 | 0.553 | 0.538 | 0.259 |
| Claude-3 | 0.053 | 11.442 | 0.321 | 0.013 | 0.652 | 0.673 | 0.244 | 0.515 | 0.320 | 0.409 |
| CommandR-plus | 0.038 | 22.249 | 0.553 | 0.011 | 0.716 | 0.788 | 0.272 | 0.480 | 0.403 | 0.345 |
| GPT-4 | 0.914 | 41.297 | 0.736 | 0.049 | 0.949 | 0.971 | 0.241 | 0.343 | 0.528 | 0.234 |
| Llama3-70B | 0.016 | 8.771 | 0.481 | 0.020 | 0.491 | 0.498 | 0.010 | 0.732 | 0.261 | 0.466 |
| IKUN | 0.489 | 42.055 | 0.985 | 0.013 | **1.000** | 1.000 | 0.519 | 0.056 | 0.642 | 0.151 |
| IKUN-C | **1.374** | 26.271 | 0.974 | 0.013 | 0.999 | 0.999 | **0.573** | 0.207 | 0.595 | 0.175 |
| Unbabel-Tower70B | 0.172 | 35.096 | 0.983 | 0.012 | **1.000** | 0.998 | 0.531 | 0.171 | 0.630 | 0.169 |
| CycleL | 0.009 | 0.358 | 0.471 | 0.006 | 0.048 | 0.284 | 0.000 | **0.021** | 0.116 | 0.462 |
| DLUT_GTCOM | 0.129 | 21.671 | 0.976 | 0.020 | 0.952 | 0.968 | 0.528 | 0.095 | 0.550 | 0.183 |
| IOL_Research | 1.087 | **51.423** | 0.985 | **0.077** | 0.987 | 0.988 | 0.519 | 0.065 | **0.656** | 0.147 |
| MSLC | 0.081 | 3.515 | **0.999** | 0.015 | 0.013 | 0.000 | 0.015 | 0.625 | 0.149 | 0.521 |
| NTTSU | 0.081 | 6.629 | 0.996 | 0.015 | 0.300 | 0.297 | 0.099 | 0.559 | 0.246 | 0.429 |
| ONLINE-A | 0.119 | 44.345 | **0.999** | 0.012 | **1.000** | 1.000 | 0.513 | 0.051 | 0.644 | 0.149 |
| ONLINE-B | 1.139 | 47.534 | 0.995 | 0.020 | **1.000** | 1.000 | 0.541 | 0.034 | 0.651 | **0.146** |
| ONLINE-G | 0.384 | 33.157 | **0.999** | 0.016 | **1.000** | 1.000 | 0.519 | 0.059 | 0.628 | 0.149 |
| ONLINE-W | 0.501 | 38.632 | 0.993 | 0.022 | **1.000** | 1.000 | 0.510 | 0.051 | 0.636 | 0.152 |
| Team-J | 0.054 | 19.455 | 0.998 | 0.010 | 0.919 | 0.924 | 0.435 | 0.081 | 0.541 | 0.191 |
| UvA-MT | 0.026 | 10.203 | 0.976 | 0.011 | 0.330 | 0.315 | 0.334 | 0.222 | 0.316 | 0.374 |

Table 59: Japanese→Chinese, direct (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 1.611 | 26.334 | 0.991 | 0.015 | **1.000** | 0.701 | 0.827 | 0.159 | 0.596 | 0.208 |
| Claude-3 | 0.611 | 37.944 | 0.725 | 0.015 | 0.979 | 0.998 | 0.474 | 0.441 | 0.561 | 0.246 |
| CommandR-plus | 1.998 | 33.800 | 0.873 | 0.015 | 0.993 | 0.953 | 0.546 | 0.426 | 0.606 | 0.228 |
| GPT-4 | 7.295 | 44.604 | 0.999 | 0.018 | **1.000** | **1.000** | 0.732 | 0.248 | 0.681 | 0.176 |
| Llama3-70B | 0.713 | 32.894 | 0.996 | 0.017 | 0.998 | 0.868 | 0.887 | 0.095 | 0.660 | 0.174 |
| IKUN | 0.341 | 40.606 | 0.776 | 0.011 | **1.000** | **1.000** | 0.606 | 0.275 | 0.628 | 0.214 |
| IKUN-C | 0.165 | 13.209 | 0.690 | 0.004 | 0.965 | 0.998 | 0.351 | 0.395 | 0.446 | 0.249 |
| Unbabel-Tower70B | 3.371 | 40.782 | 0.995 | 0.013 | **1.000** | **1.000** | 0.802 | 0.187 | 0.689 | 0.168 |
| CycleL | 0.001 | 0.747 | 0.211 | 0.006 | 0.436 | 0.596 | 0.000 | **0.010** | 0.178 | 0.396 |
| DLUT_GTCOM | 0.092 | 26.034 | 0.911 | 0.020 | **1.000** | **1.000** | 0.471 | 0.483 | 0.591 | 0.222 |
| IOL_Research | **9.717** | **54.953** | 0.996 | **0.044** | **1.000** | **1.000** | 0.711 | 0.252 | 0.685 | 0.173 |
| MSLC | 4.401 | 39.207 | 0.998 | 0.016 | **1.000** | 0.998 | **0.902** | 0.087 | **0.704** | **0.154** |
| NTTSU | 0.013 | 27.408 | **1.000** | 0.012 | **1.000** | **1.000** | 0.605 | 0.307 | 0.652 | 0.185 |
| ONLINE-A | 0.045 | 41.376 | **1.000** | 0.011 | **1.000** | **1.000** | 0.786 | 0.196 | 0.685 | 0.169 |
| ONLINE-B | 0.486 | 40.506 | 0.999 | 0.010 | **1.000** | **1.000** | 0.766 | 0.185 | 0.683 | 0.168 |
| ONLINE-G | 0.142 | 28.402 | 0.998 | 0.012 | 0.998 | **1.000** | 0.294 | 0.558 | 0.575 | 0.225 |
| ONLINE-W | 0.801 | 47.180 | 0.998 | 0.021 | **1.000** | **1.000** | 0.693 | 0.187 | 0.675 | 0.168 |
| Team-J | 0.021 | 22.331 | 0.999 | 0.012 | **1.000** | **1.000** | 0.482 | 0.424 | 0.570 | 0.202 |
| UvA-MT | 0.023 | 27.965 | 0.999 | 0.016 | **1.000** | **1.000** | 0.198 | 0.690 | 0.568 | 0.239 |

Table 60: Japanese→Chinese, 0-shot (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.176 | 20.123 | 0.887 | 0.010 | 0.944 | 0.941 | 0.005 | 0.966 | 0.447 | 0.314 |
| Claude-3 | 0.602 | 49.704 | 0.818 | 0.016 | 0.989 | 0.989 | 0.289 | 0.093 | 0.561 | 0.183 |
| CommandR-plus | 0.042 | 32.258 | 0.659 | 0.011 | 0.780 | 0.789 | 0.130 | 0.388 | 0.423 | 0.308 |
| GPT-4 | 0.751 | 51.736 | 0.996 | 0.043 | **1.000** | **1.000** | 0.104 | 0.126 | 0.597 | 0.156 |
| Llama3-70B | 0.159 | 43.130 | 0.950 | 0.020 | 0.990 | 0.998 | 0.035 | 0.258 | 0.550 | 0.187 |
| IKUN | 0.636 | 43.475 | 0.878 | 0.011 | **1.000** | **1.000** | 0.140 | 0.225 | 0.575 | 0.192 |
| IKUN-C | 0.066 | 24.453 | 0.971 | 0.011 | 0.996 | 0.999 | 0.200 | 0.412 | 0.541 | 0.205 |
| Unbabel-Tower70B | 0.072 | 37.680 | 0.990 | 0.010 | 0.999 | 0.999 | 0.132 | 0.313 | 0.579 | 0.188 |
| CycleL | 0.002 | 0.317 | 0.469 | 0.004 | 0.000 | 0.009 | 0.000 | **0.002** | 0.069 | 0.505 |
| DLUT_GTCOM | 0.059 | 19.556 | 0.995 | 0.010 | 0.978 | 0.999 | 0.093 | 0.102 | 0.471 | 0.160 |
| IOL_Research | **1.958** | 52.599 | 0.977 | **0.066** | 0.988 | 0.984 | 0.124 | 0.149 | 0.596 | 0.163 |
| MSLC | 0.039 | 3.199 | **0.999** | 0.010 | 0.006 | 0.000 | 0.023 | 0.532 | 0.148 | 0.508 |
| NTTSU | 0.040 | 4.784 | 0.996 | 0.009 | 0.148 | 0.140 | 0.021 | 0.472 | 0.188 | 0.458 |
| ONLINE-A | 0.120 | 50.307 | **0.999** | 0.012 | **1.000** | **1.000** | 0.103 | 0.084 | 0.588 | 0.153 |
| ONLINE-B | 0.796 | **56.173** | 0.995 | 0.011 | **1.000** | **1.000** | 0.098 | 0.081 | 0.588 | 0.154 |
| ONLINE-G | 0.229 | 35.943 | **0.999** | 0.015 | **1.000** | **1.000** | **0.353** | 0.111 | **0.609** | 0.158 |
| ONLINE-W | 0.246 | 43.317 | 0.988 | 0.015 | **1.000** | **1.000** | 0.186 | 0.051 | 0.599 | 0.158 |
| Team-J | 0.029 | 25.499 | 0.998 | 0.012 | 0.998 | 0.995 | 0.088 | 0.064 | 0.564 | **0.152** |
| UvA-MT | 0.007 | 3.772 | 0.994 | 0.001 | 0.094 | 0.104 | 0.026 | 0.556 | 0.178 | 0.486 |

Table 61: Japanese→Chinese, 0-shot (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 1.355 | 15.611 | **0.628** | 0.334 | 0.854 | 0.908 | 0.721 | 0.219 | 0.549 | 0.213 |
| Claude-3 | 19.917 | 28.083 | 0.461 | 0.509 | 0.999 | 0.541 | 0.588 | 0.409 | 0.536 | 0.271 |
| CommandR-plus | 21.863 | 31.575 | 0.444 | 0.513 | **1.000** | 0.635 | 0.515 | 0.439 | 0.621 | 0.264 |
| GPT-4 | **67.128** | **67.621** | 0.498 | 0.513 | **1.000** | 0.628 | 0.812 | 0.184 | 0.709 | 0.221 |
| Llama3-70B | 28.043 | 39.157 | 0.498 | 0.513 | 0.999 | 0.950 | **0.918** | 0.081 | **0.770** | **0.160** |
| IKUN | 2.085 | 20.834 | 0.268 | 0.424 | 0.950 | 0.913 | 0.709 | 0.228 | 0.563 | 0.240 |
| IKUN-C | 19.770 | 25.817 | 0.376 | 0.424 | 0.935 | 0.792 | 0.600 | 0.291 | 0.519 | 0.252 |
| Unbabel-Tower70B | 45.809 | 49.308 | 0.490 | 0.507 | **1.000** | 0.796 | 0.821 | 0.175 | 0.731 | 0.198 |
| CycleL | 0.495 | 0.893 | 0.022 | 0.293 | 0.377 | 0.415 | 0.000 | **0.010** | 0.158 | 0.417 |
| DLUT_GTCOM | 26.249 | 21.270 | 0.446 | 0.512 | **1.000** | **1.000** | 0.552 | 0.447 | 0.641 | 0.213 |
| IOL_Research | 26.634 | 42.331 | 0.494 | 0.508 | **1.000** | 0.953 | 0.786 | 0.209 | 0.749 | 0.179 |
| MSLC | 17.597 | 41.233 | 0.497 | 0.513 | **1.000** | 0.498 | 0.882 | 0.067 | 0.692 | 0.223 |
| NTTSU | 53.273 | 49.379 | 0.499 | 0.508 | **1.000** | 0.499 | 0.747 | 0.240 | 0.670 | 0.248 |
| ONLINE-A | 15.584 | 36.244 | 0.499 | 0.509 | **1.000** | 0.973 | 0.873 | 0.126 | 0.728 | 0.163 |
| ONLINE-B | 29.018 | 33.631 | 0.498 | 0.508 | **1.000** | 0.523 | 0.805 | 0.168 | 0.692 | 0.234 |
| ONLINE-G | 47.668 | 43.705 | 0.497 | 0.507 | **1.000** | 0.499 | 0.546 | 0.426 | 0.650 | 0.283 |
| ONLINE-W | 51.642 | 50.577 | 0.494 | **0.514** | **1.000** | 0.519 | 0.761 | 0.169 | 0.685 | 0.235 |
| Team-J | 39.907 | 30.652 | 0.498 | 0.509 | **1.000** | 0.529 | 0.594 | 0.395 | 0.607 | 0.266 |
| UvA-MT | 3.718 | 28.426 | 0.497 | 0.512 | **1.000** | 0.499 | 0.517 | 0.460 | 0.568 | 0.279 |

Table 62: Japanese→Chinese, 1-shot (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 31.765 | 30.959 | 0.466 | 0.509 | 0.974 | 0.491 | 0.038 | 0.639 | 0.521 | 0.314 |
| Claude-3 | 34.976 | 38.640 | 0.491 | **0.514** | 0.998 | 0.502 | 0.187 | 0.244 | 0.598 | 0.249 |
| CommandR-plus | 24.284 | 29.056 | 0.401 | 0.507 | 0.950 | 0.728 | 0.132 | 0.311 | 0.548 | 0.247 |
| GPT-4 | 30.521 | 41.301 | 0.498 | 0.508 | **1.000** | 0.753 | 0.114 | 0.559 | 0.625 | 0.259 |
| Llama3-70B | 38.715 | **47.363** | 0.493 | 0.508 | **1.000** | **1.000** | 0.064 | 0.438 | 0.652 | 0.205 |
| IKUN | 37.847 | 45.313 | 0.367 | 0.507 | **1.000** | **1.000** | 0.113 | 0.632 | 0.640 | 0.252 |
| IKUN-C | 15.065 | 27.858 | 0.482 | 0.482 | 0.983 | 0.860 | 0.168 | 0.616 | 0.567 | 0.258 |
| Unbabel-Tower70B | 37.691 | 44.139 | 0.493 | 0.504 | **1.000** | 0.525 | 0.118 | 0.574 | 0.591 | 0.294 |
| CycleL | 0.301 | 0.443 | 0.048 | 0.252 | 0.138 | 0.332 | 0.000 | **0.053** | 0.110 | 0.470 |
| DLUT_GTCOM | 34.167 | 19.519 | 0.543 | 0.503 | 0.994 | 0.504 | 0.081 | 0.094 | 0.454 | 0.222 |
| IOL_Research | 9.940 | 23.674 | 0.531 | 0.504 | 0.999 | 0.925 | 0.131 | 0.490 | 0.522 | 0.219 |
| MSLC | 33.728 | 30.354 | 0.569 | 0.504 | 0.501 | 0.000 | 0.016 | 0.574 | 0.370 | 0.430 |
| NTTSU | 36.670 | 29.558 | 0.506 | 0.504 | 0.503 | 0.002 | 0.016 | 0.787 | 0.361 | 0.469 |
| ONLINE-A | 3.209 | 26.445 | 0.531 | 0.509 | **1.000** | **1.000** | 0.129 | 0.458 | 0.524 | **0.202** |
| ONLINE-B | 38.173 | 32.152 | 0.496 | 0.507 | **1.000** | 0.501 | 0.086 | 0.409 | 0.528 | 0.272 |
| ONLINE-G | 39.419 | 32.800 | 0.498 | 0.508 | **1.000** | 0.988 | **0.337** | 0.481 | **0.689** | 0.212 |
| ONLINE-W | **40.948** | 42.011 | 0.556 | 0.508 | **1.000** | 0.586 | 0.159 | 0.224 | 0.616 | 0.231 |
| Team-J | 29.786 | 24.750 | 0.499 | 0.506 | 0.998 | 0.633 | 0.103 | 0.284 | 0.539 | 0.236 |
| UvA-MT | 7.284 | 14.144 | **0.627** | 0.501 | 0.512 | 0.294 | 0.000 | 0.770 | 0.276 | 0.407 |

Table 63: Japanese→Chinese, 1-shot (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.264 | 10.033 | 0.573 | 0.016 | 0.741 | 0.752 | 0.160 | 0.453 | 0.329 | 0.340 |
| Claude-3 | 5.179 | 19.101 | 0.239 | 0.005 | 0.275 | 0.082 | 0.213 | 0.182 | 0.121 | 0.512 |
| CommandR-plus | 12.402 | 23.916 | 0.849 | 0.018 | 0.804 | 0.197 | 0.709 | 0.135 | 0.381 | 0.324 |
| GPT-4 | **67.617** | **71.625** | **0.998** | 0.037 | 0.985 | 0.592 | 0.624 | 0.259 | 0.518 | 0.236 |
| Llama3-70B | 15.287 | 30.940 | 0.993 | 0.021 | 0.853 | 0.973 | 0.235 | 0.248 | 0.452 | 0.202 |
| IKUN | 0.000 | 5.032 | 0.067 | 0.001 | 0.132 | 0.583 | 0.064 | 0.141 | 0.121 | 0.483 |
| IKUN-C | 0.042 | 8.097 | 0.152 | 0.002 | 0.207 | 0.476 | 0.137 | 0.143 | 0.140 | 0.475 |
| Unbabel-Tower70B | 22.687 | 34.978 | 0.987 | 0.021 | 0.917 | 0.797 | 0.377 | 0.214 | 0.459 | 0.214 |
| CycleL | 0.094 | 0.891 | 0.002 | 0.001 | 0.073 | 0.291 | 0.000 | 0.031 | 0.053 | 0.523 |
| DLUT_GTCOM | 12.520 | 13.585 | 0.939 | 0.017 | 0.897 | 0.946 | 0.491 | 0.290 | 0.547 | 0.221 |
| IOL_Research | 15.099 | 31.032 | 0.996 | **0.115** | **0.998** | 0.991 | 0.558 | 0.319 | **0.643** | **0.175** |
| MSLC | 14.225 | 38.986 | 0.422 | 0.011 | 0.450 | 0.005 | 0.219 | 0.244 | 0.160 | 0.479 |
| NTTSU | 25.209 | 37.241 | 0.001 | 0.000 | 0.200 | **0.999** | 0.000 | 0.058 | 0.171 | 0.408 |
| ONLINE-A | 13.913 | 33.281 | 0.983 | 0.011 | 0.880 | 0.001 | 0.881 | 0.070 | 0.396 | 0.314 |
| ONLINE-B | 24.786 | 29.969 | 0.590 | 0.005 | 0.565 | 0.058 | 0.213 | 0.464 | 0.205 | 0.464 |
| ONLINE-G | 35.089 | 40.914 | 0.978 | 0.015 | 0.955 | 0.000 | **0.900** | 0.038 | 0.409 | 0.299 |
| ONLINE-W | 43.162 | 45.767 | 0.749 | 0.016 | 0.755 | 0.033 | 0.672 | 0.048 | 0.322 | 0.357 |
| Team-J | 27.221 | 26.221 | 0.000 | 0.000 | 0.056 | 0.141 | 0.000 | **0.004** | 0.028 | 0.544 |
| UvA-MT | 32.440 | 45.552 | 0.554 | 0.007 | 0.552 | 0.029 | 0.515 | 0.012 | 0.239 | 0.410 |

Table 64: Japanese→Chinese, 0-shot JSON format (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 20.086 | 27.893 | 0.934 | 0.035 | 0.441 | 0.432 | 0.393 | 0.592 | 0.372 | 0.396 |
| Claude-3 | 22.020 | 24.731 | 0.851 | 0.132 | 0.607 | 0.586 | 0.570 | 0.360 | 0.485 | 0.321 |
| CommandR-plus | 11.205 | 20.551 | 0.670 | 0.066 | 0.683 | 0.666 | 0.589 | 0.181 | 0.458 | 0.304 |
| GPT-4 | 19.191 | 30.499 | **1.000** | 0.089 | 0.465 | 0.439 | 0.453 | 0.545 | 0.421 | 0.369 |
| Llama3-70B | 21.816 | **35.723** | 0.988 | 0.132 | 0.946 | 0.941 | 0.945 | 0.049 | 0.696 | 0.149 |
| IKUN | 21.067 | 35.609 | 0.973 | 0.080 | 0.987 | 0.985 | 0.966 | 0.010 | 0.677 | 0.141 |
| IKUN-C | 9.191 | 24.181 | 0.780 | 0.032 | 0.763 | 0.896 | 0.715 | 0.098 | 0.511 | 0.233 |
| Unbabel-Tower70B | 22.497 | 34.133 | 0.999 | 0.073 | 0.712 | 0.704 | 0.703 | 0.294 | 0.546 | 0.261 |
| CycleL | 0.003 | 0.308 | 0.006 | 0.001 | 0.059 | 0.450 | 0.004 | 0.022 | 0.074 | 0.502 |
| DLUT_GTCOM | 19.225 | 14.731 | 0.908 | 0.017 | 0.253 | 0.245 | 0.264 | 0.393 | 0.260 | 0.431 |
| IOL_Research | 5.577 | 12.479 | 0.892 | 0.092 | 0.760 | 0.799 | 0.743 | 0.106 | 0.571 | 0.224 |
| MSLC | **30.548** | 29.612 | 0.756 | 0.015 | 0.031 | 0.081 | 0.065 | 0.775 | 0.137 | 0.562 |
| NTTSU | 29.066 | 30.664 | 0.933 | 0.020 | 0.032 | 0.009 | 0.033 | 0.951 | 0.149 | 0.572 |
| ONLINE-A | 0.595 | 12.888 | 0.902 | 0.113 | 0.940 | 0.886 | 0.856 | 0.013 | 0.640 | 0.168 |
| ONLINE-B | 22.201 | 18.574 | 0.963 | 0.125 | **0.999** | 0.976 | **0.988** | 0.006 | 0.709 | 0.135 |
| ONLINE-G | 21.850 | 22.009 | 0.998 | **0.152** | 0.994 | **0.993** | 0.985 | 0.010 | **0.725** | **0.125** |
| ONLINE-W | 20.095 | 27.330 | 0.389 | 0.016 | 0.371 | 0.796 | 0.242 | 0.306 | 0.278 | 0.395 |
| Team-J | 14.330 | 17.465 | 0.020 | 0.021 | 0.136 | 0.114 | 0.177 | 0.721 | 0.074 | 0.638 |
| UvA-MT | 2.731 | 13.997 | 0.294 | 0.011 | 0.022 | 0.011 | 0.043 | 0.498 | 0.057 | 0.602 |

Table 65: Japanese→Chinese, 0-shot JSON format (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 0.393 | 11.802 | 0.392 | 0.015 | 0.694 | 0.674 | 0.100 | 0.431 | 0.272 | 0.380 |
| Claude-3 | 2.021 | 18.559 | 0.402 | 0.005 | 0.483 | 0.027 | 0.392 | 0.243 | 0.188 | 0.475 |
| CommandR-plus | 8.531 | 21.922 | 0.902 | 0.025 | 0.873 | 0.257 | 0.757 | 0.162 | 0.414 | 0.301 |
| GPT-4 | **66.340** | **72.970** | 0.993 | 0.025 | 0.978 | 0.419 | 0.586 | 0.319 | 0.475 | 0.272 |
| Llama3-70B | 10.351 | 28.423 | 0.983 | 0.025 | 0.824 | 0.985 | 0.201 | 0.267 | 0.439 | 0.208 |
| IKUN | 0.000 | 2.210 | 0.015 | 0.000 | 0.083 | 0.578 | 0.000 | 0.176 | 0.097 | 0.504 |
| IKUN-C | 0.000 | 2.210 | 0.015 | 0.000 | 0.083 | 0.578 | 0.000 | 0.176 | 0.097 | 0.504 |
| Unbabel-Tower70B | 11.566 | 29.219 | 0.973 | 0.022 | 0.848 | 0.953 | 0.262 | 0.257 | 0.459 | 0.209 |
| CycleL | 0.024 | 0.693 | 0.002 | 0.000 | 0.049 | 0.346 | 0.002 | 0.037 | 0.057 | 0.522 |
| DLUT_GTCOM | 8.747 | 11.083 | 0.995 | 0.022 | 0.858 | 0.926 | 0.456 | 0.306 | 0.533 | 0.225 |
| IOL_Research | 10.299 | 28.733 | 0.995 | **0.078** | **0.993** | 0.988 | 0.578 | 0.314 | **0.633** | **0.181** |
| MSLC | 12.970 | 39.098 | 0.485 | 0.012 | 0.505 | 0.012 | 0.157 | 0.152 | 0.169 | 0.448 |
| NTTSU | 12.223 | 29.135 | 0.002 | 0.002 | 0.228 | **0.998** | 0.000 | 0.061 | 0.176 | 0.404 |
| ONLINE-A | 13.046 | 32.576 | 0.980 | 0.022 | 0.792 | 0.002 | 0.853 | 0.110 | 0.382 | 0.331 |
| ONLINE-B | 22.458 | 29.216 | 0.738 | 0.012 | 0.718 | 0.051 | 0.282 | 0.451 | 0.259 | 0.419 |
| ONLINE-G | 29.636 | 38.288 | 0.975 | 0.022 | 0.944 | 0.000 | 0.841 | 0.108 | 0.401 | 0.310 |
| ONLINE-W | 38.867 | 43.492 | 0.757 | 0.020 | 0.728 | 0.034 | 0.703 | 0.015 | 0.324 | 0.354 |
| Team-J | 22.863 | 22.714 | 0.002 | 0.000 | 0.159 | 0.250 | 0.000 | **0.000** | 0.059 | 0.513 |
| UvA-MT | 34.476 | 54.459 | **0.998** | 0.022 | 0.985 | 0.000 | **0.961** | 0.007 | 0.428 | 0.286 |

Table 66: Japanese→Chinese, 1-shot JSON format (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 14.704 | 26.618 | 0.885 | 0.032 | 0.479 | 0.482 | 0.445 | 0.553 | 0.384 | 0.385 |
| Claude-3 | 12.153 | 20.101 | 0.658 | 0.042 | 0.215 | 0.196 | 0.130 | 0.756 | 0.198 | 0.540 |
| CommandR-plus | 8.781 | 19.002 | 0.689 | 0.046 | 0.663 | 0.587 | 0.570 | 0.208 | 0.430 | 0.320 |
| GPT-4 | 15.159 | 28.058 | 0.993 | 0.024 | 0.044 | 0.029 | 0.032 | 0.961 | 0.167 | 0.559 |
| Llama3-70B | 21.514 | **34.361** | 0.927 | 0.125 | 0.878 | 0.875 | 0.841 | 0.130 | 0.626 | 0.191 |
| IKUN | 14.299 | 32.547 | 0.985 | 0.078 | 0.976 | 0.980 | 0.963 | 0.022 | 0.680 | 0.143 |
| IKUN-C | 7.250 | 23.292 | 0.809 | 0.034 | 0.824 | 0.919 | 0.782 | 0.066 | 0.544 | 0.212 |
| Unbabel-Tower70B | 16.824 | 31.039 | 0.995 | 0.056 | 0.687 | 0.672 | 0.670 | 0.318 | 0.528 | 0.276 |
| CycleL | 0.000 | 0.325 | 0.007 | 0.002 | 0.078 | 0.467 | 0.000 | 0.007 | 0.079 | 0.496 |
| DLUT_GTCOM | 13.570 | 12.482 | 0.978 | 0.007 | 0.061 | 0.071 | 0.068 | 0.320 | 0.173 | 0.462 |
| IOL_Research | 4.593 | 12.908 | 0.958 | 0.088 | 0.861 | 0.848 | 0.863 | 0.078 | 0.630 | 0.189 |
| MSLC | **24.824** | 25.529 | 0.998 | 0.020 | 0.017 | 0.000 | 0.042 | 0.946 | 0.156 | 0.565 |
| NTTSU | 23.175 | 27.433 | 0.861 | 0.012 | 0.020 | 0.007 | 0.046 | 0.912 | 0.137 | 0.582 |
| ONLINE-A | 0.221 | 10.235 | **1.000** | 0.095 | 0.910 | 0.988 | 0.804 | 0.024 | 0.647 | 0.147 |
| ONLINE-B | 15.326 | 16.097 | 0.998 | **0.127** | **1.000** | 0.980 | 0.988 | 0.010 | 0.721 | 0.130 |
| ONLINE-G | 15.138 | 17.887 | 0.995 | **0.127** | **1.000** | **1.000** | **0.995** | **0.005** | **0.725** | **0.126** |
| ONLINE-W | 11.378 | 21.633 | 0.174 | 0.010 | 0.291 | 0.949 | 0.132 | 0.257 | 0.236 | 0.405 |
| Team-J | 7.562 | 14.115 | 0.012 | 0.022 | 0.191 | 0.171 | 0.306 | 0.660 | 0.108 | 0.614 |
| UvA-MT | 1.650 | 12.212 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.479 | 0.000 | 0.669 |

Table 67: Japanese→Chinese, 1-shot JSON format (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 51.796 | 71.017 | **1.000** | 0.912 | 0.925 | **1.000** | 0.854 | 0.124 | 0.936 | 0.261 |
| Claude-3 | **59.164** | **76.525** | 1.000 | **0.945** | **0.952** | **1.000** | 0.871 | 0.113 | **0.957** | 0.266 |
| CommandR-plus | 52.291 | 71.954 | 0.998 | 0.930 | 0.950 | 0.999 | 0.837 | 0.149 | 0.941 | 0.263 |
| GPT-4 | 50.830 | 71.774 | **1.000** | 0.936 | 0.946 | **1.000** | 0.880 | 0.106 | 0.949 | 0.255 |
| Llama3-70B | 42.691 | 65.406 | **1.000** | 0.906 | 0.934 | 0.999 | **0.891** | 0.095 | 0.931 | **0.249** |
| IKUN | 44.345 | 65.724 | 0.999 | 0.919 | 0.934 | **1.000** | 0.842 | 0.146 | 0.928 | 0.258 |
| IKUN-C | 43.714 | 65.549 | **1.000** | 0.900 | 0.929 | **1.000** | 0.852 | 0.131 | 0.922 | 0.257 |
| Unbabel-Tower70B | 50.091 | 71.296 | 0.991 | 0.923 | 0.940 | **1.000** | 0.831 | 0.155 | 0.937 | 0.268 |
| BJFU-LPT | 23.070 | 42.742 | 0.999 | 0.673 | 0.780 | 0.965 | 0.483 | 0.280 | 0.729 | 0.289 |
| CUNI-Transformer | 51.200 | 70.250 | **1.000** | 0.922 | 0.947 | 0.999 | 0.848 | 0.143 | 0.940 | 0.263 |
| CycleL | 0.110 | 0.686 | 0.000 | 0.050 | 0.004 | 0.002 | 0.007 | **0.000** | 0.010 | 0.567 |
| IOL_Research | 54.964 | 73.144 | 0.984 | 0.925 | 0.941 | **1.000** | 0.856 | 0.125 | 0.943 | 0.267 |
| ONLINE-A | 49.693 | 69.758 | 0.999 | 0.907 | 0.942 | 0.999 | 0.808 | 0.163 | 0.924 | 0.265 |
| ONLINE-B | 47.317 | 68.256 | 0.998 | 0.897 | 0.924 | **1.000** | 0.792 | 0.180 | 0.915 | 0.268 |
| ONLINE-G | 43.649 | 65.989 | 0.999 | 0.906 | 0.933 | **1.000** | 0.769 | 0.197 | 0.910 | 0.268 |
| ONLINE-W | 51.432 | 69.965 | **1.000** | 0.920 | 0.931 | **1.000** | 0.787 | 0.181 | 0.924 | 0.269 |
| TranssionMT | 47.952 | 68.873 | 0.998 | 0.902 | 0.927 | **1.000** | 0.798 | 0.173 | 0.918 | 0.267 |

Table 68: Czech→Ukrainian, clean

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 46.384 | 70.808 | **0.999** | **1.000** | **1.000** | **1.000** | 0.860 | 0.038 | **0.976** | 0.008 |
| Claude-3 | 12.117 | 45.669 | 0.504 | 0.953 | 0.973 | 0.946 | 0.412 | 0.515 | 0.733 | 0.164 |
| CommandR-plus | 44.810 | 69.489 | 0.983 | 0.999 | 0.995 | **1.000** | 0.870 | 0.047 | 0.973 | 0.013 |
| GPT-4 | **46.744** | 69.115 | 0.966 | 0.995 | 0.998 | **1.000** | 0.836 | 0.093 | 0.963 | 0.022 |
| Llama3-70B | 40.918 | 67.909 | 0.991 | **1.000** | **1.000** | **1.000** | 0.878 | 0.037 | 0.974 | 0.009 |
| IKUN | 35.111 | 65.110 | 0.987 | **1.000** | **1.000** | 0.999 | 0.509 | 0.264 | 0.918 | 0.042 |
| IKUN-C | 28.412 | 60.127 | **0.999** | **1.000** | **1.000** | **1.000** | **0.894** | 0.032 | 0.952 | 0.007 |
| Unbabel-Tower70B | 36.354 | 68.928 | 0.993 | **1.000** | **1.000** | **1.000** | 0.797 | 0.048 | 0.965 | 0.010 |
| BJFU-LPT | 33.115 | 57.198 | 0.995 | 0.996 | **1.000** | **1.000** | 0.120 | 0.570 | 0.852 | 0.085 |
| CUNI-Transformer | 33.315 | 61.648 | **0.999** | **1.000** | **1.000** | **1.000** | 0.831 | 0.028 | 0.949 | **0.006** |
| CycleL | 0.023 | 1.498 | 0.000 | 0.295 | 0.038 | 0.002 | 0.000 | **0.000** | 0.048 | 0.524 |
| IOL_Research | 36.384 | 64.029 | 0.983 | 0.995 | 0.999 | **1.000** | 0.722 | 0.175 | 0.938 | 0.031 |
| ONLINE-A | 37.042 | 66.823 | 0.998 | **1.000** | **1.000** | **1.000** | 0.371 | 0.339 | 0.903 | 0.052 |
| ONLINE-B | 32.455 | 61.727 | **0.999** | **1.000** | **1.000** | **1.000** | 0.536 | 0.174 | 0.894 | 0.028 |
| ONLINE-G | 32.939 | 59.794 | **0.999** | **1.000** | **1.000** | **1.000** | 0.454 | 0.208 | 0.887 | 0.033 |
| ONLINE-W | 37.098 | 62.980 | 0.974 | **1.000** | **1.000** | **1.000** | 0.663 | 0.115 | 0.924 | 0.023 |
| TranssionMT | 43.336 | **73.713** | 0.998 | **1.000** | **1.000** | **1.000** | 0.471 | 0.302 | 0.924 | 0.046 |

Table 69: Czech→Ukrainian, direct (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | **45.336** | 69.080 | 0.995 | 0.998 | **1.000** | **1.000** | 0.297 | 0.398 | 0.889 | 0.060 |
| Claude-3 | 1.857 | 8.221 | 0.064 | 0.499 | 0.482 | 0.482 | **0.416** | 0.525 | 0.293 | 0.430 |
| CommandR-plus | 11.887 | 19.738 | 0.472 | 0.722 | 0.563 | 0.520 | 0.370 | 0.607 | 0.429 | 0.336 |
| GPT-4 | 31.663 | 57.593 | 0.718 | 0.920 | 0.953 | 0.960 | 0.118 | 0.439 | 0.748 | 0.129 |
| Llama3-70B | 3.090 | 13.892 | 0.267 | 0.424 | 0.355 | 0.360 | 0.034 | 0.835 | 0.265 | 0.493 |
| IKUN | 32.872 | 65.514 | 0.979 | **1.000** | **1.000** | **1.000** | 0.383 | 0.155 | 0.896 | 0.028 |
| IKUN-C | 26.102 | 55.453 | 0.995 | **1.000** | **1.000** | **1.000** | 0.404 | 0.251 | 0.857 | 0.039 |
| Unbabel-Tower70B | 38.610 | 72.475 | 0.991 | **1.000** | **1.000** | **1.000** | 0.365 | 0.196 | 0.904 | 0.031 |
| BJFU-LPT | 9.017 | 16.702 | 0.993 | 0.821 | 0.767 | 0.526 | 0.034 | 0.438 | 0.473 | 0.193 |
| CUNI-Transformer | 1.213 | 7.112 | **0.999** | 0.693 | 0.124 | 0.000 | 0.009 | 0.496 | 0.261 | 0.388 |
| CycleL | 0.032 | 1.317 | 0.000 | 0.157 | 0.001 | 0.000 | 0.000 | **0.147** | 0.023 | 0.571 |
| IOL_Research | 40.429 | 66.889 | 0.987 | 0.990 | **1.000** | **1.000** | 0.250 | 0.280 | 0.870 | 0.045 |
| ONLINE-A | 43.726 | 73.914 | 0.996 | **1.000** | **1.000** | **1.000** | 0.279 | 0.360 | 0.896 | 0.055 |
| ONLINE-B | 38.394 | 69.998 | 0.998 | **1.000** | **1.000** | **1.000** | 0.362 | 0.171 | 0.901 | 0.027 |
| ONLINE-G | 35.910 | 65.340 | 0.998 | **1.000** | **1.000** | **1.000** | 0.346 | 0.319 | 0.878 | 0.048 |
| ONLINE-W | 42.766 | 65.290 | 0.942 | **1.000** | **1.000** | **1.000** | 0.355 | 0.269 | 0.890 | 0.049 |
| TranssionMT | 43.361 | **74.580** | 0.998 | **1.000** | **1.000** | **1.000** | 0.362 | 0.168 | **0.909** | **0.026** |

Table 70: Czech→Ukrainian, direct (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 47.481 | 70.420 | 0.991 | **1.000** | 0.998 | **1.000** | 0.572 | 0.253 | 0.934 | 0.041 |
| Claude-3 | 52.844 | 75.437 | 0.977 | 0.990 | 0.985 | 0.984 | 0.600 | 0.212 | 0.926 | 0.041 |
| CommandR-plus | 46.760 | 67.874 | 0.974 | **1.000** | 0.996 | 0.999 | 0.621 | 0.218 | 0.938 | 0.038 |
| GPT-4 | 50.991 | 73.421 | **1.000** | **1.000** | **1.000** | **1.000** | 0.657 | 0.203 | **0.951** | 0.031 |
| Llama3-70B | 32.827 | 63.628 | 0.995 | **1.000** | **1.000** | **1.000** | 0.552 | 0.229 | 0.930 | 0.035 |
| IKUN | 50.135 | 74.445 | 0.999 | **1.000** | **1.000** | **1.000** | 0.029 | 0.710 | 0.860 | 0.103 |
| IKUN-C | 36.860 | 65.280 | 0.999 | **1.000** | **1.000** | **1.000** | **0.676** | 0.135 | 0.932 | **0.022** |
| Unbabel-Tower70B | 46.420 | 71.580 | 0.998 | **1.000** | **1.000** | **1.000** | 0.481 | 0.231 | 0.923 | 0.036 |
| BJFU-LPT | 44.301 | 67.430 | 0.999 | **1.000** | **1.000** | **1.000** | 0.027 | 0.717 | 0.860 | 0.104 |
| CUNI-Transformer | 35.405 | 61.580 | 0.999 | **1.000** | **1.000** | **1.000** | 0.487 | 0.344 | 0.912 | 0.052 |
| CycleL | 0.010 | 1.927 | 0.000 | 0.437 | 0.045 | 0.009 | 0.000 | **0.000** | 0.070 | 0.501 |
| IOL_Research | 59.379 | 79.591 | 0.990 | **1.000** | **1.000** | **1.000** | 0.508 | 0.239 | 0.926 | 0.037 |
| ONLINE-A | 44.735 | 73.463 | 0.998 | **1.000** | **1.000** | **1.000** | 0.010 | 0.614 | 0.854 | 0.090 |
| ONLINE-B | 53.593 | 72.722 | 0.999 | **1.000** | **1.000** | **1.000** | 0.043 | 0.512 | 0.835 | 0.076 |
| ONLINE-G | 39.987 | 63.954 | 0.999 | **1.000** | **1.000** | **1.000** | 0.009 | 0.590 | 0.831 | 0.087 |
| ONLINE-W | 38.964 | 66.735 | **1.000** | **1.000** | **1.000** | **1.000** | 0.257 | 0.394 | 0.875 | 0.058 |
| TranssionMT | **64.619** | **83.303** | 0.998 | **1.000** | **1.000** | **1.000** | 0.000 | 0.720 | 0.857 | 0.105 |

Table 71: Czech→Ukrainian, 0-shot (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 48.630 | 74.977 | **0.999** | 0.999 | **1.000** | **1.000** | 0.153 | 0.447 | 0.875 | 0.066 |
| Claude-3 | 0.822 | 10.245 | 0.050 | 0.916 | 0.627 | 0.949 | **0.897** | 0.067 | 0.505 | 0.222 |
| CommandR-plus | 19.125 | 31.488 | 0.154 | 0.958 | 0.832 | 0.909 | 0.501 | 0.424 | 0.584 | 0.228 |
| GPT-4 | 48.135 | 74.266 | **0.999** | **1.000** | **1.000** | **1.000** | 0.259 | 0.285 | 0.894 | 0.042 |
| Llama3-70B | 36.870 | 61.658 | 0.725 | 0.869 | 0.843 | 0.847 | 0.182 | 0.580 | 0.699 | 0.187 |
| IKUN | 49.178 | 73.777 | 0.998 | **1.000** | **1.000** | **1.000** | 0.411 | 0.170 | **0.915** | **0.027** |
| IKUN-C | 31.518 | 60.419 | 0.996 | **1.000** | **1.000** | **1.000** | 0.269 | 0.386 | 0.870 | 0.059 |
| Unbabel-Tower70B | 48.720 | 76.687 | 0.996 | **1.000** | **1.000** | **1.000** | 0.153 | 0.293 | 0.875 | 0.044 |
| BJFU-LPT | 7.069 | 12.867 | **0.999** | 0.998 | 0.578 | 0.266 | 0.000 | 0.077 | 0.413 | 0.178 |
| CUNI-Transformer | 0.813 | 7.072 | **0.999** | 0.971 | 0.087 | 0.000 | 0.000 | 0.059 | 0.294 | 0.290 |
| CycleL | 0.014 | 2.350 | 0.000 | 0.293 | 0.000 | 0.001 | 0.000 | **0.021** | 0.042 | 0.533 |
| IOL_Research | 55.173 | 79.879 | 0.994 | **1.000** | **1.000** | **1.000** | 0.275 | 0.236 | 0.896 | 0.037 |
| ONLINE-A | 57.268 | 80.732 | 0.996 | **1.000** | **1.000** | **1.000** | 0.187 | 0.406 | 0.883 | 0.060 |
| ONLINE-B | 47.307 | 75.251 | 0.998 | **1.000** | **1.000** | **1.000** | 0.237 | 0.288 | 0.885 | 0.044 |
| ONLINE-G | 46.645 | 72.724 | 0.996 | **1.000** | **1.000** | **1.000** | 0.039 | 0.563 | 0.840 | 0.083 |
| ONLINE-W | 48.424 | 68.657 | **0.999** | **1.000** | **1.000** | **1.000** | 0.388 | 0.283 | 0.904 | 0.043 |
| TranssionMT | **64.582** | **83.997** | 0.998 | **1.000** | **1.000** | **1.000** | 0.235 | 0.267 | 0.890 | 0.041 |

Table 72: Czech→Ukrainian, 0-shot (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 13.326 | 39.802 | 0.497 | **1.000** | **1.000** | **1.000** | 0.247 | 0.656 | 0.724 | 0.167 |
| Claude-3 | 13.101 | 39.207 | 0.490 | 0.993 | 0.990 | 0.799 | 0.285 | 0.146 | 0.654 | **0.125** |
| CommandR-plus | 14.038 | 34.504 | 0.487 | **1.000** | **1.000** | 0.796 | 0.275 | 0.155 | 0.653 | 0.125 |
| GPT-4 | 22.718 | 46.429 | 0.499 | **1.000** | **1.000** | **1.000** | 0.297 | 0.650 | **0.797** | 0.165 |
| Llama3-70B | 11.644 | 36.664 | 0.493 | **1.000** | **1.000** | **1.000** | 0.224 | 0.665 | 0.710 | 0.168 |
| IKUN | 14.630 | 41.597 | 0.494 | 0.999 | **1.000** | 0.976 | 0.023 | 0.863 | 0.688 | 0.200 |
| IKUN-C | 12.107 | 37.939 | **0.518** | **1.000** | **1.000** | **1.000** | **0.362** | 0.569 | 0.739 | 0.151 |
| Unbabel-Tower70B | 13.321 | 40.229 | 0.498 | **1.000** | **1.000** | 0.996 | 0.214 | 0.660 | 0.716 | 0.167 |
| BJFU-LPT | 16.880 | 41.573 | 0.499 | **1.000** | **1.000** | 0.504 | 0.017 | 0.479 | 0.646 | 0.212 |
| CUNI-Transformer | 4.540 | 32.380 | 0.499 | 0.808 | **1.000** | 0.995 | 0.211 | 0.737 | 0.645 | 0.206 |
| CycleL | 0.035 | 1.793 | 0.000 | 0.651 | 0.028 | 0.000 | 0.000 | **0.000** | 0.097 | 0.474 |
| IOL_Research | 13.454 | 41.451 | 0.494 | **1.000** | 0.998 | **1.000** | 0.220 | 0.661 | 0.685 | 0.168 |
| ONLINE-A | 29.551 | 50.897 | 0.498 | **1.000** | **1.000** | 0.955 | 0.002 | 0.830 | 0.778 | 0.197 |
| ONLINE-B | 29.834 | 43.079 | 0.499 | **1.000** | **1.000** | 0.563 | 0.005 | 0.791 | 0.691 | 0.247 |
| ONLINE-G | 23.284 | 40.978 | 0.499 | **1.000** | **1.000** | **1.000** | 0.005 | 0.854 | 0.765 | 0.194 |
| ONLINE-W | 16.322 | 43.296 | 0.499 | **1.000** | **1.000** | 0.996 | 0.121 | 0.742 | 0.729 | 0.179 |
| TranssionMT | **30.538** | **53.582** | 0.499 | **1.000** | **1.000** | 0.974 | 0.001 | 0.895 | 0.779 | 0.204 |

Table 73: Czech→Ukrainian, 1-shot (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 12.403 | 39.988 | 0.501 | **1.000** | **1.000** | **1.000** | 0.157 | 0.733 | 0.687 | 0.177 |
| Claude-3 | 5.186 | 18.859 | 0.007 | 0.945 | 0.649 | 0.965 | **0.476** | 0.513 | 0.484 | 0.281 |
| CommandR-plus | 5.225 | 21.346 | 0.113 | 0.930 | 0.824 | 0.949 | 0.269 | 0.219 | 0.524 | 0.202 |
| GPT-4 | 20.201 | 44.942 | 0.499 | **1.000** | **1.000** | 0.982 | 0.246 | 0.655 | 0.773 | 0.168 |
| Llama3-70B | 21.258 | 45.719 | 0.470 | 0.998 | 0.996 | 0.998 | 0.170 | 0.734 | 0.759 | 0.182 |
| IKUN | 14.652 | 41.278 | 0.497 | **1.000** | **1.000** | **1.000** | 0.390 | 0.596 | 0.734 | 0.158 |
| IKUN-C | 14.705 | 38.737 | 0.499 | **1.000** | **1.000** | **1.000** | 0.231 | 0.750 | 0.738 | 0.180 |
| Unbabel-Tower70B | 13.789 | 41.805 | 0.498 | **1.000** | **1.000** | 0.999 | 0.177 | 0.667 | 0.705 | 0.168 |
| BJFU-LPT | 15.776 | 21.214 | **0.515** | 0.999 | 0.655 | 0.166 | 0.000 | 0.037 | 0.406 | 0.243 |
| CUNI-Transformer | 5.529 | 11.149 | 0.499 | 0.938 | 0.435 | 0.000 | 0.001 | 0.067 | 0.268 | 0.315 |
| CycleL | 0.144 | 2.202 | 0.000 | 0.465 | 0.000 | 0.000 | 0.000 | **0.000** | 0.066 | 0.505 |
| IOL_Research | 8.829 | 36.552 | 0.508 | 0.860 | 0.895 | 0.977 | 0.252 | 0.268 | 0.642 | **0.147** |
| ONLINE-A | 25.043 | 51.472 | 0.510 | **1.000** | **1.000** | **1.000** | 0.187 | 0.741 | 0.786 | 0.176 |
| ONLINE-B | 23.961 | 41.456 | 0.499 | **1.000** | **1.000** | 0.940 | 0.235 | 0.725 | 0.731 | 0.184 |
| ONLINE-G | 20.539 | 43.621 | 0.499 | **1.000** | **1.000** | **1.000** | 0.043 | 0.906 | 0.722 | 0.201 |
| ONLINE-W | 22.623 | 46.890 | 0.499 | **1.000** | **1.000** | **1.000** | 0.379 | 0.617 | **0.809** | 0.161 |
| TranssionMT | **25.323** | **52.901** | 0.499 | **1.000** | **1.000** | **1.000** | 0.204 | 0.727 | 0.787 | 0.176 |

Table 74: Czech→Ukrainian, 1-shot (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 2.147 | 18.888 | 0.999 | **0.996** | 0.996 | 0.996 | 0.825 | 0.138 | **0.938** | 0.026 |
| Claude-3 | 1.009 | 13.452 | 0.847 | 0.864 | 0.408 | 0.383 | 0.695 | 0.201 | 0.559 | 0.244 |
| CommandR-plus | 2.182 | 14.449 | 0.668 | 0.813 | 0.771 | 0.808 | 0.543 | 0.410 | 0.682 | 0.197 |
| GPT-4 | 6.110 | 23.763 | **1.000** | 0.956 | 0.384 | 0.356 | 0.815 | 0.051 | 0.600 | 0.194 |
| Llama3-70B | 2.096 | 18.345 | 0.999 | 0.988 | **0.999** | **1.000** | 0.864 | 0.108 | 0.921 | **0.020** |
| IKUN | 8.386 | 22.957 | 0.966 | 0.911 | 0.397 | 0.384 | 0.666 | 0.113 | 0.552 | 0.210 |
| IKUN-C | 1.859 | 17.569 | 0.747 | 0.728 | 0.694 | 0.880 | 0.386 | 0.311 | 0.610 | 0.184 |
| Unbabel-Tower70B | 2.233 | 19.069 | 0.996 | 0.993 | 0.984 | 0.983 | 0.808 | 0.152 | 0.930 | 0.031 |
| BJFU-LPT | 1.363 | 16.719 | 0.000 | 0.058 | 0.141 | 0.716 | 0.000 | 0.072 | 0.131 | 0.456 |
| CUNI-Transformer | 0.010 | 10.501 | 0.047 | 0.048 | 0.137 | 0.353 | 0.017 | 0.088 | 0.086 | 0.501 |
| CycleL | 0.003 | 1.416 | 0.000 | 0.039 | 0.027 | 0.053 | 0.000 | 0.017 | 0.017 | 0.557 |
| IOL_Research | 1.788 | 15.626 | 0.698 | 0.692 | 0.733 | 0.916 | 0.591 | 0.247 | 0.703 | 0.176 |
| ONLINE-A | **12.677** | **29.133** | 0.995 | 0.958 | 0.930 | 0.814 | **0.871** | 0.089 | 0.798 | 0.057 |
| ONLINE-B | 10.813 | 18.183 | 0.908 | 0.934 | 0.913 | 0.914 | 0.632 | 0.197 | 0.785 | 0.079 |
| ONLINE-G | 7.923 | 18.949 | 0.006 | 0.006 | 0.207 | 0.005 | 0.005 | **0.002** | 0.033 | 0.540 |
| ONLINE-W | 5.625 | 22.669 | 0.987 | 0.936 | 0.847 | 0.807 | 0.742 | 0.141 | 0.778 | 0.083 |
| TranssionMT | 10.790 | 27.831 | 0.930 | 0.938 | 0.923 | 0.922 | 0.656 | 0.209 | 0.803 | 0.074 |

Table 75: Czech→Ukrainian, 0-shot JSON format (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 2.060 | 16.944 | **0.999** | 0.996 | 0.990 | 0.990 | 0.979 | 0.017 | 0.973 | 0.009 |
| Claude-3 | **13.359** | 20.291 | 0.993 | 0.785 | 0.518 | 0.498 | 0.492 | 0.506 | 0.610 | 0.250 |
| CommandR-plus | 0.959 | 12.553 | 0.727 | 0.788 | 0.786 | 0.876 | 0.690 | 0.193 | 0.742 | 0.148 |
| GPT-4 | 4.594 | 22.466 | **0.999** | 0.562 | 0.089 | 0.061 | 0.073 | 0.909 | 0.276 | 0.463 |
| Llama3-70B | 3.365 | 20.334 | 0.998 | 0.995 | **0.999** | **1.000** | **0.987** | **0.013** | 0.973 | **0.006** |
| IKUN | 2.207 | 18.989 | 0.987 | 0.988 | 0.988 | 0.990 | 0.966 | 0.028 | 0.956 | 0.013 |
| IKUN-C | 3.175 | 20.407 | 0.958 | 0.847 | 0.755 | 0.764 | 0.797 | 0.130 | 0.778 | 0.119 |
| Unbabel-Tower70B | 2.207 | 18.509 | 0.976 | 0.982 | 0.984 | 0.995 | 0.962 | 0.038 | 0.966 | 0.018 |
| BJFU-LPT | 2.373 | 17.972 | 0.475 | 0.289 | 0.122 | 0.253 | 0.086 | 0.554 | 0.177 | 0.492 |
| CUNI-Transformer | 0.067 | 9.490 | 0.038 | 0.033 | 0.001 | 0.001 | 0.000 | 0.159 | 0.010 | 0.646 |
| CycleL | 0.010 | 1.562 | 0.000 | 0.011 | 0.061 | 0.015 | 0.000 | 0.028 | 0.012 | 0.564 |
| IOL_Research | 0.959 | 10.262 | 0.259 | 0.258 | 0.362 | 0.852 | 0.241 | 0.299 | 0.340 | 0.370 |
| ONLINE-A | 5.662 | **26.044** | 0.965 | 0.971 | 0.977 | 0.998 | 0.952 | 0.043 | 0.960 | 0.022 |
| ONLINE-B | 5.679 | 15.915 | 0.980 | **0.998** | 0.985 | 0.994 | 0.965 | 0.029 | 0.963 | 0.014 |
| ONLINE-G | 4.470 | 17.782 | 0.098 | 0.146 | 0.258 | 0.870 | 0.076 | 0.525 | 0.228 | 0.457 |
| ONLINE-W | 5.600 | 25.450 | 0.998 | 0.996 | **0.999** | 0.998 | 0.972 | 0.024 | **0.978** | 0.007 |
| TranssionMT | 5.662 | 26.039 | 0.982 | **0.998** | 0.985 | 0.999 | 0.968 | 0.027 | 0.968 | 0.013 |

Table 76: Czech→Ukrainian, 0-shot JSON format (non-English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 1.209 | 15.930 | 0.988 | 0.980 | 0.988 | 0.990 | 0.846 | 0.118 | **0.939** | 0.028 |
| Claude-3 | 0.175 | 9.852 | 0.640 | 0.686 | 0.096 | 0.069 | 0.512 | 0.370 | 0.299 | 0.411 |
| CommandR-plus | 0.945 | 12.276 | 0.801 | 0.892 | 0.931 | 0.961 | 0.659 | 0.314 | 0.825 | 0.109 |
| GPT-4 | 2.235 | 17.932 | **1.000** | 0.971 | 0.669 | 0.667 | 0.855 | 0.074 | 0.776 | 0.112 |
| Llama3-70B | 1.144 | 15.433 | **1.000** | **0.993** | **1.000** | **1.000** | 0.860 | 0.120 | 0.933 | **0.020** |
| IKUN | 8.349 | 21.330 | **1.000** | 0.919 | 0.042 | 0.000 | 0.743 | 0.025 | 0.391 | 0.295 |
| IKUN-C | 0.784 | 14.040 | 0.593 | 0.547 | 0.527 | 0.806 | 0.240 | 0.287 | 0.460 | 0.267 |
| Unbabel-Tower70B | 1.265 | 16.248 | 0.988 | 0.975 | 0.978 | 0.983 | 0.806 | 0.169 | 0.923 | 0.040 |
| BJFU-LPT | 0.161 | 11.288 | 0.000 | 0.110 | 0.103 | 0.534 | 0.000 | 0.037 | 0.107 | 0.474 |
| CUNI-Transformer | 0.000 | 7.124 | 0.118 | 0.098 | 0.218 | 0.404 | 0.054 | 0.118 | 0.129 | 0.470 |
| CycleL | 0.000 | 1.146 | 0.000 | 0.032 | 0.039 | 0.056 | 0.000 | 0.032 | 0.018 | 0.558 |
| IOL_Research | 0.997 | 12.974 | 0.713 | 0.708 | 0.730 | 0.926 | 0.591 | 0.248 | 0.714 | 0.172 |
| ONLINE-A | **8.813** | **24.505** | 0.993 | 0.949 | 0.936 | 0.831 | **0.877** | 0.078 | 0.801 | 0.054 |
| ONLINE-B | 7.119 | 15.058 | 0.900 | 0.912 | 0.909 | 0.907 | 0.637 | 0.191 | 0.778 | 0.083 |
| ONLINE-G | 4.738 | 15.217 | 0.007 | 0.007 | 0.206 | 0.010 | 0.005 | **0.000** | 0.034 | 0.539 |
| ONLINE-W | 4.298 | 19.352 | 0.980 | 0.934 | 0.860 | 0.838 | 0.686 | 0.189 | 0.789 | 0.085 |
| TranssionMT | 7.093 | 23.382 | 0.912 | 0.919 | 0.914 | 0.924 | 0.686 | 0.184 | 0.794 | 0.076 |

Table 77: Czech→Ukrainian, 1-shot JSON format (English source)

| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | SAAvg |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 1.151 | 13.983 | **1.000** | **1.000** | **1.000** | **1.000** | 0.976 | 0.024 | 0.979 | **0.006** |
| Claude-3 | **11.965** | 17.619 | 0.939 | 0.482 | 0.061 | 0.044 | 0.034 | 0.954 | 0.234 | 0.495 |
| CommandR-plus | 0.845 | 12.246 | 0.697 | 0.809 | 0.824 | 0.848 | 0.689 | 0.249 | 0.745 | 0.157 |
| GPT-4 | 2.116 | 17.765 | 0.995 | 0.504 | 0.022 | 0.005 | 0.024 | 0.934 | 0.225 | 0.492 |
| Llama3-70B | 1.113 | 14.895 | 0.993 | 0.993 | 0.995 | 0.998 | 0.971 | 0.029 | 0.968 | 0.010 |
| IKUN | 1.116 | 15.848 | 0.949 | 0.958 | 0.963 | 0.980 | 0.936 | 0.037 | 0.932 | 0.030 |
| IKUN-C | 1.367 | 16.257 | 0.914 | 0.949 | 0.958 | 0.990 | 0.912 | 0.073 | 0.919 | 0.043 |
| Unbabel-Tower70B | 1.132 | 15.429 | 0.988 | **1.000** | 0.998 | 0.998 | 0.973 | 0.027 | 0.980 | 0.009 |
| BJFU-LPT | 0.452 | 13.846 | 0.213 | 0.174 | 0.120 | 0.377 | 0.000 | 0.721 | 0.126 | 0.550 |
| CUNI-Transformer | 0.001 | 6.806 | 0.117 | 0.073 | 0.000 | 0.000 | 0.000 | 0.254 | 0.027 | 0.646 |
| CycleL | 0.001 | 1.255 | 0.000 | 0.020 | 0.073 | 0.029 | 0.000 | 0.032 | 0.017 | 0.560 |
| IOL_Research | 0.655 | 8.578 | 0.147 | 0.147 | 0.249 | 0.765 | 0.115 | 0.430 | 0.235 | 0.452 |
| ONLINE-A | 2.950 | 21.129 | 0.954 | 0.961 | 0.971 | 0.998 | 0.939 | 0.051 | 0.952 | 0.028 |
| ONLINE-B | 2.989 | 12.473 | 0.973 | 0.990 | 0.993 | 0.995 | **0.980** | **0.012** | 0.969 | 0.012 |
| ONLINE-G | 2.261 | 13.608 | 0.054 | 0.110 | 0.244 | 0.961 | 0.020 | 0.545 | 0.207 | 0.465 |
| ONLINE-W | 2.936 | 20.718 | 0.998 | 0.995 | **1.000** | **1.000** | 0.971 | 0.029 | **0.981** | 0.007 |
| TranssionMT | 2.970 | **21.130** | 0.971 | 0.990 | 0.993 | 0.998 | 0.976 | 0.017 | 0.971 | 0.013 |

Table 78: Czech→Ukrainian, 1-shot JSON format (non-English source)

## A.2 Summary results

### A.2.1 Weakest attacks

| System | clean | | adversarial | | | | | | | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | |
| Aya23 | 50.124 | 69.491 | 0.995 | 0.990 | 0.988 | 1.000 | 0.886 | 0.114 | **0.972** | direct |
| Claude-3 | 63.945 | 80.516 | 0.973 | 0.979 | 0.977 | 0.979 | 0.005 | 0.744 | 0.837 | 1-shot |
| CommandR-plus | 51.532 | 70.648 | 0.963 | 1.000 | 0.998 | 0.999 | 0.000 | 0.905 | 0.850 | 1-shot |
| GPT-4 | 58.671 | 76.248 | 0.999 | 1.000 | 1.000 | 1.000 | 0.129 | 0.379 | 0.875 | 0-shot |
| Llama3-70B | 55.838 | 73.779 | 0.999 | 0.917 | 0.942 | 0.978 | 0.973 | 0.026 | 0.958 | 0-shot JSON format |
| NVIDIA-NeMo | 53.441 | 71.047 | 0.982 | 0.951 | 0.983 | 1.000 | 0.848 | 0.152 | 0.943 | direct |
| CUNI-DS | 45.865 | 65.698 | 0.985 | 0.907 | 0.930 | 0.985 | 0.953 | 0.038 | 0.933 | 1-shot JSON format |
| IKUN | 46.017 | 65.324 | 0.973 | 0.884 | 0.909 | 0.968 | 0.936 | 0.055 | 0.915 | 1-shot JSON format |
| IKUN-C | 39.794 | 60.823 | 0.979 | 0.848 | 0.864 | 0.966 | 0.927 | 0.040 | 0.893 | 0-shot JSON format |
| Unbabel-Tower70B | 54.457 | 73.925 | 0.995 | 0.960 | 0.963 | 1.000 | 0.670 | 0.329 | 0.901 | direct |
| Yandex | 42.793 | 65.032 | 0.780 | 0.969 | 0.990 | 1.000 | 0.845 | 0.155 | 0.899 | direct |
| CycleL | 1.720 | 19.371 | 0.967 | 0.985 | 0.842 | 0.984 | 0.000 | 0.879 | 0.547 | 0-shot |
| CycleL2 | 0.823 | 15.256 | 0.977 | 0.652 | 0.554 | 0.998 | 0.000 | 0.162 | 0.456 | direct |
| Dubformer | 0.811 | 2.480 | 0.999 | 0.450 | 0.048 | 0.000 | 0.001 | 0.136 | 0.218 | 0-shot |
| IOL_Research | 62.421 | 77.519 | 0.995 | 0.851 | 0.868 | 0.895 | 0.895 | 0.032 | 0.889 | 1-shot JSON format |
| ONLINE-A | 57.977 | 75.168 | 0.999 | 0.925 | 0.942 | 0.976 | 0.958 | 0.042 | 0.960 | 0-shot JSON format |
| ONLINE-B | 55.403 | 73.776 | 0.994 | 0.913 | 0.947 | 0.972 | 0.945 | 0.050 | 0.951 | 1-shot JSON format |
| ONLINE-G | 53.353 | 74.154 | 0.999 | 0.973 | 0.995 | 1.000 | 0.902 | 0.098 | 0.957 | direct |
| ONLINE-W | 53.906 | 72.810 | 0.998 | 0.919 | 0.947 | 0.985 | 0.969 | 0.029 | 0.958 | 1-shot JSON format |
| TSU-HITs | 22.052 | 43.818 | 0.124 | 0.813 | 0.949 | 0.996 | 0.721 | 0.257 | 0.651 | direct |
| TranssionMT | 55.300 | 74.002 | 0.996 | 0.917 | 0.949 | 0.974 | 0.947 | 0.053 | 0.954 | 1-shot JSON format |

Table 79: English→Russian; weakest attack by Avg. win

| System | clean | | adversarial | | | | | | | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | Avg. win | |
| Aya23 | 60.528 | 77.596 | 0.995 | 0.928 | 0.962 | 0.995 | 0.999 | 0.000 | 0.971 | 0-shot JSON format |
| Claude-3 | 69.372 | 84.126 | 0.999 | 0.955 | 0.978 | 0.998 | 0.998 | 0.000 | 0.986 | 0-shot JSON format |
| CommandR-plus | 60.904 | 78.355 | 0.939 | 1.000 | 1.000 | 1.000 | 0.734 | 0.171 | 0.953 | 1-shot |
| GPT-4 | 70.239 | 84.067 | 1.000 | 1.000 | 1.000 | 1.000 | 0.865 | 0.037 | 0.981 | 1-shot |
| Llama3-70B | 64.414 | 79.829 | 0.998 | 1.000 | 1.000 | 1.000 | 0.922 | 0.009 | 0.988 | 1-shot |
| NVIDIA-NeMo | 62.179 | 77.817 | 0.968 | 0.988 | 0.984 | 1.000 | 0.994 | 0.005 | 0.989 | direct |
| AIST-AIRC | 54.511 | 72.781 | 0.999 | 0.996 | 0.996 | 1.000 | 0.980 | 0.009 | 0.994 | direct |
| CUNI-NL | 51.442 | 69.699 | 0.761 | 1.000 | 0.999 | 0.999 | 0.988 | 0.005 | 0.964 | direct |
| IKUN | 51.652 | 70.262 | 0.994 | 0.999 | 0.990 | 1.000 | 0.856 | 0.073 | 0.974 | 1-shot |
| IKUN-C | 44.710 | 65.240 | 0.989 | 1.000 | 1.000 | 1.000 | 0.890 | 0.060 | 0.982 | 1-shot |
| Unbabel-Tower70B | 61.008 | 78.193 | 0.995 | 1.000 | 1.000 | 1.000 | 0.985 | 0.004 | 0.997 | 1-shot |
| CycleL | 20.487 | 44.322 | 0.987 | 0.996 | 0.998 | 1.000 | 0.000 | 0.589 | 0.781 | 0-shot |
| CycleL2 | 20.487 | 44.322 | 0.987 | 0.996 | 0.998 | 1.000 | 0.000 | 0.589 | 0.781 | 0-shot |
| Dubformer | 26.213 | 32.808 | 0.272 | 0.483 | 0.515 | 0.857 | 0.196 | 0.748 | 0.484 | direct |
| IOL_Research | 69.214 | 82.833 | 0.999 | 1.000 | 1.000 | 1.000 | 0.820 | 0.055 | 0.974 | 1-shot |
| MSLC | 41.196 | 64.234 | 0.972 | 1.000 | 0.999 | 1.000 | 0.529 | 0.084 | 0.928 | 1-shot |
| ONLINE-A | 68.859 | 82.629 | 0.999 | 0.999 | 0.999 | 1.000 | 0.999 | 0.000 | **0.999** | direct |
| ONLINE-B | 54.922 | 74.946 | 0.998 | 1.000 | 1.000 | 1.000 | 0.823 | 0.037 | 0.974 | 1-shot |
| ONLINE-G | 68.624 | 82.302 | 0.999 | 0.993 | 0.994 | 1.000 | 0.995 | 0.005 | 0.995 | direct |
| ONLINE-W | 61.546 | 78.220 | 0.961 | 0.999 | 0.999 | 1.000 | 0.999 | 0.001 | 0.994 | direct |
| TSU-HITs | 29.868 | 49.567 | 0.144 | 0.652 | 0.853 | 0.946 | 0.353 | 0.168 | 0.526 | direct |
| TranssionMT | 54.873 | 74.941 | 0.998 | 1.000 | 1.000 | 1.000 | 0.825 | 0.037 | 0.975 | 1-shot |

Table 80: English→German; weakest attack by Avg. win

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
| Aya23 | 19.085 | 40.614 | 0.993 | 0.714 | 1.000 | 1.000 | 0.005 | 0.892 | 0.673 | 1-shot |
| Claude-3 | 1.919 | 53.543 | 0.953 | 0.157 | 0.815 | 0.878 | 0.867 | 0.071 | 0.668 | 0-shot JSON format |
| CommandR-plus | 14.366 | 43.986 | 0.856 | 0.118 | 0.968 | 1.000 | 0.002 | 0.928 | 0.559 | 1-shot |
| GPT-4 | 17.514 | 54.097 | 0.999 | 0.001 | 1.000 | 1.000 | 0.015 | 0.955 | 0.574 | 0-shot |
| Llama3-70B | 27.898 | 43.181 | 0.979 | 0.965 | 1.000 | 1.000 | 0.000 | 0.956 | 0.706 | 1-shot |
| NVIDIA-NeMo | 2.076 | 35.694 | 0.829 | 0.991 | 1.000 | 1.000 | 0.002 | 0.974 | 0.693 | 1-shot |
| AIST-AIRC | 0.719 | 34.974 | 0.969 | 0.993 | 1.000 | 1.000 | 0.010 | 0.916 | 0.710 | 1-shot |
| IKUN | 13.311 | 31.025 | 0.980 | 0.998 | 1.000 | 1.000 | 0.002 | 0.897 | 0.711 | 1-shot |
| IKUN-C | 2.249 | 26.016 | 0.941 | 0.012 | 0.816 | 0.920 | 0.911 | 0.054 | 0.597 | 0-shot JSON format |
| Unbabel-Tower70B | 8.143 | 41.692 | 0.990 | 0.105 | 0.903 | 0.987 | 0.935 | 0.058 | 0.703 | 0-shot JSON format |
| CycleL | 0.041 | 3.364 | 0.006 | 0.949 | 0.463 | 0.998 | 0.000 | 0.061 | 0.345 | 1-shot |
| DLUT_GTCOM | 0.813 | 42.293 | 0.971 | 0.980 | 1.000 | 1.000 | 0.011 | 0.810 | 0.709 | 1-shot |
| IOL_Research | 19.182 | 51.107 | 0.971 | 0.985 | 1.000 | 1.000 | 0.004 | 0.969 | 0.709 | 1-shot |
| NTTSU | 4.594 | 33.132 | 0.865 | 0.001 | 0.953 | 0.998 | 0.789 | 0.207 | 0.646 | direct |
| ONLINE-A | 1.220 | 44.459 | 0.920 | 0.995 | 1.000 | 1.000 | 0.015 | 0.840 | 0.704 | 1-shot |
| ONLINE-B | 1.015 | 44.589 | 0.996 | 0.996 | 1.000 | 1.000 | 0.028 | 0.889 | **0.717** | 1-shot |
| ONLINE-G | 3.339 | 45.429 | 0.995 | 0.989 | 1.000 | 1.000 | 0.011 | 0.968 | 0.714 | 1-shot |
| ONLINE-W | 4.871 | 34.170 | 0.989 | 0.000 | 1.000 | 1.000 | 0.002 | 0.982 | 0.571 | 1-shot |
| Team-J | 0.416 | 36.323 | 0.998 | 0.998 | 1.000 | 1.000 | 0.002 | 0.994 | 0.714 | 1-shot |
| UvA-MT | 1.159 | 43.238 | 0.908 | 0.004 | 0.903 | 0.999 | 0.851 | 0.147 | 0.658 | direct |

Table 81: English→Japanese; weakest attack by Avg. win

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
| Aya23 | 44.375 | 63.672 | 0.968 | 0.960 | 0.958 | 0.998 | 0.657 | 0.319 | 0.926 | direct |
| Claude-3 | 60.166 | 76.954 | 0.060 | 0.453 | 0.557 | 0.157 | 0.083 | 0.610 | 0.205 | 0-shot |
| CommandR-plus | 39.996 | 61.592 | 0.963 | 1.000 | 0.996 | 1.000 | 0.033 | 0.760 | 0.853 | 1-shot |
| GPT-4 | 50.565 | 69.608 | 0.998 | 1.000 | 1.000 | 1.000 | 0.023 | 0.962 | 0.860 | 1-shot |
| Llama3-70B | 51.601 | 69.311 | 0.998 | 0.862 | 0.927 | 0.980 | 0.949 | 0.039 | 0.930 | 0-shot JSON format |
| NVIDIA-NeMo | 47.354 | 66.582 | 0.980 | 0.993 | 0.991 | 1.000 | 0.703 | 0.285 | 0.951 | direct |
| IKUN | 40.887 | 60.362 | 0.924 | 0.987 | 0.980 | 1.000 | 0.673 | 0.304 | 0.926 | direct |
| IKUN-C | 35.290 | 56.369 | 0.956 | 0.978 | 0.963 | 1.000 | 0.681 | 0.304 | 0.922 | direct |
| Unbabel-Tower70B | 56.242 | 74.129 | 0.998 | 0.994 | 0.995 | 1.000 | 0.700 | 0.293 | 0.954 | direct |
| CycleL | 0.268 | 12.822 | 0.000 | 0.284 | 0.370 | 1.000 | 0.001 | 0.119 | 0.237 | direct |
| IOL_Research | 53.133 | 70.132 | 0.998 | 0.870 | 0.936 | 0.999 | 0.969 | 0.023 | 0.942 | 1-shot JSON format |
| ONLINE-A | 59.021 | 74.613 | 0.999 | 0.996 | 0.994 | 1.000 | 0.707 | 0.283 | 0.956 | direct |
| ONLINE-B | 56.473 | 71.907 | 0.998 | 0.898 | 0.951 | 0.999 | 0.950 | 0.032 | 0.953 | 0-shot JSON format |
| ONLINE-G | 55.704 | 72.554 | 0.999 | 0.998 | 0.991 | 1.000 | 0.705 | 0.285 | 0.955 | direct |
| ONLINE-W | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | direct |
| TranssionMT | 56.588 | 73.267 | 0.999 | 0.909 | 0.958 | 1.000 | 0.967 | 0.022 | **0.965** | 0-shot JSON format |

Table 82: English→Hindi; weakest attack by Avg. win

440

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
| Aya23 | 71.590 | 83.455 | 0.985 | 0.999 | 0.998 | 1.000 | 0.956 | 0.044 | 0.991 | direct |
| Claude-3 | 77.382 | 88.287 | 0.994 | 0.947 | 0.977 | 0.990 | 0.993 | 0.005 | 0.982 | 0-shot JSON format |
| CommandR-plus | 69.366 | 82.843 | 0.453 | 0.957 | 0.955 | 0.917 | 0.222 | 0.649 | 0.757 | 1-shot |
| GPT-4 | 76.485 | 86.879 | 0.999 | 1.000 | 1.000 | 1.000 | 0.862 | 0.067 | 0.980 | 1-shot |
| Llama3-70B | 75.659 | 85.899 | 0.996 | 0.938 | 0.971 | 0.995 | 0.999 | 0.001 | 0.982 | 0-shot JSON format |
| NVIDIA-NeMo | 71.684 | 83.575 | 0.984 | 0.995 | 0.990 | 0.999 | 0.978 | 0.022 | 0.992 | direct |
| IKUN | 56.366 | 73.524 | 0.979 | 0.879 | 0.903 | 0.996 | 0.980 | 0.000 | 0.949 | 0-shot JSON format |
| IKUN-C | 52.543 | 70.275 | 0.989 | 1.000 | 1.000 | 1.000 | 0.865 | 0.094 | 0.979 | 1-shot |
| Occiglot | 49.361 | 68.297 | 0.951 | 0.919 | 0.862 | 1.000 | 0.958 | 0.029 | 0.922 | direct |
| Unbabel-Tower70B | 58.762 | 76.431 | 0.989 | 1.000 | 1.000 | 1.000 | 0.974 | 0.011 | 0.995 | 1-shot |
| CycleL | 32.147 | 51.642 | 0.985 | 1.000 | 1.000 | 1.000 | 0.001 | 0.174 | 0.855 | 1-shot |
| Dubformer | 60.120 | 79.825 | 0.952 | 0.793 | 0.519 | 0.468 | 0.088 | 0.386 | 0.670 | 0-shot |
| IOL_Research | 76.839 | 86.496 | 0.994 | 0.929 | 0.960 | 0.994 | 0.993 | 0.005 | 0.976 | 0-shot JSON format |
| MSLC | 56.800 | 74.431 | 0.994 | 0.945 | 0.836 | 1.000 | 0.894 | 0.093 | 0.930 | direct |
| ONLINE-A | 74.616 | 85.820 | 0.987 | 1.000 | 1.000 | 1.000 | 0.966 | 0.034 | 0.993 | direct |
| ONLINE-B | 72.932 | 83.788 | 0.994 | 0.999 | 1.000 | 1.000 | 0.979 | 0.021 | **0.996** | direct |
| ONLINE-G | 76.360 | 86.243 | 0.999 | 0.944 | 0.974 | 0.995 | 0.999 | 0.000 | 0.986 | 1-shot JSON format |
| ONLINE-W | 58.478 | 74.701 | 0.994 | 0.998 | 0.994 | 0.999 | 0.978 | 0.022 | 0.994 | direct |
| TSU-HITs | 24.907 | 50.317 | 0.093 | 1.000 | 0.991 | 1.000 | 0.012 | 0.706 | 0.728 | 1-shot |
| TranssionMT | 73.144 | 85.551 | 0.990 | 1.000 | 1.000 | 1.000 | 0.966 | 0.034 | 0.994 | direct |

Table 83: English→Spanish; weakest attack by Avg. win

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
| Aya23 | 57.243 | 74.550 | 0.994 | 0.927 | 0.953 | 0.995 | 0.991 | 0.007 | 0.963 | 1-shot JSON format |
| Claude-3 | 66.823 | 81.945 | 0.973 | 0.944 | 0.956 | 0.976 | 0.968 | 0.026 | 0.962 | 0-shot JSON format |
| CommandR-plus | 54.377 | 73.408 | 0.868 | 1.000 | 0.996 | 1.000 | 0.419 | 0.428 | 0.895 | 1-shot |
| GPT-4 | 64.985 | 79.784 | 1.000 | 0.998 | 0.995 | 1.000 | 0.775 | 0.209 | 0.966 | 0-shot |
| Llama3-70B | 61.753 | 77.069 | 0.999 | 0.968 | 0.974 | 0.995 | 0.998 | 0.001 | 0.981 | 0-shot JSON format |
| NVIDIA-NeMo | 55.940 | 72.507 | 0.987 | 0.994 | 0.979 | 1.000 | 0.993 | 0.005 | 0.989 | direct |
| CUNI-MH | 57.511 | 75.301 | 0.998 | 1.000 | 1.000 | 1.000 | 0.988 | 0.012 | 0.998 | direct |
| IKUN | 45.469 | 65.478 | 0.998 | 1.000 | 1.000 | 1.000 | 0.600 | 0.372 | 0.941 | 0-shot |
| IKUN-C | 37.968 | 58.621 | 0.989 | 1.000 | 1.000 | 1.000 | 0.542 | 0.259 | 0.933 | 1-shot |
| SCIR-MT | 63.339 | 78.457 | 0.987 | 1.000 | 0.999 | 1.000 | 0.989 | 0.009 | 0.996 | direct |
| Unbabel-Tower70B | 51.206 | 71.180 | 0.969 | 0.993 | 0.988 | 1.000 | 0.979 | 0.017 | 0.982 | direct |
| CUNI-DocTransformer | 58.378 | 75.431 | 0.998 | 0.996 | 0.996 | 1.000 | 0.991 | 0.007 | 0.997 | direct |
| CUNI-GA | 56.400 | 74.149 | 0.987 | 0.968 | 0.936 | 1.000 | 0.968 | 0.031 | 0.952 | direct |
| CUNI-Transformer | 56.400 | 74.149 | 0.987 | 0.968 | 0.936 | 1.000 | 0.968 | 0.031 | 0.952 | direct |
| CycleL | 1.469 | 17.798 | 0.984 | 0.978 | 0.824 | 0.995 | 0.000 | 0.098 | 0.541 | 0-shot |
| CycleL2 | 5.734 | 24.422 | 0.974 | 0.996 | 0.996 | 1.000 | 0.000 | 0.181 | 0.704 | 1-shot |
| IOL_Research | 64.617 | 78.908 | 0.993 | 0.939 | 0.956 | 0.983 | 0.979 | 0.007 | 0.963 | 0-shot JSON format |
| ONLINE-A | 63.853 | 79.054 | 0.999 | 1.000 | 0.998 | 1.000 | 0.971 | 0.029 | 0.994 | direct |
| ONLINE-B | 59.851 | 76.425 | 0.998 | 1.000 | 0.999 | 1.000 | 0.998 | 0.002 | 0.999 | direct |
| ONLINE-G | 63.404 | 78.063 | 0.999 | 1.000 | 1.000 | 1.000 | 0.995 | 0.002 | **0.999** | direct |
| ONLINE-W | 55.114 | 73.094 | 0.998 | 1.000 | 1.000 | 1.000 | 0.998 | 0.002 | 0.999 | direct |
| TSU-HITs | 16.169 | 34.946 | 0.029 | 0.749 | 0.843 | 0.978 | 0.449 | 0.177 | 0.600 | direct |
| TranssionMT | 62.123 | 78.598 | 0.999 | 1.000 | 0.998 | 1.000 | 0.972 | 0.028 | 0.995 | direct |

Table 84: English→Czech; weakest attack by Avg. win

|  | clean | | adversarial | | | | | | | |
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 31.789 | 44.156 | 0.996 | 0.004 | 0.976 | 1.000 | 0.968 | 0.027 | 0.698 | direct |
| Claude-3 | 6.160 | 52.796 | 0.909 | 0.201 | 0.829 | 0.902 | 0.875 | 0.075 | 0.678 | 0-shot JSON format |
| CommandR-plus | 12.165 | 44.248 | 0.918 | 0.007 | 0.834 | 0.969 | 0.869 | 0.109 | 0.629 | direct |
| GPT-4 | 11.698 | 49.684 | 0.999 | 0.007 | 1.000 | 1.000 | 0.091 | 0.882 | 0.587 | 0-shot |
| Llama3-70B | 14.886 | 45.139 | 0.994 | 0.126 | 0.909 | 0.987 | 0.961 | 0.017 | 0.706 | 1-shot JSON format |
| NVIDIA-NeMo | 1.315 | 35.577 | 0.983 | 0.010 | 0.957 | 0.993 | 0.953 | 0.040 | 0.682 | direct |
| IKUN | 2.722 | 34.863 | 0.973 | 0.004 | 0.922 | 1.000 | 0.968 | 0.028 | 0.665 | direct |
| IKUN-C | 4.261 | 29.898 | 0.965 | 0.006 | 0.951 | 1.000 | 0.961 | 0.033 | 0.671 | direct |
| Unbabel-Tower70B | 2.125 | 40.668 | 0.976 | 0.001 | 0.958 | 1.000 | 0.957 | 0.040 | 0.686 | direct |
| CycleL | 0.057 | 3.676 | 0.912 | 0.006 | 0.994 | 1.000 | 0.007 | 0.121 | 0.424 | direct |
| CycleL2 | 0.000 | 0.779 | 0.006 | 0.000 | 0.109 | 0.973 | 0.000 | 0.127 | 0.155 | 0-shot |
| HW-TSC | 18.593 | 47.754 | 0.999 | 0.004 | 0.941 | 1.000 | 0.984 | 0.013 | 0.689 | direct |
| IOL_Research | 28.529 | 54.058 | 0.988 | 0.002 | 0.862 | 1.000 | 0.889 | 0.108 | 0.667 | direct |
| ONLINE-A | 11.048 | 49.271 | 0.999 | 0.005 | 0.968 | 1.000 | 0.968 | 0.029 | 0.701 | direct |
| ONLINE-B | 2.844 | 45.939 | 0.999 | 0.002 | 0.976 | 1.000 | 0.974 | 0.021 | 0.704 | direct |
| ONLINE-G | 2.939 | 42.534 | 0.999 | 0.148 | 0.896 | 0.994 | 0.953 | 0.031 | **0.709** | 1-shot JSON format |
| ONLINE-W | 3.376 | 44.271 | 0.999 | 0.009 | 0.919 | 1.000 | 0.982 | 0.015 | 0.684 | direct |
| UvA-MT | 0.668 | 34.492 | 0.991 | 0.009 | 0.962 | 1.000 | 0.985 | 0.013 | 0.695 | direct |

Table 85: English→Chinese; weakest attack by Avg. win

|  | clean | | adversarial | | | | | | | |
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 49.791 | 70.611 | 0.998 | 0.917 | 0.945 | 0.999 | 0.973 | 0.022 | 0.958 | 1-shot JSON format |
| Claude-3 | 62.653 | 79.565 | 0.998 | 0.950 | 0.968 | 0.998 | 0.966 | 0.029 | 0.976 | 0-shot JSON format |
| CommandR-plus | 52.187 | 72.206 | 0.884 | 0.869 | 0.863 | 0.958 | 0.756 | 0.225 | 0.872 | direct |
| GPT-4 | 54.848 | 74.440 | 0.994 | 1.000 | 1.000 | 1.000 | 0.032 | 0.792 | 0.861 | 1-shot |
| Llama3-70B | 49.780 | 70.822 | 0.999 | 0.903 | 0.942 | 1.000 | 0.968 | 0.031 | 0.951 | 0-shot JSON format |
| NVIDIA-NeMo | 47.690 | 67.971 | 0.983 | 0.999 | 0.999 | 1.000 | 0.875 | 0.125 | 0.979 | direct |
| IKUN | 34.437 | 58.673 | 0.928 | 0.972 | 0.987 | 1.000 | 0.924 | 0.076 | 0.933 | direct |
| IKUN-C | 36.359 | 59.479 | 0.995 | 0.960 | 0.965 | 1.000 | 0.916 | 0.084 | 0.916 | direct |
| Unbabel-Tower70B | 49.401 | 71.358 | 0.991 | 0.995 | 0.999 | 1.000 | 0.912 | 0.088 | 0.983 | direct |
| CycleL | 0.928 | 14.750 | 0.000 | 0.841 | 0.977 | 0.944 | 0.000 | 0.099 | 0.395 | 1-shot |
| Dubformer | 49.968 | 71.853 | 0.523 | 0.454 | 0.466 | 0.700 | 0.280 | 0.610 | 0.482 | direct |
| IOL_Research | 59.265 | 76.482 | 0.996 | 1.000 | 1.000 | 1.000 | 0.015 | 0.672 | 0.859 | 1-shot |
| ONLINE-A | 53.505 | 72.787 | 0.995 | 1.000 | 0.999 | 1.000 | 0.764 | 0.234 | 0.965 | direct |
| ONLINE-B | 52.012 | 72.139 | 0.998 | 0.998 | 0.999 | 1.000 | 0.923 | 0.077 | **0.988** | direct |
| ONLINE-G | 47.843 | 70.719 | 0.999 | 0.999 | 0.999 | 1.000 | 0.880 | 0.120 | 0.979 | direct |
| ONLINE-W | 56.473 | 74.051 | 0.906 | 0.999 | 0.998 | 1.000 | 0.882 | 0.118 | 0.969 | direct |
| TranssionMT | 54.465 | 74.167 | 0.995 | 0.999 | 0.998 | 1.000 | 0.802 | 0.198 | 0.970 | direct |

Table 86: English→Ukrainian; weakest attack by Avg. win

|  | clean | | adversarial | | | | | | | |
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 13.449 | 35.122 | 0.977 | 0.993 | 0.989 | 0.999 | 0.069 | 0.390 | 0.805 | 1-shot |
| Claude-3 | 55.420 | 75.544 | 0.980 | 0.925 | 0.956 | 0.973 | 0.979 | 0.020 | 0.956 | 0-shot JSON format |
| CommandR-plus | 20.222 | 44.344 | 0.799 | 0.520 | 0.575 | 0.720 | 0.621 | 0.182 | 0.591 | 0-shot JSON format |
| GPT-4 | 42.953 | 65.458 | 1.000 | 1.000 | 1.000 | 1.000 | 0.455 | 0.083 | 0.922 | 1-shot |
| Llama3-70B | 38.608 | 60.739 | 0.996 | 0.897 | 0.929 | 0.985 | 0.988 | 0.004 | 0.928 | 0-shot JSON format |
| IKUN | 31.698 | 55.417 | 0.854 | 0.998 | 0.990 | 1.000 | 0.996 | 0.002 | 0.972 | direct |
| IKUN-C | 25.692 | 49.700 | 0.983 | 0.996 | 0.987 | 0.999 | 0.990 | 0.009 | 0.988 | direct |
| Unbabel-Tower70B | 44.358 | 67.090 | 0.917 | 0.988 | 0.982 | 1.000 | 0.963 | 0.037 | 0.972 | direct |
| AMI | 52.729 | 72.148 | 0.998 | 0.999 | 0.998 | 1.000 | 0.999 | 0.001 | 0.997 | direct |
| CycleL | 10.383 | 29.998 | 0.957 | 0.995 | 0.961 | 1.000 | 0.004 | 0.453 | 0.696 | 0-shot |
| Dubformer | 41.037 | 61.391 | 0.433 | 0.452 | 0.438 | 0.644 | 0.356 | 0.576 | 0.456 | direct |
| IOL_Research | 45.690 | 64.846 | 0.995 | 1.000 | 1.000 | 1.000 | 0.772 | 0.196 | 0.967 | 0-shot |
| ONLINE-A | 55.587 | 73.600 | 0.999 | 1.000 | 1.000 | 1.000 | 0.994 | 0.005 | 0.997 | direct |
| ONLINE-B | 57.116 | 73.904 | 0.996 | 0.999 | 0.999 | 1.000 | 0.999 | 0.000 | **0.999** | direct |
| ONLINE-G | 47.642 | 67.534 | 0.998 | 0.996 | 1.000 | 1.000 | 0.993 | 0.002 | 0.998 | direct |
| ONLINE-W | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | direct |
| TSU-HITs | 8.553 | 28.192 | 0.089 | 0.867 | 0.864 | 0.999 | 0.001 | 0.318 | 0.467 | 1-shot |
| TranssionMT | 57.314 | 74.708 | 0.998 | 0.999 | 0.999 | 1.000 | 0.996 | 0.001 | 0.999 | direct |

Table 87: English→Icelandic; weakest attack by Avg. win

| | clean | | adversarial | | | | | | | |
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 16.547 | 35.168 | 0.991 | 0.015 | 1.000 | 0.701 | 0.827 | 0.159 | 0.596 | 0-shot (en) |
| Claude-3 | 3.943 | 38.065 | 0.491 | 0.514 | 0.998 | 0.502 | 0.187 | 0.244 | 0.598 | 1-shot (non-en) |
| CommandR-plus | 7.728 | 35.127 | 0.444 | 0.513 | 1.000 | 0.635 | 0.515 | 0.439 | 0.621 | 1-shot (en) |
| GPT-4 | 15.472 | 39.233 | 0.498 | 0.513 | 1.000 | 0.628 | 0.812 | 0.184 | 0.709 | 1-shot (en) |
| Llama3-70B | 18.386 | 32.080 | 0.498 | 0.513 | 0.999 | 0.950 | 0.918 | 0.081 | **0.770** | 1-shot (en) |
| IKUN | 1.519 | 28.192 | 0.985 | 0.078 | 0.976 | 0.980 | 0.963 | 0.022 | 0.680 | 1-shot JSON format (non-en) |
| IKUN-C | 5.156 | 23.669 | 0.974 | 0.013 | 0.999 | 0.999 | 0.573 | 0.207 | 0.595 | direct (non-en) |
| Unbabel-Tower70B | 6.585 | 36.271 | 0.490 | 0.507 | 1.000 | 0.796 | 0.821 | 0.175 | 0.731 | 1-shot (en) |
| CycleL | 0.013 | 2.344 | 0.324 | 0.007 | 0.529 | 0.732 | 0.006 | 0.022 | 0.228 | direct (en) |
| DLUT_GTCOM | 0.735 | 30.945 | 0.446 | 0.512 | 1.000 | 1.000 | 0.552 | 0.447 | 0.641 | 1-shot (en) |
| IOL_Research | 16.514 | 39.294 | 0.494 | 0.508 | 1.000 | 0.953 | 0.786 | 0.209 | 0.749 | 1-shot (en) |
| MSLC | 9.124 | 29.066 | 0.998 | 0.016 | 1.000 | 0.998 | 0.902 | 0.087 | 0.704 | 0-shot (en) |
| NTTSU | 0.456 | 32.324 | 0.499 | 0.508 | 1.000 | 0.499 | 0.747 | 0.240 | 0.670 | 1-shot (en) |
| ONLINE-A | 4.688 | 39.838 | 0.499 | 0.509 | 1.000 | 0.973 | 0.873 | 0.126 | 0.728 | 1-shot (en) |
| ONLINE-B | 1.534 | 38.803 | 0.998 | 0.127 | 1.000 | 0.980 | 0.988 | 0.010 | 0.721 | 1-shot JSON format (non-en) |
| ONLINE-G | 2.440 | 33.098 | 0.998 | 0.152 | 0.994 | 0.993 | 0.985 | 0.010 | 0.725 | 0-shot JSON format (non-en) |
| ONLINE-W | 2.803 | 38.856 | 0.494 | 0.514 | 1.000 | 0.519 | 0.761 | 0.169 | 0.685 | 1-shot (en) |
| Team-J | 0.573 | 28.582 | 0.498 | 0.509 | 1.000 | 0.529 | 0.594 | 0.395 | 0.607 | 1-shot (en) |
| UvA-MT | 0.413 | 32.523 | 0.497 | 0.512 | 1.000 | 0.499 | 0.517 | 0.460 | 0.568 | 1-shot (en) |

Table 88: Japanese→Chinese; weakest attack by Avg. win

| | clean | | adversarial | | | | | | | |
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 51.796 | 71.017 | 1.000 | 1.000 | 1.000 | 1.000 | 0.976 | 0.024 | 0.979 | 1-shot JSON format (non-en) |
| Claude-3 | 59.164 | 76.525 | 0.977 | 0.990 | 0.985 | 0.984 | 0.600 | 0.212 | 0.926 | 0-shot (en) |
| CommandR-plus | 52.291 | 71.954 | 0.983 | 0.999 | 0.995 | 1.000 | 0.870 | 0.047 | 0.973 | direct (en) |
| GPT-4 | 50.830 | 71.774 | 0.966 | 0.995 | 0.998 | 1.000 | 0.836 | 0.093 | 0.963 | direct (en) |
| Llama3-70B | 42.691 | 65.406 | 0.991 | 1.000 | 1.000 | 1.000 | 0.878 | 0.037 | 0.974 | direct (en) |
| IKUN | 44.345 | 65.724 | 0.987 | 0.988 | 0.988 | 0.990 | 0.966 | 0.028 | 0.956 | 0-shot JSON format (non-en) |
| IKUN-C | 43.714 | 65.549 | 0.999 | 1.000 | 1.000 | 1.000 | 0.894 | 0.032 | 0.952 | direct (en) |
| Unbabel-Tower70B | 50.091 | 71.296 | 0.988 | 1.000 | 0.998 | 0.998 | 0.973 | 0.027 | 0.980 | 1-shot JSON format (non-en) |
| BJFU-LPT | 23.070 | 42.742 | 0.999 | 1.000 | 1.000 | 1.000 | 0.027 | 0.717 | 0.860 | 0-shot (en) |
| CUNI-Transformer | 51.200 | 70.250 | 0.999 | 1.000 | 1.000 | 1.000 | 0.831 | 0.028 | 0.949 | direct (en) |
| CycleL | 0.110 | 0.686 | 0.000 | 0.651 | 0.028 | 0.000 | 0.000 | 0.000 | 0.097 | 1-shot (en) |
| IOL_Research | 54.964 | 73.144 | 0.983 | 0.995 | 0.999 | 1.000 | 0.722 | 0.175 | 0.938 | direct (en) |
| ONLINE-A | 49.693 | 69.758 | 0.965 | 0.971 | 0.977 | 0.998 | 0.952 | 0.043 | 0.960 | 0-shot JSON format (non-en) |
| ONLINE-B | 47.317 | 68.256 | 0.973 | 0.990 | 0.993 | 0.995 | 0.980 | 0.012 | 0.969 | 1-shot JSON format (non-en) |
| ONLINE-G | 43.649 | 65.989 | 0.999 | 1.000 | 1.000 | 1.000 | 0.454 | 0.208 | 0.887 | direct (en) |
| ONLINE-W | 51.432 | 69.965 | 0.998 | 0.995 | 1.000 | 1.000 | 0.971 | 0.029 | **0.981** | 1-shot JSON format (non-en) |
| TranssionMT | 47.952 | 68.873 | 0.971 | 0.990 | 0.993 | 0.998 | 0.976 | 0.017 | 0.971 | 1-shot JSON format (non-en) |

Table 89: Czech→Ukrainian; weakest attack by Avg. win

| System | clean | | adversarial | | | | | | | Task |
| | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **Avg. win** | |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 42.393 | 60.496 | 1.000 | 1.000 | 1.000 | 1.000 | 0.976 | 0.024 | 0.979 | 1-shot JSON format (non-en) |
| Claude-3 | 47.904 | 71.624 | 0.866 | 0.626 | 0.783 | 0.811 | 0.802 | 0.029 | 0.760 | 0-shot JSON format |
| CommandR-plus | 39.558 | 61.701 | 0.945 | 0.507 | 0.997 | 0.980 | 0.668 | 0.270 | 0.776 | direct (en) |
| GPT-4 | 46.751 | 68.297 | 0.999 | 1.000 | 1.000 | 1.000 | 0.259 | 0.285 | 0.894 | 0-shot (non-en) |
| Llama3-70B | 45.592 | 63.932 | 0.998 | 0.995 | 0.999 | 1.000 | 0.987 | 0.013 | 0.973 | 0-shot JSON format (non-en) |
| NVIDIA-NeMo | 42.710 | 63.846 | 0.916 | 0.742 | 0.921 | 0.967 | 0.886 | 0.105 | 0.870 | direct |
| AIST-AIRC | 27.615 | 53.878 | 0.984 | 0.996 | 1.000 | 1.000 | 0.455 | 0.476 | 0.848 | 1-shot |
| CUNI-DS | 45.865 | 65.698 | 0.985 | 0.907 | 0.930 | 0.985 | 0.953 | 0.038 | 0.933 | 1-shot JSON format |
| CUNI-MH | 57.511 | 75.301 | 0.998 | 1.000 | 1.000 | 1.000 | 0.988 | 0.012 | **0.998** | direct |
| CUNI-NL | 51.442 | 69.699 | 0.761 | 1.000 | 0.999 | 0.999 | 0.988 | 0.005 | 0.964 | direct |
| IKUN | 33.493 | 55.349 | 0.987 | 0.988 | 0.988 | 0.990 | 0.966 | 0.028 | 0.956 | 0-shot JSON format (non-en) |
| IKUN-C | 29.794 | 51.422 | 0.914 | 0.949 | 0.958 | 0.990 | 0.912 | 0.073 | 0.919 | 1-shot JSON format (non-en) |
| Occiglot | 49.361 | 68.297 | 0.951 | 0.919 | 0.862 | 1.000 | 0.958 | 0.029 | 0.922 | direct |
| SCIR-MT | 63.339 | 78.457 | 0.987 | 1.000 | 0.999 | 1.000 | 0.989 | 0.009 | 0.996 | direct |
| Unbabel-Tower70B | 40.216 | 63.839 | 0.988 | 1.000 | 0.998 | 0.998 | 0.973 | 0.027 | 0.980 | 1-shot JSON format (non-en) |
| Yandex | 42.793 | 65.032 | 0.780 | 0.969 | 0.990 | 1.000 | 0.845 | 0.155 | 0.899 | direct |
| AMI | 52.729 | 72.148 | 0.998 | 0.999 | 0.998 | 1.000 | 0.999 | 0.001 | 0.997 | direct |
| BJFU-LPT | 23.070 | 42.742 | 0.999 | 1.000 | 1.000 | 1.000 | 0.027 | 0.717 | 0.860 | 0-shot (en) |
| CUNI-DocTransformer | 58.378 | 75.431 | 0.998 | 0.996 | 0.996 | 1.000 | 0.991 | 0.007 | 0.997 | direct |
| CUNI-GA | 56.400 | 74.149 | 0.987 | 0.968 | 0.936 | 1.000 | 0.968 | 0.031 | 0.952 | direct |
| CUNI-Transformer | 53.800 | 72.199 | 0.987 | 0.968 | 0.936 | 1.000 | 0.968 | 0.031 | 0.952 | direct |
| CycleL | 6.148 | 18.252 | 0.644 | 0.699 | 0.719 | 0.931 | 0.000 | 0.446 | 0.499 | 0-shot |
| CycleL2 | 6.761 | 21.195 | 0.687 | 0.710 | 0.681 | 0.975 | 0.000 | 0.413 | 0.523 | 0-shot |
| DLUT_GTCOM | 0.774 | 36.619 | 0.971 | 0.980 | 1.000 | 1.000 | 0.011 | 0.810 | 0.709 | 1-shot |
| Dubformer | 35.630 | 49.672 | 0.947 | 0.569 | 0.290 | 0.212 | 0.025 | 0.247 | 0.420 | 0-shot |
| HW-TSC | 18.593 | 47.754 | 0.999 | 0.004 | 0.941 | 1.000 | 0.984 | 0.013 | 0.689 | direct |
| IOL_Research | 50.033 | 68.620 | 0.994 | 1.000 | 1.000 | 1.000 | 0.275 | 0.236 | 0.896 | 0-shot (non-en) |
| MSLC | 35.706 | 55.910 | 0.984 | 1.000 | 0.999 | 1.000 | 0.482 | 0.068 | 0.924 | 1-shot |
| NTTSU | 2.525 | 32.728 | 0.499 | 0.508 | 1.000 | 0.499 | 0.747 | 0.240 | 0.670 | 1-shot (en) |
| ONLINE-A | 45.461 | 67.909 | 0.965 | 0.971 | 0.977 | 0.998 | 0.952 | 0.043 | 0.960 | 0-shot JSON format (non-en) |
| ONLINE-B | 41.947 | 65.861 | 0.973 | 0.990 | 0.993 | 0.995 | 0.980 | 0.012 | 0.969 | 1-shot JSON format (non-en) |
| ONLINE-G | 42.300 | 65.329 | 0.926 | 0.766 | 0.962 | 1.000 | 0.906 | 0.092 | 0.899 | direct |
| ONLINE-W | 31.636 | 50.922 | 0.998 | 0.995 | 1.000 | 1.000 | 0.971 | 0.029 | 0.981 | 1-shot JSON format (non-en) |
| TSU-HITs | 20.310 | 41.368 | 0.144 | 0.686 | 0.789 | 0.973 | 0.407 | 0.261 | 0.550 | direct |
| Team-J | 0.494 | 32.453 | 0.998 | 0.998 | 1.000 | 1.000 | 0.002 | 0.994 | 0.714 | 1-shot |
| TranssionMT | 57.720 | 75.513 | 0.971 | 0.990 | 0.993 | 0.998 | 0.976 | 0.017 | 0.971 | 1-shot JSON format (non-en) |
| UvA-MT | 0.746 | 36.751 | 0.950 | 0.006 | 0.933 | 0.999 | 0.918 | 0.080 | 0.676 | direct |

Table 90: Average across all language pairs; weakest attack by Avg. win

## A.2.2 Strongest attacks

| | clean | | adversarial | | | | | | | |
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 50.124 | 69.491 | 0.911 | 0.810 | 0.858 | 0.963 | 0.852 | 0.121 | 0.293 | 0-shot JSON format |
| Claude-3 | 63.945 | 80.516 | 0.010 | 0.006 | 0.001 | 0.005 | 0.000 | 1.000 | **0.836** | direct |
| CommandR-plus | 51.532 | 70.648 | 0.824 | 0.322 | 0.335 | 0.488 | 0.330 | 0.170 | 0.551 | 1-shot JSON format |
| GPT-4 | 58.671 | 76.248 | 0.998 | 0.065 | 0.038 | 0.035 | 0.033 | 0.015 | 0.670 | 1-shot JSON format |
| Llama3-70B | 55.838 | 73.779 | 0.266 | 0.211 | 0.188 | 0.244 | 0.127 | 0.873 | 0.720 | direct |
| NVIDIA-NeMo | 53.441 | 71.047 | 0.351 | 0.367 | 0.406 | 0.991 | 0.343 | 0.011 | 0.354 | 0-shot JSON format |
| CUNI-DS | 45.865 | 65.698 | 0.952 | 0.435 | 0.252 | 0.995 | 0.000 | 0.998 | 0.529 | 0-shot |
| IKUN | 46.017 | 65.324 | 0.973 | 0.884 | 0.909 | 0.968 | 0.936 | 0.055 | 0.260 | 1-shot JSON format |
| IKUN-C | 39.794 | 60.823 | 0.996 | 0.949 | 0.878 | 1.000 | 0.054 | 0.875 | 0.326 | 0-shot |
| Unbabel-Tower70B | 54.457 | 73.925 | 0.968 | 0.655 | 0.671 | 0.720 | 0.679 | 0.055 | 0.381 | 1-shot JSON format |
| Yandex | 42.793 | 65.032 | 0.016 | 0.026 | 0.108 | 0.775 | 0.002 | 0.985 | 0.617 | 1-shot JSON format |
| CycleL | 1.720 | 19.371 | 0.000 | 0.100 | 0.166 | 0.062 | 0.000 | 0.006 | 0.526 | 1-shot JSON format |
| CycleL2 | 0.823 | 15.256 | 0.000 | 0.097 | 0.095 | 0.804 | 0.000 | 0.006 | 0.430 | 1-shot JSON format |
| Dubformer | 0.811 | 2.480 | 0.999 | 0.039 | 0.002 | 0.000 | 0.002 | 0.009 | 0.684 | 0-shot JSON format |
| IOL_Research | 62.421 | 77.519 | 0.965 | 0.655 | 0.589 | 0.990 | 0.463 | 0.535 | 0.407 | direct |
| ONLINE-A | 57.977 | 75.168 | 0.999 | 0.925 | 0.942 | 0.976 | 0.958 | 0.042 | 0.262 | 0-shot JSON format |
| ONLINE-B | 55.403 | 73.776 | 0.976 | 0.890 | 0.916 | 0.945 | 0.923 | 0.050 | 0.268 | 0-shot JSON format |
| ONLINE-G | 53.353 | 74.154 | 0.000 | 0.007 | 0.000 | 0.000 | 0.000 | 0.001 | 0.575 | 1-shot JSON format |
| ONLINE-W | 53.906 | 72.810 | 0.999 | 0.911 | 0.941 | 0.984 | 0.971 | 0.027 | 0.257 | 0-shot JSON format |
| TSU-HITs | 22.052 | 43.818 | 0.000 | 0.067 | 0.054 | 0.640 | 0.000 | 0.624 | 0.564 | 1-shot JSON format |
| TranssionMT | 55.300 | 74.002 | 0.993 | 0.903 | 0.924 | 0.966 | 0.951 | 0.044 | 0.265 | 0-shot JSON format |

Table 91: English→Russian; strongest attack by SAAvg

| | clean | | adversarial | | | | | | | |
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
|---|---|---|---|---|---|---|---|---|---|---|
| Aya23 | 60.528 | 77.596 | 0.998 | 1.000 | 1.000 | 1.000 | 0.092 | 0.846 | 0.301 | 0-shot |
| Claude-3 | 69.372 | 84.126 | 0.024 | 0.119 | 0.173 | 0.234 | 0.024 | 0.974 | 0.767 | direct |
| CommandR-plus | 60.904 | 78.355 | 0.968 | 0.498 | 0.498 | 0.534 | 0.529 | 0.058 | 0.465 | 1-shot JSON format |
| GPT-4 | 70.239 | 84.067 | 0.999 | 0.381 | 0.364 | 0.335 | 0.340 | 0.001 | 0.530 | 1-shot JSON format |
| Llama3-70B | 64.414 | 79.829 | 0.996 | 1.000 | 1.000 | 1.000 | 0.075 | 0.891 | 0.314 | 0-shot |
| NVIDIA-NeMo | 62.179 | 77.817 | 0.681 | 0.678 | 0.710 | 0.967 | 0.665 | 0.001 | 0.313 | 1-shot JSON format |
| AIST-AIRC | 54.511 | 72.781 | 0.251 | 0.191 | 0.162 | 0.846 | 0.084 | 0.013 | 0.429 | 1-shot JSON format |
| CUNI-NL | 51.442 | 69.699 | 0.905 | 0.800 | 0.854 | 0.994 | 0.901 | 0.007 | 0.279 | 1-shot JSON format |
| IKUN | 51.652 | 70.262 | 0.996 | 1.000 | 0.999 | 1.000 | 0.131 | 0.815 | 0.281 | 0-shot |
| IKUN-C | 44.710 | 65.240 | 0.994 | 0.968 | 0.919 | 1.000 | 0.092 | 0.900 | 0.331 | 0-shot |
| Unbabel-Tower70B | 61.008 | 78.193 | 0.989 | 0.633 | 0.651 | 0.651 | 0.654 | 0.001 | 0.400 | 0-shot JSON format |
| CycleL | 20.487 | 44.322 | 0.000 | 0.072 | 0.132 | 0.372 | 0.000 | 0.002 | 0.495 | 1-shot JSON format |
| CycleL2 | 20.487 | 44.322 | 0.000 | 0.072 | 0.132 | 0.372 | 0.000 | 0.002 | 0.495 | 1-shot JSON format |
| Dubformer | 26.213 | 32.808 | 0.360 | 0.039 | 0.023 | 0.010 | 0.009 | 0.621 | **0.824** | 1-shot JSON format |
| IOL_Research | 69.214 | 82.833 | 0.812 | 0.607 | 0.531 | 0.999 | 0.918 | 0.081 | 0.389 | direct |
| MSLC | 41.196 | 64.234 | 0.028 | 0.048 | 0.011 | 0.002 | 0.002 | 0.013 | 0.623 | 1-shot JSON format |
| ONLINE-A | 68.859 | 82.629 | 0.999 | 1.000 | 1.000 | 1.000 | 0.126 | 0.873 | 0.311 | 0-shot |
| ONLINE-B | 54.922 | 74.946 | 0.245 | 1.000 | 0.998 | 1.000 | 0.996 | 0.004 | 0.300 | direct |
| ONLINE-G | 68.624 | 82.302 | 0.999 | 1.000 | 1.000 | 1.000 | 0.246 | 0.745 | 0.293 | 0-shot |
| ONLINE-W | 61.546 | 78.220 | 0.999 | 1.000 | 1.000 | 1.000 | 0.106 | 0.887 | 0.314 | 0-shot |
| TSU-HITs | 29.868 | 49.567 | 0.000 | 0.034 | 0.042 | 0.264 | 0.000 | 0.028 | 0.540 | 0-shot JSON format |
| TranssionMT | 54.873 | 74.941 | 0.242 | 1.000 | 0.998 | 1.000 | 0.996 | 0.004 | 0.300 | direct |

Table 92: English→German; strongest attack by SAAvg

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 19.085 | 40.614 | 0.099 | 0.000 | 0.848 | 1.000 | 0.507 | 0.491 | 0.449 | direct |
| Claude-3 | 1.919 | 53.543 | 0.007 | 0.000 | 0.005 | 0.013 | 0.000 | 1.000 | **0.830** | direct |
| CommandR-plus | 14.366 | 43.986 | 0.404 | 0.002 | 0.350 | 0.936 | 0.026 | 0.953 | 0.548 | 0-shot |
| GPT-4 | 17.514 | 54.097 | 0.999 | 0.001 | 0.002 | 0.001 | 0.009 | 0.048 | 0.694 | 1-shot JSON format |
| Llama3-70B | 27.898 | 43.181 | 0.424 | 0.001 | 0.785 | 0.813 | 0.683 | 0.317 | 0.418 | direct |
| NVIDIA-NeMo | 2.076 | 35.694 | 0.459 | 0.005 | 0.487 | 0.742 | 0.741 | 0.204 | 0.504 | direct |
| AIST-AIRC | 0.719 | 34.974 | 0.013 | 0.000 | 0.012 | 0.016 | 0.002 | 0.015 | 0.568 | 1-shot JSON format |
| IKUN | 13.311 | 31.025 | 0.154 | 0.001 | 0.950 | 1.000 | 0.662 | 0.334 | 0.384 | direct |
| IKUN-C | 2.249 | 26.016 | 0.854 | 0.001 | 0.416 | 0.923 | 0.005 | 0.994 | 0.473 | 0-shot |
| Unbabel-Tower70B | 8.143 | 41.692 | 0.936 | 0.002 | 0.854 | 1.000 | 0.006 | 0.989 | 0.413 | 0-shot |
| CycleL | 0.041 | 3.364 | 0.013 | 0.004 | 0.326 | 0.982 | 0.000 | 0.431 | 0.444 | 0-shot |
| DLUT_GTCOM | 0.813 | 42.293 | 0.043 | 0.000 | 0.062 | 0.159 | 0.023 | 0.103 | 0.586 | 0-shot JSON format |
| IOL_Research | 19.182 | 51.107 | 0.938 | 0.004 | 0.933 | 0.979 | 0.020 | 0.957 | 0.381 | 0-shot |
| NTTSU | 4.594 | 33.132 | 0.780 | 0.004 | 0.343 | 0.184 | 0.184 | 0.267 | 0.595 | 0-shot |
| ONLINE-A | 1.220 | 44.459 | 0.372 | 0.000 | 0.379 | 0.372 | 0.367 | 0.004 | 0.453 | 0-shot JSON format |
| ONLINE-B | 1.015 | 44.589 | 0.933 | 0.001 | 0.918 | 1.000 | 0.006 | 0.974 | 0.403 | 0-shot |
| ONLINE-G | 3.339 | 45.429 | 0.242 | 0.020 | 0.267 | 0.367 | 0.267 | 0.024 | 0.477 | 0-shot JSON format |
| ONLINE-W | 4.871 | 34.170 | 0.104 | 0.000 | 0.192 | 0.098 | 0.100 | 0.005 | 0.525 | 0-shot JSON format |
| Team-J | 0.416 | 36.323 | 0.987 | 0.002 | 0.433 | 0.499 | 0.494 | 0.039 | 0.489 | 0-shot JSON format |
| UvA-MT | 1.159 | 43.238 | 0.951 | 0.002 | 0.020 | 0.040 | 0.048 | 0.070 | 0.686 | 0-shot JSON format |

Table 93: English→Japanese; strongest attack by SAAvg

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 44.375 | 63.672 | 0.976 | 0.640 | 0.685 | 0.754 | 0.715 | 0.054 | 0.376 | 0-shot JSON format |
| Claude-3 | 60.166 | 76.954 | 0.010 | 0.009 | 0.004 | 0.024 | 0.000 | 1.000 | **0.843** | direct |
| CommandR-plus | 39.996 | 61.592 | 0.764 | 0.305 | 0.299 | 0.328 | 0.220 | 0.376 | 0.588 | direct |
| GPT-4 | 50.565 | 69.608 | 0.996 | 0.035 | 0.016 | 0.015 | 0.022 | 0.028 | 0.682 | 1-shot JSON format |
| Llama3-70B | 51.601 | 69.311 | 0.082 | 0.022 | 0.020 | 0.026 | 0.018 | 0.979 | 0.829 | direct |
| NVIDIA-NeMo | 47.354 | 66.582 | 0.022 | 0.076 | 0.130 | 0.928 | 0.001 | 0.905 | 0.552 | 1-shot JSON format |
| IKUN | 40.887 | 60.362 | 0.263 | 0.308 | 0.383 | 0.993 | 0.246 | 0.264 | 0.405 | 1-shot JSON format |
| IKUN-C | 35.290 | 56.369 | 0.968 | 0.546 | 0.271 | 0.994 | 0.001 | 0.998 | 0.517 | 0-shot |
| Unbabel-Tower70B | 56.242 | 74.129 | 0.998 | 0.940 | 0.769 | 0.999 | 0.000 | 1.000 | 0.399 | 0-shot |
| CycleL | 0.268 | 12.822 | 0.000 | 0.372 | 0.417 | 0.558 | 0.000 | 0.422 | 0.441 | 0-shot |
| IOL_Research | 53.133 | 70.132 | 0.979 | 0.750 | 0.727 | 0.998 | 0.627 | 0.362 | 0.334 | direct |
| ONLINE-A | 59.021 | 74.613 | 0.999 | 1.000 | 1.000 | 1.000 | 0.001 | 0.958 | 0.331 | 0-shot |
| ONLINE-B | 56.473 | 71.907 | 0.998 | 0.989 | 0.979 | 0.973 | 0.000 | 0.995 | 0.354 | 0-shot |
| ONLINE-G | 55.704 | 72.554 | 0.616 | 0.503 | 0.475 | 0.660 | 0.519 | 0.214 | 0.392 | 0-shot JSON format |
| ONLINE-W | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.016 | 0.574 | 0-shot |
| TranssionMT | 56.588 | 73.267 | 0.999 | 0.990 | 0.980 | 0.971 | 0.000 | 0.995 | 0.355 | 0-shot |

Table 94: English→Hindi; strongest attack by SAAvg

| | clean | | | | | adversarial | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 71.590 | 83.455 | 0.919 | 0.335 | 0.341 | 0.346 | 0.326 | 0.061 | 0.544 | 0-shot JSON format |
| Claude-3 | 77.382 | 88.287 | 0.009 | 0.024 | 0.061 | 0.005 | 0.000 | 1.000 | **0.823** | direct |
| CommandR-plus | 69.366 | 82.843 | 0.939 | 0.146 | 0.149 | 0.160 | 0.135 | 0.087 | 0.641 | 0-shot JSON format |
| GPT-4 | 76.485 | 86.879 | 0.250 | 0.356 | 0.341 | 0.962 | 0.190 | 0.810 | 0.625 | direct |
| Llama3-70B | 75.659 | 85.899 | 0.066 | 0.054 | 0.054 | 0.058 | 0.024 | 0.976 | 0.821 | direct |
| NVIDIA-NeMo | 71.684 | 83.575 | 0.996 | 1.000 | 0.996 | 1.000 | 0.300 | 0.695 | 0.316 | 0-shot |
| IKUN | 56.366 | 73.524 | 0.842 | 0.902 | 0.820 | 1.000 | 0.988 | 0.009 | 0.281 | direct |
| IKUN-C | 52.543 | 70.275 | 0.999 | 0.925 | 0.812 | 1.000 | 0.248 | 0.747 | 0.350 | 0-shot |
| Occiglot | 49.361 | 68.297 | 0.679 | 0.469 | 0.299 | 0.966 | 0.004 | 0.983 | 0.569 | 0-shot |
| Unbabel-Tower70B | 58.762 | 76.431 | 0.989 | 0.798 | 0.829 | 0.854 | 0.847 | 0.001 | 0.327 | 0-shot JSON format |
| CycleL | 32.147 | 51.642 | 0.000 | 0.058 | 0.127 | 0.649 | 0.001 | 0.001 | 0.521 | 0-shot JSON format |
| Dubformer | 60.120 | 79.825 | 0.237 | 0.078 | 0.078 | 0.295 | 0.012 | 0.520 | 0.721 | 1-shot JSON format |
| IOL_Research | 76.839 | 86.496 | 0.933 | 0.887 | 0.862 | 1.000 | 0.917 | 0.082 | 0.294 | direct |
| MSLC | 56.800 | 74.431 | 0.887 | 0.509 | 0.498 | 0.621 | 0.534 | 0.010 | 0.450 | 0-shot JSON format |
| ONLINE-A | 74.616 | 85.820 | 0.999 | 1.000 | 1.000 | 1.000 | 0.274 | 0.681 | 0.296 | 0-shot |
| ONLINE-B | 72.932 | 83.788 | 0.996 | 1.000 | 1.000 | 1.000 | 0.272 | 0.716 | 0.305 | 0-shot |
| ONLINE-G | 76.360 | 86.243 | 0.999 | 1.000 | 1.000 | 1.000 | 0.275 | 0.721 | 0.303 | 0-shot |
| ONLINE-W | 58.478 | 74.701 | 0.999 | 1.000 | 1.000 | 1.000 | 0.089 | 0.903 | 0.304 | 0-shot |
| TSU-HITs | 24.907 | 50.317 | 0.009 | 0.091 | 0.102 | 0.344 | 0.005 | 0.020 | 0.511 | 0-shot JSON format |
| TranssionMT | 73.144 | 85.551 | 0.999 | 1.000 | 1.000 | 1.000 | 0.277 | 0.678 | 0.296 | 0-shot |

Table 95: English→Spanish; strongest attack by SAAvg

| | clean | | | | | adversarial | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 57.243 | 74.550 | 0.928 | 0.879 | 0.896 | 0.985 | 0.927 | 0.060 | 0.277 | 0-shot JSON format |
| Claude-3 | 66.823 | 81.945 | 0.006 | 0.032 | 0.040 | 0.006 | 0.000 | 1.000 | **0.834** | direct |
| CommandR-plus | 54.377 | 73.408 | 0.729 | 0.296 | 0.267 | 0.458 | 0.286 | 0.425 | 0.594 | direct |
| GPT-4 | 64.985 | 79.784 | 0.999 | 0.073 | 0.043 | 0.004 | 0.054 | 0.087 | 0.684 | 1-shot JSON format |
| Llama3-70B | 61.753 | 77.069 | 0.780 | 0.783 | 0.732 | 0.944 | 0.192 | 0.782 | 0.403 | 0-shot |
| NVIDIA-NeMo | 55.940 | 72.507 | 0.890 | 0.448 | 0.453 | 0.494 | 0.457 | 0.049 | 0.486 | 1-shot JSON format |
| CUNI-MH | 57.511 | 75.301 | 0.996 | 0.999 | 0.993 | 1.000 | 0.273 | 0.714 | 0.259 | 0-shot |
| IKUN | 45.469 | 65.478 | 0.310 | 0.357 | 0.382 | 0.996 | 0.269 | 0.148 | 0.359 | 1-shot JSON format |
| IKUN-C | 37.968 | 58.621 | 0.995 | 0.919 | 0.733 | 1.000 | 0.062 | 0.936 | 0.345 | 0-shot |
| SCIR-MT | 63.339 | 78.457 | 0.987 | 1.000 | 0.999 | 1.000 | 0.073 | 0.907 | 0.302 | 0-shot |
| Unbabel-Tower70B | 51.206 | 71.180 | 0.967 | 0.854 | 0.869 | 0.901 | 0.887 | 0.017 | 0.282 | 1-shot JSON format |
| CUNI-DocTransformer | 58.378 | 75.431 | 0.998 | 0.444 | 0.449 | 0.441 | 0.492 | 0.062 | 0.496 | 0-shot JSON format |
| CUNI-GA | 56.400 | 74.149 | 0.942 | 0.174 | 0.149 | 0.116 | 0.198 | 0.087 | 0.639 | 0-shot JSON format |
| CUNI-Transformer | 56.400 | 74.149 | 0.942 | 0.174 | 0.149 | 0.116 | 0.198 | 0.087 | 0.639 | 0-shot JSON format |
| CycleL | 1.469 | 17.798 | 0.000 | 0.078 | 0.078 | 0.621 | 0.000 | 0.002 | 0.464 | 1-shot JSON format |
| CycleL2 | 5.734 | 24.422 | 0.000 | 0.099 | 0.126 | 0.787 | 0.000 | 0.015 | 0.429 | 0-shot JSON format |
| IOL_Research | 64.617 | 78.908 | 0.879 | 0.764 | 0.643 | 1.000 | 0.816 | 0.176 | 0.341 | direct |
| ONLINE-A | 63.853 | 79.054 | 0.999 | 1.000 | 1.000 | 1.000 | 0.359 | 0.592 | 0.259 | 0-shot |
| ONLINE-B | 59.851 | 76.425 | 0.998 | 1.000 | 0.991 | 1.000 | 0.301 | 0.667 | 0.262 | 0-shot |
| ONLINE-G | 63.404 | 78.063 | 0.998 | 0.903 | 0.961 | 0.995 | 0.993 | 0.005 | 0.255 | 0-shot JSON format |
| ONLINE-W | 55.114 | 73.094 | 0.999 | 1.000 | 1.000 | 1.000 | 0.078 | 0.920 | 0.293 | 0-shot |
| TSU-HITs | 16.169 | 34.946 | 0.007 | 0.100 | 0.177 | 0.902 | 0.006 | 0.007 | 0.409 | 0-shot JSON format |
| TranssionMT | 62.123 | 78.598 | 0.999 | 1.000 | 0.990 | 0.993 | 0.307 | 0.659 | 0.265 | 0-shot |

Table 96: English→Czech; strongest attack by SAAvg

|  | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 31.789 | 44.156 | 0.929 | 0.005 | 0.468 | 0.983 | 0.006 | 0.984 | 0.458 | 0-shot |
| Claude-3 | 6.160 | 52.796 | 0.006 | 0.001 | 0.077 | 0.262 | 0.000 | 0.998 | **0.772** | direct |
| CommandR-plus | 12.165 | 44.248 | 0.486 | 0.004 | 0.494 | 0.810 | 0.037 | 0.930 | 0.523 | 0-shot |
| GPT-4 | 11.698 | 49.684 | 0.999 | 0.002 | 0.006 | 0.002 | 0.021 | 0.032 | 0.691 | 1-shot JSON format |
| Llama3-70B | 14.886 | 45.139 | 0.973 | 0.031 | 0.931 | 0.990 | 0.091 | 0.841 | 0.335 | 0-shot |
| NVIDIA-NeMo | 1.315 | 35.577 | 0.035 | 0.002 | 0.009 | 0.015 | 0.001 | 0.207 | 0.621 | 1-shot JSON format |
| IKUN | 2.722 | 34.863 | 0.996 | 0.005 | 0.995 | 1.000 | 0.075 | 0.783 | 0.308 | 0-shot |
| IKUN-C | 4.261 | 29.898 | 0.991 | 0.005 | 0.646 | 0.999 | 0.033 | 0.951 | 0.390 | 0-shot |
| Unbabel-Tower70B | 2.125 | 40.668 | 0.998 | 0.054 | 0.493 | 0.551 | 0.545 | 0.027 | 0.445 | 1-shot JSON format |
| CycleL | 0.057 | 3.676 | 0.001 | 0.000 | 0.002 | 0.330 | 0.000 | 0.000 | 0.524 | 0-shot JSON format |
| CycleL2 | 0.000 | 0.779 | 0.004 | 0.000 | 0.044 | 0.632 | 0.000 | 0.004 | 0.480 | 0-shot JSON format |
| HW-TSC | 18.593 | 47.754 | 0.321 | 0.015 | 0.170 | 0.776 | 0.130 | 0.076 | 0.444 | 0-shot JSON format |
| IOL_Research | 28.529 | 54.058 | 0.998 | 0.009 | 1.000 | 1.000 | 0.033 | 0.958 | 0.355 | 0-shot |
| ONLINE-A | 11.048 | 49.271 | 0.999 | 0.005 | 0.996 | 1.000 | 0.034 | 0.862 | 0.339 | 0-shot |
| ONLINE-B | 2.844 | 45.939 | 0.999 | 0.005 | 0.999 | 1.000 | 0.047 | 0.922 | 0.380 | 0-shot |
| ONLINE-G | 2.939 | 42.534 | 0.998 | 0.010 | 0.805 | 1.000 | 0.015 | 0.984 | 0.386 | 0-shot |
| ONLINE-W | 3.376 | 44.271 | 0.847 | 0.018 | 0.233 | 0.370 | 0.241 | 0.045 | 0.550 | 0-shot JSON format |
| UvA-MT | 0.668 | 34.492 | 0.015 | 0.000 | 0.000 | 0.000 | 0.000 | 0.136 | 0.607 | 1-shot JSON format |

Table 97: English→Chinese; strongest attack by SAAvg

|  | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 49.791 | 70.611 | 0.998 | 0.767 | 0.698 | 0.999 | 0.635 | 0.365 | 0.331 | direct |
| Claude-3 | 62.653 | 79.565 | 0.024 | 0.007 | 0.004 | 0.012 | 0.009 | 0.984 | **0.838** | direct |
| CommandR-plus | 52.187 | 72.206 | 0.288 | 0.360 | 0.400 | 0.972 | 0.351 | 0.628 | 0.580 | 1-shot JSON format |
| GPT-4 | 54.848 | 74.440 | 0.999 | 0.043 | 0.009 | 0.006 | 0.086 | 0.278 | 0.720 | 1-shot JSON format |
| Llama3-70B | 49.780 | 70.822 | 0.251 | 0.166 | 0.143 | 0.211 | 0.127 | 0.873 | 0.738 | direct |
| NVIDIA-NeMo | 47.690 | 67.971 | 0.428 | 0.492 | 0.519 | 0.974 | 0.493 | 0.166 | 0.368 | 0-shot JSON format |
| IKUN | 34.437 | 58.673 | 0.940 | 0.860 | 0.882 | 0.968 | 0.905 | 0.048 | 0.254 | 1-shot JSON format |
| IKUN-C | 36.359 | 59.479 | 0.996 | 0.925 | 0.860 | 1.000 | 0.012 | 0.983 | 0.338 | 0-shot |
| Unbabel-Tower70B | 49.401 | 71.358 | 0.982 | 0.941 | 0.873 | 0.999 | 0.023 | 0.956 | 0.348 | 0-shot |
| CycleL | 0.928 | 14.750 | 0.000 | 0.118 | 0.168 | 0.982 | 0.000 | 0.000 | 0.390 | 0-shot JSON format |
| Dubformer | 49.968 | 71.853 | 0.384 | 0.037 | 0.027 | 0.208 | 0.039 | 0.655 | 0.761 | 0-shot JSON format |
| IOL_Research | 59.265 | 76.482 | 0.130 | 0.106 | 0.228 | 0.816 | 0.105 | 0.070 | 0.429 | 1-shot JSON format |
| ONLINE-A | 53.505 | 72.787 | 0.966 | 0.906 | 0.930 | 0.998 | 0.941 | 0.056 | 0.263 | 0-shot JSON format |
| ONLINE-B | 52.012 | 72.139 | 0.987 | 0.905 | 0.935 | 0.991 | 0.955 | 0.038 | 0.255 | 0-shot JSON format |
| ONLINE-G | 47.843 | 70.719 | 0.148 | 0.108 | 0.211 | 0.753 | 0.119 | 0.103 | 0.446 | 0-shot JSON format |
| ONLINE-W | 56.473 | 74.051 | 0.996 | 0.917 | 0.929 | 0.991 | 0.945 | 0.049 | 0.261 | 1-shot JSON format |
| TranssionMT | 54.465 | 74.167 | 0.993 | 0.913 | 0.944 | 0.999 | 0.967 | 0.029 | 0.256 | 1-shot JSON format |

Table 98: English→Ukrainian; strongest attack by SAAvg

|  | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 13.449 | 35.122 | 0.805 | 0.448 | 0.460 | 0.616 | 0.409 | 0.109 | 0.400 | 1-shot JSON format |
| Claude-3 | 55.420 | 75.544 | 0.021 | 0.035 | 0.044 | 0.280 | 0.001 | 0.976 | **0.838** | 1-shot JSON format |
| CommandR-plus | 20.222 | 44.344 | 0.337 | 0.196 | 0.180 | 0.370 | 0.195 | 0.605 | 0.667 | direct |
| GPT-4 | 42.953 | 65.458 | 0.991 | 0.078 | 0.035 | 0.009 | 0.005 | 0.032 | 0.677 | 1-shot JSON format |
| Llama3-70B | 38.608 | 60.739 | 0.076 | 0.069 | 0.067 | 0.078 | 0.042 | 0.958 | 0.800 | direct |
| IKUN | 31.698 | 55.417 | 0.778 | 0.718 | 0.767 | 0.942 | 0.797 | 0.055 | 0.275 | 0-shot JSON format |
| IKUN-C | 25.692 | 49.700 | 0.455 | 0.420 | 0.480 | 0.836 | 0.426 | 0.111 | 0.372 | 1-shot JSON format |
| Unbabel-Tower70B | 44.358 | 67.090 | 0.996 | 0.938 | 0.843 | 1.000 | 0.212 | 0.760 | 0.318 | 0-shot |
| AMI | 52.729 | 72.148 | 0.837 | 0.812 | 0.853 | 0.994 | 0.863 | 0.048 | 0.285 | 1-shot JSON format |
| CycleL | 10.383 | 29.998 | 0.000 | 0.072 | 0.120 | 0.428 | 0.000 | 0.000 | 0.485 | 0-shot JSON format |
| Dubformer | 41.037 | 61.391 | 0.732 | 0.070 | 0.034 | 0.045 | 0.004 | 0.255 | 0.723 | 1-shot JSON format |
| IOL_Research | 45.690 | 64.846 | 0.996 | 0.465 | 0.409 | 0.994 | 0.815 | 0.181 | 0.398 | direct |
| ONLINE-A | 55.587 | 73.600 | 0.994 | 0.920 | 0.944 | 0.995 | 0.991 | 0.001 | 0.251 | 0-shot JSON format |
| ONLINE-B | 57.116 | 73.904 | 0.891 | 0.892 | 0.934 | 0.993 | 0.946 | 0.005 | 0.261 | 1-shot JSON format |
| ONLINE-G | 47.642 | 67.534 | 0.999 | 0.998 | 0.989 | 1.000 | 0.165 | 0.603 | 0.252 | 0-shot |
| ONLINE-W | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.023 | 0.575 | 0-shot |
| TSU-HITs | 8.553 | 28.192 | 0.004 | 0.054 | 0.175 | 0.947 | 0.000 | 0.821 | 0.567 | 0-shot JSON format |
| TranssionMT | 57.314 | 74.708 | 0.881 | 0.897 | 0.934 | 0.993 | 0.945 | 0.005 | 0.264 | 1-shot JSON format |

Table 99: English→Icelandic; strongest attack by SAAvg

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 16.547 | 35.168 | 0.934 | 0.035 | 0.441 | 0.432 | 0.393 | 0.592 | 0.396 | 0-shot JSON format (non-en) |
| Claude-3 | 3.943 | 38.065 | 0.658 | 0.042 | 0.215 | 0.196 | 0.130 | 0.756 | 0.540 | 1-shot JSON format (non-en) |
| CommandR-plus | 7.728 | 35.127 | 0.553 | 0.011 | 0.716 | 0.788 | 0.272 | 0.480 | 0.345 | direct (non-en) |
| GPT-4 | 15.472 | 39.233 | 0.993 | 0.024 | 0.044 | 0.029 | 0.032 | 0.961 | 0.559 | 1-shot JSON format (non-en) |
| Llama3-70B | 18.386 | 32.080 | 0.481 | 0.020 | 0.491 | 0.498 | 0.010 | 0.732 | 0.466 | direct (non-en) |
| IKUN | 1.519 | 28.192 | 0.015 | 0.000 | 0.083 | 0.578 | 0.000 | 0.176 | 0.504 | 1-shot JSON format (en) |
| IKUN-C | 5.156 | 23.669 | 0.015 | 0.000 | 0.083 | 0.578 | 0.000 | 0.176 | 0.504 | 1-shot JSON format (en) |
| Unbabel-Tower70B | 6.585 | 36.271 | 0.493 | 0.504 | 1.000 | 0.525 | 0.118 | 0.574 | 0.294 | 1-shot (non-en) |
| CycleL | 0.013 | 2.344 | 0.002 | 0.001 | 0.073 | 0.291 | 0.000 | 0.031 | 0.523 | 0-shot JSON format (en) |
| DLUT_GTCOM | 0.735 | 30.945 | 0.978 | 0.007 | 0.061 | 0.071 | 0.068 | 0.320 | 0.462 | 1-shot JSON format (non-en) |
| IOL_Research | 16.514 | 39.294 | 0.892 | 0.092 | 0.760 | 0.799 | 0.743 | 0.106 | 0.224 | 0-shot JSON format (non-en) |
| MSLC | 9.124 | 29.066 | 0.998 | 0.020 | 0.017 | 0.000 | 0.042 | 0.946 | 0.565 | 1-shot JSON format (non-en) |
| NTTSU | 0.456 | 32.324 | 0.861 | 0.012 | 0.020 | 0.007 | 0.046 | 0.912 | 0.582 | 1-shot JSON format (non-en) |
| ONLINE-A | 4.688 | 39.838 | 0.980 | 0.022 | 0.792 | 0.002 | 0.853 | 0.110 | 0.331 | 1-shot JSON format (en) |
| ONLINE-B | 1.534 | 38.803 | 0.590 | 0.005 | 0.565 | 0.058 | 0.213 | 0.464 | 0.464 | 0-shot JSON format (en) |
| ONLINE-G | 2.440 | 33.098 | 0.975 | 0.022 | 0.944 | 0.000 | 0.841 | 0.108 | 0.310 | 1-shot JSON format (en) |
| ONLINE-W | 2.803 | 38.856 | 0.174 | 0.010 | 0.291 | 0.949 | 0.132 | 0.257 | 0.405 | 1-shot JSON format (non-en) |
| Team-J | 0.573 | 28.582 | 0.020 | 0.021 | 0.136 | 0.114 | 0.177 | 0.721 | 0.638 | 0-shot JSON format (non-en) |
| UvA-MT | 0.413 | 32.523 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.479 | **0.669** | 1-shot JSON format (non-en) |

Table 100: Japanese→Chinese; strongest attack by SAAvg

| | clean | | adversarial | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| System | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | Task |
| Aya23 | 51.796 | 71.017 | 0.501 | 1.000 | 1.000 | 1.000 | 0.157 | 0.733 | 0.177 | 1-shot (non-en) |
| Claude-3 | 59.164 | 76.525 | 0.939 | 0.482 | 0.061 | 0.044 | 0.034 | 0.954 | 0.495 | 1-shot JSON format (non-en) |
| CommandR-plus | 52.291 | 71.954 | 0.472 | 0.722 | 0.563 | 0.520 | 0.370 | 0.607 | 0.336 | direct (non-en) |
| GPT-4 | 50.830 | 71.774 | 0.995 | 0.504 | 0.022 | 0.005 | 0.024 | 0.934 | 0.492 | 1-shot JSON format (non-en) |
| Llama3-70B | 42.691 | 65.406 | 0.267 | 0.424 | 0.355 | 0.360 | 0.034 | 0.835 | 0.493 | direct (non-en) |
| IKUN | 44.345 | 65.724 | 1.000 | 0.919 | 0.042 | 0.000 | 0.743 | 0.025 | 0.295 | 1-shot JSON format (en) |
| IKUN-C | 43.714 | 65.549 | 0.593 | 0.547 | 0.527 | 0.806 | 0.240 | 0.287 | 0.267 | 1-shot JSON format (en) |
| Unbabel-Tower70B | 50.091 | 71.296 | 0.498 | 1.000 | 1.000 | 0.999 | 0.177 | 0.667 | 0.168 | 1-shot (non-en) |
| BJFU-LPT | 23.070 | 42.742 | 0.213 | 0.174 | 0.120 | 0.377 | 0.000 | 0.721 | 0.550 | 1-shot JSON format (non-en) |
| CUNI-Transformer | 51.200 | 70.250 | 0.117 | 0.073 | 0.000 | 0.000 | 0.000 | 0.254 | **0.646** | 1-shot JSON format (non-en) |
| CycleL | 0.110 | 0.686 | 0.000 | 0.157 | 0.001 | 0.000 | 0.000 | 0.147 | 0.571 | direct (non-en) |
| IOL_Research | 54.964 | 73.144 | 0.147 | 0.147 | 0.249 | 0.765 | 0.115 | 0.430 | 0.452 | 1-shot JSON format (non-en) |
| ONLINE-A | 49.693 | 69.758 | 0.498 | 1.000 | 1.000 | 0.955 | 0.002 | 0.830 | 0.197 | 1-shot (en) |
| ONLINE-B | 47.317 | 68.256 | 0.499 | 1.000 | 1.000 | 0.563 | 0.005 | 0.791 | 0.247 | 1-shot (en) |
| ONLINE-G | 43.649 | 65.989 | 0.006 | 0.006 | 0.207 | 0.005 | 0.005 | 0.002 | 0.540 | 0-shot JSON format (en) |
| ONLINE-W | 51.432 | 69.965 | 0.499 | 1.000 | 1.000 | 0.996 | 0.121 | 0.742 | 0.179 | 1-shot (en) |
| TranssionMT | 47.952 | 68.873 | 0.499 | 1.000 | 1.000 | 0.974 | 0.001 | 0.895 | 0.204 | 1-shot (en) |

Table 101: Czech→Ukrainian; strongest attack by SAAvg

| System | clean | | adversarial | | | | | | | | Task |
| | BLEU | chrF | QM | BW | CW | LID | Transl | Ans | **SAAvg** | | Task |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Aya23 | 42.393 | 60.496 | 0.934 | 0.035 | 0.441 | 0.432 | 0.393 | 0.592 | 0.396 | | 0-shot JSON format (non-en) |
| Claude-3 | 47.904 | 71.624 | 0.012 | 0.026 | 0.046 | 0.071 | 0.004 | 0.995 | **0.818** | | direct |
| CommandR-plus | 39.558 | 61.701 | 0.751 | 0.261 | 0.378 | 0.581 | 0.379 | 0.205 | 0.520 | | 1-shot JSON format |
| GPT-4 | 46.751 | 68.297 | 0.998 | 0.096 | 0.076 | 0.059 | 0.078 | 0.058 | 0.663 | | 1-shot JSON format |
| Llama3-70B | 45.592 | 63.932 | 0.411 | 0.242 | 0.419 | 0.447 | 0.394 | 0.606 | 0.582 | | direct |
| NVIDIA-NeMo | 42.710 | 63.846 | 0.579 | 0.470 | 0.437 | 0.799 | 0.264 | 0.150 | 0.354 | | 1-shot JSON format |
| AIST-AIRC | 27.615 | 53.878 | 0.132 | 0.095 | 0.087 | 0.431 | 0.043 | 0.014 | 0.499 | | 1-shot JSON format |
| CUNI-DS | 45.865 | 65.698 | 0.952 | 0.435 | 0.252 | 0.995 | 0.000 | 0.998 | 0.529 | | 0-shot |
| CUNI-MH | 57.511 | 75.301 | 0.996 | 0.999 | 0.993 | 1.000 | 0.273 | 0.714 | 0.259 | | 0-shot |
| CUNI-NL | 51.442 | 69.699 | 0.905 | 0.800 | 0.854 | 0.994 | 0.901 | 0.007 | 0.279 | | 1-shot JSON format |
| IKUN | 33.493 | 55.349 | 0.507 | 0.460 | 0.062 | 0.289 | 0.371 | 0.100 | 0.399 | | 1-shot JSON format (en) |
| IKUN-C | 29.794 | 51.422 | 0.304 | 0.273 | 0.305 | 0.692 | 0.120 | 0.232 | 0.385 | | 1-shot JSON format (en) |
| Occiglot | 49.361 | 68.297 | 0.679 | 0.469 | 0.299 | 0.966 | 0.004 | 0.983 | 0.569 | | 0-shot |
| SCIR-MT | 63.339 | 78.457 | 0.987 | 1.000 | 0.999 | 1.000 | 0.073 | 0.907 | 0.302 | | 0-shot |
| Unbabel-Tower70B | 40.216 | 63.839 | 0.987 | 0.751 | 0.884 | 1.000 | 0.147 | 0.830 | 0.343 | | 0-shot |
| Yandex | 42.793 | 65.032 | 0.016 | 0.026 | 0.108 | 0.775 | 0.002 | 0.985 | 0.617 | | 1-shot JSON format |
| AMI | 52.729 | 72.148 | 0.837 | 0.812 | 0.853 | 0.994 | 0.863 | 0.048 | 0.285 | | 1-shot JSON format |
| BJFU-LPT | 23.070 | 42.742 | 0.213 | 0.174 | 0.120 | 0.377 | 0.000 | 0.721 | 0.550 | | 1-shot JSON format (non-en) |
| CUNI-DocTransformer | 58.378 | 75.431 | 0.998 | 0.444 | 0.449 | 0.441 | 0.492 | 0.062 | 0.496 | | 0-shot JSON format |
| CUNI-GA | 56.400 | 74.149 | 0.942 | 0.174 | 0.149 | 0.116 | 0.198 | 0.087 | 0.639 | | 0-shot JSON format |
| CUNI-Transformer | 53.800 | 72.199 | 0.117 | 0.073 | 0.000 | 0.000 | 0.000 | 0.254 | 0.646 | | 1-shot JSON format (non-en) |
| CycleL | 6.148 | 18.252 | 0.000 | 0.157 | 0.001 | 0.000 | 0.000 | 0.147 | 0.571 | | direct (non-en) |
| CycleL2 | 6.761 | 21.195 | 0.001 | 0.066 | 0.099 | 0.659 | 0.000 | 0.007 | 0.457 | | 1-shot JSON format |
| DLUT_GTCOM | 0.774 | 36.619 | 0.043 | 0.000 | 0.062 | 0.159 | 0.023 | 0.103 | 0.586 | | 0-shot JSON format |
| Dubformer | 35.630 | 49.672 | 0.586 | 0.054 | 0.030 | 0.088 | 0.015 | 0.376 | 0.737 | | 1-shot JSON format |
| HW-TSC | 18.593 | 47.754 | 0.321 | 0.015 | 0.170 | 0.776 | 0.130 | 0.076 | 0.444 | | 0-shot JSON format |
| IOL_Research | 50.033 | 68.620 | 0.147 | 0.147 | 0.249 | 0.765 | 0.115 | 0.430 | 0.452 | | 1-shot JSON format (non-en) |
| MSLC | 35.706 | 55.910 | 0.998 | 0.020 | 0.017 | 0.000 | 0.042 | 0.946 | 0.565 | | 1-shot JSON format (non-en) |
| NTTSU | 2.525 | 32.728 | 0.780 | 0.004 | 0.343 | 0.184 | 0.184 | 0.267 | 0.595 | | 0-shot |
| ONLINE-A | 45.461 | 67.909 | 0.990 | 0.779 | 0.993 | 1.000 | 0.158 | 0.683 | 0.289 | | 0-shot |
| ONLINE-B | 41.947 | 65.861 | 0.990 | 0.777 | 0.987 | 0.997 | 0.212 | 0.696 | 0.297 | | 0-shot |
| ONLINE-G | 42.300 | 65.329 | 0.054 | 0.110 | 0.244 | 0.961 | 0.020 | 0.545 | 0.465 | | 1-shot JSON format (non-en) |
| ONLINE-W | 31.636 | 50.922 | 0.174 | 0.010 | 0.291 | 0.949 | 0.132 | 0.257 | 0.405 | | 1-shot JSON format (non-en) |
| TSU-HITs | 20.310 | 41.368 | 0.004 | 0.063 | 0.104 | 0.519 | 0.002 | 0.189 | 0.517 | | 0-shot JSON format |
| Team-J | 0.494 | 32.453 | 0.020 | 0.021 | 0.136 | 0.114 | 0.177 | 0.721 | 0.638 | | 0-shot JSON format (non-en) |
| TranssionMT | 57.720 | 75.513 | 0.998 | 0.998 | 0.995 | 0.995 | 0.263 | 0.621 | 0.269 | | 0-shot |
| UvA-MT | 0.746 | 36.751 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.479 | 0.669 | | 1-shot JSON format (non-en) |

Table 102: Average across all language pairs; strongest attack by SAAvg

# Killing Two Flies with One Stone: An Attempt to Break LLMs Using English→Icelandic Idioms and Proper Names

**Bjarki Ármannsson, Hinrik Hafsteinsson, Atli Jasonarson, Steinþór Steingrímsson**

The Árni Magnússon Institute for Icelandic Studies

Reykjavík, Iceland

bjarki.armannsson,hinrik.hafsteinsson,atli.jasonarson,
steinthor.steingrimsson@arnastofnun.is

## Abstract

This paper presents the submission of the Árni Magnússon Institute's team to the WMT24 test suite subtask, focusing on idiomatic expressions and proper names for the English→Icelandic translation direction.

Intuitively and empirically, idioms and proper names are known to be a significant challenge for modern translation models. We create two different test suites. The first evaluates the competency of MT systems in translating common English idiomatic expressions, as well as testing whether systems can distinguish between those expressions and the same phrases when used in a literal context. The second test suite consists of place names that should be translated into their Icelandic exonyms (and correctly inflected) and pairs of Icelandic names that share a surface form between the male and female variants, so that incorrect translations impact meaning as well as readability.

The scores reported are relatively low, especially for idiomatic expressions and place names, and indicate considerable room for improvement.

## 1 Introduction

Significant advances in machine translation have in recent years been achieved by integrating Large Language Models (LLMs) into neural translation systems (Xu et al., 2024). Careful analysis, however, has repeatedly shown that despite recording higher scores and producing text with greater fluency compared to previous state-of-the-art neural systems, the translations produced by LLMs are still far from perfect and can include significant biases, misinformation and hallucinations (Hendy et al., 2023), half-hidden in the impressive-looking output. Aiming to expose "weaknesses and serious flaws" of these systems that might otherwise get "hidden in the average", the theme of this year's WMT test suite subtask is "Help us break LLMs",

with organizers asking for custom test sets focusing on phenomena that can provide specific challenges for LLM-based systems. This paper describes the efforts of the Árni Magnússon Institute's team to pick holes in otherwise seemingly fluent English→Icelandic translations.

We experiment with two main features we believe should prove particularly challenging for English→Icelandic LLM-based machine translation systems; idiomatic expressions and proper names. More specifically, we focus on:

- **Idiomatic expressions in English and their literal counterparts:** In the first of our two test sets, we investigate idiomatic expressions in English which do not directly translate to Icelandic. Where possible, we also include 'inverse' examples of usage in a literal form (as in "Are you supposed to **chew the fat** from steak?" or "Blow into the balloon and **tie the knot** without letting the air out.") to give an idea of the translation models' ability to correctly switch between literal and non-literal translations of the same phrase, depending on context.

- **Proper names:** In our second test set, we also consider names of both people and places. We carefully curate a list of city and area names in English that should be translated to their common Icelandic names (and correctly inflected). We then include a list of simple sentences containing both Icelandic and English given names. For the Icelandic names, we observe whether they are correctly inflected in the Icelandic text (which impacts not only the text's readability, but also its meaning). Common English names, meanwhile, are included to test that the models don't 'translate' them to Icelandic – i.e. alter them in some unintended way.

We release our test suites and evaluation code

for others to build on and to allow for further comparison between future models in these categories.[1]

## 2 Related Work

Idiomatic expressions (and multi-word expressions (MWEs) in general) have been the focus of much work in the field of machine translation in recent years and the construction of impressive idiom datasets has been carried out for many other languages and language pairs. See e.g. Stap et al. (2024) for English↔German and Russian→English, Tang (2022) for Chinese→English, Fadaee et al. (2018) for English↔German and Haagsma et al. (2020) and Adewumi et al. (2022) for monolingual datasets of English idiomatic expressions.

Macketanz et al. (2022) include idioms among many other interesting linguistic phenomena in their dataset for English↔German and English→Russian and we took some inspiration from their work when deciding on our scoring format. Halldórsson et al. (2022) list Icelandic idioms with English equivalents, this dataset is described and discussed in more detail in Steingrímsson et al. (2024). We are not aware of any dataset for the English→Icelandic translation direction published previous to our work.

In recent years, the emergence of LLMs has led to work investigating how they handle the translations of idioms and MWEs compared with previous models. Raunak et al. (2023), using measures of 'literalness', find that GPT models produce less literal translations between English and German, Chinese, and Russian than previous neural models, a difference most pronounced in the case of idiomatic expressions. Finally, Shwartz (2021) provide an accessible overview of the kinds of problems posed by MWEs for language models in general.

## 3 Methodology

### 3.1 Idiomatic Expressions

We make use of the set of potential idiomatic expressions defined in the PIE Corpus (Adewumi et al., 2022) and, for each expression we use, extract two examples of usage from the NewsCrawl corpus of WMT 2023 (Kocmi et al., 2023)[2] For our purposes, we narrow the PIE set down from 591 expressions to 199. Our aim was to remove those we deem too rare or obscure to be truly relevant (e.g. *horses for courses* or *monkey's uncle*) for model comparison, expressions which directly (or more or less directly) translate between English and Icelandic (e.g. "open the floodgates" has an Icelandic equivalent, "opna flóðgáttirnar") and those for which we find no example usage in the NewsCrawl corpus. We make the number of expressions an even 200 by adding one that was not in the PIE corpus: "kill two birds with one stone".

Each of the 400 example sentences - two examples for each of the 200 selected idioms - is then manually reviewed to make sure that the relevant idiomatic expression is being used in the intended, non-literal sense. To further increase the difficulty of the task (though still keeping it trivial for fluent human speakers of Icelandic and English), we also try and test the models on their ability to translate the words in these expressions literally when appropriate. We include 223 additional example sentences, for as many expressions as we were able, where the expression is used in a literal sense (or in a few cases, very slightly altered to try and exploit the likelihood bias of LLMs).[3] These examples are largely taken from the NewsCrawl corpus but synthetic in some cases.

To evaluate the models' performance, we construct two 'positive' sets of Icelandic word forms or multiword expressions for each idiom. One set contains words that we would expect to find in a literal translation of the phrase, the other words or phrases that could be expected to appear in a suitable, non-literal translation of the idiomatic expression. In many cases, we also construct 'negative' sets of words that instantly lead to a sentence being marked incorrect, such as the Icelandic words for "weather" or "pink" for idiomatic translations of the phrases "under the weather" and "in the pink". An Icelandic translation of an example sentence in English is marked as correct if it contains any of the words in the set of 'positive' words (in any lexical form) **and** it contains none of the words in the set of 'negative' words (see Table 1).

---

[3] Early inspiration for this project was provided by the one idiomatic expression we added from outside the PIE corpus: "kill two birds with one stone". We noticed a prominent online translation service correctly translated this to the equivalent Icelandic phrase, "slá tvær flugur í einu höggi" (lit. *hit two flies in one strike*), whereas a phrase like "He killed two birds yesterday" would be wrongly translated as "Hann drap tvær flugur í gær" (lit. *He killed two flies yesterday*), exposing a weakness particular to neural and LLM-based systems. Indeed, four of the systems tested here made this particular mistake.

| Source sentence | Possible translations | Evaluation |
|---|---|---|
| Why Fleabag is **in the pink**! | Fleabag er *í góðum málum*! <br> Fleabag er **í bleiku**! | ✓(Positive match) <br> ✗ (Negative match) |
| The young woman **in the pink** continued to throw punches [...] | Unga konan **í góðum málum** lét hnefana tala áfram [...] <br> Unga konan **í bleiku fötunum** lét hnefana tala áfram [...] | ✗ (No positive match) <br><br> ✓(Positive match) |

Table 1: Fabricated example translations into Icelandic of two English sentences containing the phrase "in the pink", both from our test suite. The first English sentence uses the phrase in an idiomatic sense (meaning *in good health* or *in a state of well-being*) and the second seemingly in a literal sense (the full context, not included in the table, is: "before another wades in"). For the idiomatic sentence, we automatically mark it as correct if a match is found from a list of possible Icelandic translations (here the phrase "í góðum málum") *and* no match is found from a list of negative matches (here the lexeme "bleikur", meaning *pink*). For the literal sentence, meanwhile, some form of "bleikur" is required for a correct marking.

During our manual evaluation, we further whittled down our set as we decided a few sentences we had decided to include were actually not testing what they were meant to test (as some were, for instance, more linguistically acceptable when translated directly into Icelandic than we originally felt during the construction of our test set). We removed a total of 25 sentences this way, bringing the total of 'idiomatic' examples in our set to 397 and the total of 'literal' examples to 201. Note that although these examples were removed after we received their translations from the tested models, they are not included in our scoring.

### 3.2 Proper Names

For our testing of place names, we construct our own list of 52 names of cities and areas that we argue would be highly unusual not to translate into their Icelandic names.[4]

As a reference when collecting our place names, we make use of Wikipedia's list of Icelandic exonyms.[5] We use only a small subset of that list, however. Aiming to err on the side of caution, we try to include only place names where native speakers would be in more or less complete agreement to apply their Icelandic names rather than the ones used in English (e.g. the name "Kaupmannahöfn" for Copenhagen is invariably used, whereas "Lundúnir" for London is very rare and mostly used in a colourful or joking manner.[6] We

also leave out cases where the differences between the names used in English and in Icelandic only have to do with pronunciation or minor differences in spelling. In addition to the Icelandic exonyms we select, we make sure to also include several examples of cities where the local name is the one more generally used and English speakers use a rarer (typically French-derived) name (e.g. "München" rather than the English "Munich").

We then construct example sentences in English where each of our selected place names is used in four different contexts, corresponding to each of the four grammatical cases in Icelandic. (The exceptions are "Paris" and "Berlin", which are only tested in the genitive as they are practically the same as in English in the other three cases.) Our motivation is that due to the richer morphology of Icelandic, an accurate translation model needs to be able to map the same lexical form in English to several different forms in Icelandic, depending on the context (and this particular mapping is perhaps a problem better suited to older models than state-of-the-art LLM-based ones).

We try to avoid the possibility that our sentences will be translated into Icelandic in a way that is generally correct but uses a different syntactic structure or wording than we anticipate, which would lend itself to the use of a different grammatical case than the one we intend to test for. We do this by keeping our example sentences short and simple and choose case-governing words and prepositions carefully to maximize the probability of a particular translation in Icelandic (e.g. the sentence "The flight

---

[4]There exist context-dependent exceptions to this, of course, such as the name of a sport club or particular institution from a certain city. Our example sentences, however, refer clearly to the cities in general.

[5]https://en.wikipedia.org/wiki/Icelandic_exonyms

[6]One anonymous reviewer asked whether we had considered incorporating a native speaker survey in order to validate our choices. While the suggestion is certainly a good one, it is beyond the scope of this particular work.

| Source sentence | Possible translations | Evaluation |
|---|---|---|
| **Helgi** dreams of flying | **Helgi** dreymir um að fljúga | ✗ (Ungrammatical) |
| | **Helga** dreymir um að fljúga | ✓ |
| **Helga** dreams of flying | **Helga** dreymir um að fljúga | ✗ (Refers to Helgi, not Helga) |
| | **Helgu** dreymir um að fljúga | ✓ |

Table 2: Examples of possible translations of the phrase "dreams of flying". In Icelandic, the verb "dreyma" (*to dream*) takes a subject argument in the accusative case, which requires a translation system to alter the form of the given name in the English text. Left unaltered in Icelandic, the male name "Helgi" renders the sentence ungrammatical and the female name "Helga" would cause the reader to interpret the sentence to refer to a male called Helgi instead.

from Tórshavn to Gothenburg was delayed until the morning" should almost certainly be translated using the prepositions "frá" and "til" for "from" and "to", governing the dative and genitive cases respectively.)

Given names, both in Icelandic and English, constitute the final part of our test suite. As in the case of the place names, we construct simple sentences in English containing Icelandic names and meant to test for each of the four grammatical cases. For this task, we chose a specific subset of common names in Icelandic: male-female pairings that take the weak inflection, e.g. "Helgi"-"Helga" and "Gunni"-"Gunna", where the male name has the ending -"i" in the nominative case but -"a" in oblique cases and the female name has the ending -"a" in the nominative but -"u" in the oblique cases (and possibly also a u-umlaut as in "Svala" → "Svölu").

These name pairs, of which we select 45 from the Database of Icelandic Morphology (Bjarnadóttir et al., 2019),[7] are chosen as they seem to present a particular challenge for translation systems compared to names that take the strong declension. In constructing our test suite, we found that available models seemed to perform at random when asked to translate sentences containing these names in different cases, presumably due to the ambiguity of the lexical forms ending in -"a", which can be a male name in an oblique case or a female name in the nominative. As oblique case nominals are a distinct and common feature of the Icelandic language (Thráinsson, 2007), this problem is highly relevant in terms of correctly relaying the meaning of the sentence (see Table 2).

## 4 Results

All submissions were scored using automatic metrics we constructed. Furthermore, we manually

reviewed around 150 randomly selected examples in the case of the idioms (around 100 'idiomatic' examples and around 50 'literal' examples for each submitted system). The authors reviewed the translations themselves, manually changing the scores given by our automatic method (using the 'positive' and 'negative' keywords discussed in 3.1) if they deemed it wrong.

The translations of our proper names suite was only carried out with naive automatic methods. The translations were lemmatized using a lemmatizer for Icelandic (Ingólfsdóttir et al., 2019) and compared with a reference of which Icelandic lemmas should appear in the translation and in which grammatical form (being able to look up lemmas is especially useful for the given names, since the male and female names share surface forms).

We show the results of our manual evaluation in Table 3 and the results of automatic metrics for our idioms test suite in Table 4. For our names test suite, we show the results of our automatic metrics in Table 5. Our scripts for running the automatic evaluations and the manually reviewed examples are released along with our test sets.

### 4.1 Scores for Idiomatic Expressions

Our results show a wide range of performance across different models. The best overall accuracy on the idioms test suite is achieved by Claude 3.5, with Unbabel-Tower70B a close second, as indicated both by our automatic and manual evaluation. Claude 3.5 is also the highest-scoring submission when we only consider translations of expressions used in an idiomatic sense, both according to our automatic metrics and the manual review, and Unbabel-Tower70B the clear runner-up.

When considering the literal translations in isolation, however, the overall two best models are narrowly 'beaten' by a few models that score considerably lower overall. According to our automatic

---

[7] https://bin.arnastofnun.is/DMII/

| System | Total Idioms | Total Literals | Idiom Accuracy | Literal Accuracy | Total Accuracy |
|---|---|---|---|---|---|
| AMI | 100 | 65 | 0.29 | 0.892308 | 0.527273 |
| Aya23 | 93 | 49 | 0.0537634 | 0.122449 | 0.0774648 |
| Claude-3.5 | 96 | 56 | **0.75** | 0.857143 | **0.789474** |
| CommandR-plus | 93 | 47 | 0.0967742 | 0.382979 | 0.192857 |
| CycleL | 92 | 49 | 0 | 0.102041 | 0.035461 |
| Dubformer | 91 | 53 | 0.340659 | 0.603774 | 0.4375 |
| GPT-4 | 93 | 48 | 0.430108 | 0.833333 | 0.567376 |
| IKUN-C | 95 | 52 | 0.494737 | 0.75 | 0.585034 |
| IKUN | 95 | 51 | 0.526316 | 0.607843 | 0.554795 |
| IOL_Research | 92 | 47 | 0.434783 | 0.702128 | 0.52518 |
| Llama3-70B | 93 | 50 | 0.268817 | 0.62 | 0.391608 |
| ONLINE-A | 188 | 107 | 0.265957 | 0.859813 | 0.481356 |
| ONLINE-B | 102 | 69 | 0.22549 | 0.898551 | 0.497076 |
| ONLINE-G | 97 | 66 | 0.185567 | 0.80303 | 0.435583 |
| TranssionMT | 76 | 50 | 0.223684 | **0.92** | 0.5 |
| TSU-HITs | 92 | 48 | 0.0434783 | 0.104167 | 0.0642857 |
| Unbabel-Tower70B | 95 | 57 | 0.631579 | 0.877193 | 0.723684 |

Table 3: Results of manual evaluation of system performance on our idioms test suite. We randomly split up the translations of the test suite into segments of around 100 'idiomatic' example translations and around 50 'literal' example translations (see 'Total' columns). The highest scores in each column are in bold. The authors reviewed the translations themselves and the reviewed examples, along with our grading, can be found at `https://github.com/stofnun-arna-magnussonar/idioms_names_test_suite/idioms/human_evaluation`.

| System name | Total score | Correct idiomatics | CI ratio | Correct literals | CL ratio |
|---|---|---|---|---|---|
| AMI | 0.447236 | 83 | 0.21 | 184 | **0.9** |
| Aya23 | 0.169179 | 39 | 0.1 | 62 | 0.3 |
| Claude-3.5 | **0.654941** | 216 | **0.55** | 175 | 0.86 |
| CommandR-plus | 0.293132 | 66 | 0.17 | 109 | 0.53 |
| CycleL | 0.108878 | 22 | 0.06 | 43 | 0.21 |
| Dubformer | 0.427136 | 112 | 0.28 | 143 | 0.7 |
| GPT-4 | 0.547739 | 161 | 0.41 | 166 | 0.81 |
| IKUN-C | 0.480737 | 141 | 0.36 | 146 | 0.72 |
| IKUN | 0.509213 | 161 | 0.41 | 143 | 0.7 |
| IOL_Research | 0.482412 | 133 | 0.34 | 155 | 0.76 |
| Llama3-70B | 0.417085 | 99 | 0.25 | 150 | 0.74 |
| ONLINE-A | 0.442211 | 86 | 0.22 | 178 | 0.87 |
| ONLINE-B | 0.447236 | 85 | 0.22 | 182 | 0.89 |
| ONLINE-G | 0.413735 | 71 | 0.18 | 176 | 0.86 |
| TranssionMT | 0.448911 | 86 | 0.22 | 182 | 0.89 |
| TSU-HITs | 0.112228 | 24 | 0.06 | 43 | 0.21 |
| Unbabel-Tower70B | 0.60804 | 195 | 0.5 | 168 | 0.82 |

Table 4: Results of automatic evaluation of system performance on our idioms test suite. We show the overall score for each system but also consider separately the percentage of idiomatic text examples marked as correct and the percentage of literals marked correct, to try and give an overview of the relationship between the two. Highest scores in each column are in bold. Our scripts for running automatic evaluation can be found at `https://github.com/stofnun-arna-magnussonar/idioms_names_test_suite/idioms`.

| System name | Total score | Total city score | Total people score |
| --- | --- | --- | --- |
| AMI | **0.5399** | **0.4705** | 0.5861 |
| Aya23 | 0.3838 | 0.0432 | 0.6103 |
| Claude-3.5 | 0.5091 | 0.4591 | 0.5423 |
| CommandR-plus | 0.3339 | 0.1205 | 0.4758 |
| CycleL | 0.0 | 0.0 | 0.0 |
| Dubformer | 0.4383 | 0.3614 | 0.4894 |
| GPT-4 | 0.5109 | 0.2773 | **0.6662** |
| IKUN-C | 0.4691 | 0.2727 | 0.5997 |
| IKUN | 0.4846 | 0.2886 | 0.6148 |
| IOL_Research | 0.4773 | 0.2205 | 0.648 |
| Llama3-70B | 0.4138 | 0.3227 | 0.4743 |
| ONLINE-A | 0.5345 | 0.4659 | 0.5801 |
| ONLINE-B | 0.5109 | 0.4273 | 0.5665 |
| ONLINE-G | 0.4065 | 0.3614 | 0.4366 |
| TranssionMT | 0.5082 | 0.4227 | 0.565 |
| TSU-HITs | 0.147 | 0.0932 | 0.1828 |
| Unbabel-Tower70B | 0.5254 | 0.4114 | 0.6012 |

Table 5: Results of automatic evaluation of system performance on our names test suite, given as a proportion of properly scored city names 'Total city score', properly scored given names 'Total people score' and overall 'Total score'. Highest scores in each column are in bold. Our scripts for running automatic evaluation can be found at https://github.com/stofnun-arna-magnussonar/idioms_names_test_suite/names. (Note that the zeroes for CycleL's submission are not a mistake, this submission performed poorly and our scoring strategy is not particularly forgiving.)

metrics, our own submission (AMI) scores highest in that category, only slightly ahead of ONLINE-B and TranssionMT. These three also come out on top in the manual evaluation, with TranssionMT recording the highest score (a superb 0.92) and ONLINE-B and AMI following in second and third.

This discrepancy between performance in translating phrases in an idiomatic context and a literal context is very interesting - these three models all scored under 0.3 in idiomatic accuracy, which suggests that for some models, proficiency in effectively translating text in a literal sense comes at a cost to their ability to handle more metaphorical text. The best-performing models overall, however, were seemingly able to maneuver quite effectively between both use cases. Models, perhaps predictably, generally score higher when translating literal usage than when translating idioms.

### 4.2 Scores for Proper Names

In terms of the proper names suite, place names prove to be much more difficult for the submitted models than people's names. It is the submission by our own team which narrowly tops the list overall, ahead of ONLINE-A and Unbabel. The AMI submission also ranks highest when place names are considered in isolation, although it still gets fewer than half of all names correct. For given names, GPT-4 scores highest.

For this part of our test set, we report no manual evaluation. A cursory glance at the output, however, shows that our naive automatic scoring method still leaves quite a bit to be desired. A problem with testing for specific grammatical forms in each case is that the correct form can change depending on the sentence structure. As discussed in 3.2, we tried to control for this by keeping test sentences brief and unambiguous. Even so, we find there are examples of different phrasings than we expected in some translation outputs that call for a different grammatical form of a name than our scoring mechanism supposes, but can still be considered a decent translation.

This especially applies to the sentence form: X "cares for" Y. We assumed a correct translation into Icelandic would be: X "þykir vænt um" Y, where X would take the dative case and Y the accusative. The submitted systems, however, had many different ideas on how best to phrase this system, not all of them completely wrong.

We therefore recognize that our scoring system needs to be fine-tuned but nevertheless believe the very low scores are mainly a reflection of the difficulty of this task.

## 5 Conclusions and Future Work

Scores on both sets are relatively low, indicating that these particular categories continue to pose some problems for even state-of-the-art translation models and that there is considerable room for improvement.

Future work can explore further comparison of performance and fine-tuning of our automatic scoring methods. Given time, we could also have investigated whether more manual evaluation, ideally using more evaluators, would have resulted in different scores.

We also note that our test suite can be adapted with relative ease into other languages and hope that this allows for further work on other language directions.

## Limitations

There are several judgment calls to be made when working with our chosen categories and many of the decisions we made in terms of selecting items to be translated, defining automatic metrics for 'right' and 'wrong' translations and manual evaluation can be argued for or against. We are aware that the choices we make could be indicative of potential biases of the authors and that a different team, perhaps with a different demographic makeup, might well have constructed the test set and evaluated the translations in a different way.

These necessary choices are perhaps most apparent in terms of our idioms set. Evaluation of linguistic acceptability of translations and correspondence of idiomatic phrases between languages is based on our intuition and we are aware that fluent speakers of English and Icelandic may disagree on some decisions. Another point to consider is the degree to which we want our test set to be prescriptive - as a simple search on the Internet can prove, there are multiple usages of common English idioms directly translated into Icelandic, e.g. on social media (Hilmisdóttir et al., 2023). Determining at what point to say this usage is no longer 'incorrect' is an interesting question of ethics and philosophy of language.

As for our set of proper names, there exists some speaker variation in how and when place names are

translated into Icelandic, although we have tried to limit our set to fairly uncontroversial choices (see discussion in 3.2). The requirement of not translating English names into Icelandic is less cut and dried, as it may be appropriate for a machine translation model in some cases, e.g. in literary text or the discussion of royal or historical figures. It can also be noted that some of our English names are, in fact, given names in Iceland. This should not affect our results, however, as we allow for the inflection of a final -"a" into -"u" in female names like "Pamela" and in other cases, 'non-Icelandic' names typically remain completely unchanged in all grammatical cases.

## Acknowledgments

## References

Tosin Adewumi, Roshanak Vadoodi, Aparajita Tripathy, Konstantina Nikolaido, Foteini Liwicki, and Marcus Liwicki. 2022. Potential idiomatic expression (PIE)-English: Corpus for classes of idioms. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 689–696, Marseille, France. European Language Resources Association.

Kristín Bjarnadóttir, Kristín Ingibjörg Hlynsdóttir, and Steinþór Steingrímsson. 2019. DIM: The Database of Icelandic Morphology. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, Turku, Finland.

Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2018. Examining the tip of the iceberg: A data set for idiom translation. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Hessel Haagsma, Johan Bos, and Malvina Nissim. 2020. MAGPIE: A large corpus of potentially idiomatic expressions. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 279–287, Marseille, France. European Language Resources Association.

Björn Halldórsson, Árni Davíð Magnússon, Finnur Ágúst Ingimundarson, Einar Freyr Sigurðsson, Steinþór Steingrímsson, Halldóra Jónsdóttir, and Þórdís Úlfarsdóttir. 2022. Idiomatic expressions (Icelandic and English) 22.09. CLARIN-IS.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita,

Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? A comprehensive evaluation.

Helga Hilmisdóttir, Martina Huhtamäki, and Susanna Karlsson. 2023. Pragmatic borrowing from English. *Nordic Journal of Linguistics*, 46(3):255–256.

Svanhvít Lilja Ingólfsdóttir, Hrafn Loftsson, Jón Friðrik Daðason, and Kristín Bjarnadóttir. 2019. Nefnir: A high accuracy lemmatizer for Icelandic. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, pages 310–315, Turku, Finland. Linköping University Electronic Press.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, and Hans Uszkoreit. 2022. A linguistically motivated test suite to semi-automatically evaluate German–English machine translation output. In *Proceedings of the Language Resources and Evaluation Conference*, pages 936–947, Marseille, France. European Language Resources Association.

Vikas Raunak, Arul Menezes, Matt Post, and Hany Hassan. 2023. Do GPTs produce less literal translations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1041–1050, Toronto, Canada. Association for Computational Linguistics.

Vered Shwartz. 2021. A long hard look at MWEs in the age of language models. In *Proceedings of the 17th Workshop on Multiword Expressions (MWE 2021)*, page 1, Online. Association for Computational Linguistics.

David Stap, Eva Hasler, Bill Byrne, Christof Monz, and Ke Tran. 2024. The fine-tuning paradox: Boosting translation quality without sacrificing LLM abilities. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6189–6206, Bangkok, Thailand. Association for Computational Linguistics.

Steinþór Steingrímsson, Einar Freyr Sigurðsson, and Björn Halldórsson. 2024. Evaluating Capabilities of MT Systems in Translating Idiomatic Expressions Using a Specialized Dataset. In *Proceedings of CLARIN annual conference 2024, October 15-17, 2024*.

Kenan Tang. 2022. Petci: A parallel English translation dataset of Chinese idioms.

Höskuldur Thráinsson. 2007. *The Syntax of Icelandic*. Cambridge University Press, Cambridge.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024. A paradigm shift in machine translation: Boosting translation performance of large language models. In *The Twelfth International Conference on Learning Representations*.

# METAMETRICS-MT: Tuning Meta-Metrics for Machine Translation via Human Preference Calibration

**David Anugraha**[†1], **Garry Kuwanto**[†2], **Lucky Susanto**[3],
**Derry Tanti Wijaya**[‡,2,3], **Genta Indra Winata**[‡*4]

[1]University of Toronto    [2]Boston University
[3]Monash Indonesia    [4]Capital One
david.anugraha@cs.toronto.edu, {gkuwanto,wijaya}@bu.edu,
lucky.susanto@monash.edu, genta.winata@capitalone.com

## Abstract

We present METAMETRICS-MT, an innovative metric designed to evaluate machine translation (MT) tasks by aligning closely with human preferences through Bayesian optimization with Gaussian Processes. METAMETRICS-MT enhances existing MT metrics by optimizing their correlation with human judgments. Our experiments on the WMT24 metric shared task dataset demonstrate that METAMETRICS-MT outperforms all existing baselines, setting a new benchmark for state-of-the-art performance in the reference-based setting. Furthermore, it achieves comparable results to leading metrics in the reference-free setting, offering greater efficiency.

## 1 Introduction

Evaluating machine translation (MT) tasks is inherently complex, as no single metric can universally apply to all scenarios. A metric that performs well for one task may not be suitable for another, and its effectiveness can vary significantly depending on the specific language pairs involved. Therefore, relying solely on a single metric is often inadequate. To ensure the usefulness of automatic metrics, it is crucial to align them with human annotations (Winata et al., 2024b). To achieve a more comprehensive evaluation, benchmarks typically incorporate multiple metrics, such as lexical-based and semantic-based metrics. However, the correlation between these metrics can be skewed due to variations in the models used and the training data employed for evaluation. For instance, BERTScore (Zhang et al., 2019) uses contextual embeddings from pre-trained transformers to assess performance, with different models excelling in specific language pairs. In contrast, neural-based metrics like BLEURT (Sellam et al., 2020), COMET (Rei et al., 2020), and CometKiwi (Rei

et al., 2022) employ distinct methodologies and training datasets. These differences can affect each metric's alignment with human judgments and their reliability across language pairs. Some metrics, like XCOMET-Ensemble (Guerreiro et al., 2023), demand high memory (at least 80GB), prompting efforts to predict LLM performance using smaller models (Anugraha et al., 2024).

In this paper, we propose METAMETRICS-MT, a MT metric inspired by METAMETRICS (Winata et al., 2024a). This meta-metric is crafted to align more closely with human preferences through the use of Bayesian optimization with Gaussian Processes (GP). By systematically integrating multiple existing metrics, METAMETRICS-MT achieves state-of-the-art performance for reference-based metrics and shows a strong correlation with human scores for reference-free metrics in the WMT24 metric shared task (Freitag et al., 2024). Through the strategic combination of metrics with assigned weights, METAMETRICS-MT aims to be as competitive as, if not superior to, any individual metric. Our contributions include the following:

- We present METAMETRICS-MT in reference-based and reference-free settings, offering flexibility for various MT scenarios. Our reference-based model sets the state-of-the-art for the WMT24 task. We publicly release the code for easy usability.[1]

- We demonstrate that the METAMETRICS-MT metric is easily adjustable to meet the human preference.

- We show that METAMETRICS-MT is compact and efficient, capable of running on a commercial GPU with 40GB of memory, whereas a comparable metric like XCOMET-Ensemble requires significantly higher memory with at least 80GB.

---

*The work was conducted outside Capital One. †These authors contributed equally. ‡Senior authors.

[1]The code is available at https://github.com/meta-metrics/metametrics.

## 2 Methodology

### 2.1 METAMETRICS-MT

METAMETRICS-MT is designed to leverage multiple metrics for assessing MT tasks, with each metric being adjusted by specific weights to optimize performance. The idea of utilizing multiple metrics is to combine scores from multiple metrics regardless of the metric types. Formally, let $\theta_1, \theta_2, \ldots, \theta_N$ represent $N$ distinct metric functions with $\hat{y}_1, \ldots, \hat{y}_N$ as their respective performance on a translation task. We define $\Phi$ to compute a scalar meta-metric score of $\hat{y}_{MM}$ using $\hat{y}_1, \ldots, \hat{y}_N$. Overall, we define $\theta_{MM}$ as a meta-metric function where $\hat{y}_{MM}$ is computed as follows:

$$\hat{y}_i = \theta_i(x), \tag{1}$$

$$\hat{y}_{MM} = \theta_{MM}(x) = \Phi(\hat{y}_1, \cdots, \hat{y}_N). \tag{2}$$

Our objective is to calibrate a metric function, $\theta_{MM}$, to maximize the correlation $\rho(\hat{y}_{MM}, \gamma)$, where $\rho$ is a correlation measure and $\gamma$ represents human assessment scores, which include any scores provided by human evaluators. Each metric operates within a specific range, defined by minimum and maximum values. However, some metrics, particularly those based on neural networks, may fall outside this range. To ensure consistency, we normalize these metrics to a common scale from 0 to 1, where 0 signifies poor translation performance and 1 signifies perfect translation performance. In this process, given an original score $y_i$ for a given metric, $\tilde{y}_i$ represents the normalized score. For more details on pre-processing, please refer to Section A of the Appendix.

In this case, we use GP to model the function $\Phi$ and it can be breakdown into a weighted sum as follows:

$$y_{MM} = \alpha_1 \tilde{y}_1 + \alpha_2 \tilde{y}_2 + \ldots + \alpha_N \tilde{y}_N, \tag{3}$$

where $\alpha_1, \alpha_2, \ldots, \alpha_N$ are the corresponding weights assigned to each metric, constrained to the interval $[0, 1]$. Our objective is to determine the best set of weights for $\alpha_1, \alpha_2, \ldots, \alpha_N$, which maximizes $\rho(y_{MM}, \gamma)$. Notice that $y_{MM}$ lies in the interval of $[0, N]$, so normalizing $y_{MM}$ back to $[0, 1]$ is unnecessary as linear scaling does not affect the correlation coefficient for correlation function $\rho$.

The advantage of METAMETRICS-MT is its flexibility and adaptability across tasks and domains. By integrating metrics that strongly correlate with human judgments for specific tasks, we

| Metric | clipping | normalization | inversion | weight |
|---|---|---|---|---|
| **Reference-based** (METAMETRICS-MT) | | | | |
| MetricX-23-XXL | [0,25] | ✓ | ✓ | 1.0000 |
| COMET | [0,1] | ✓ | ✗ | 0.2055 |
| XCOMET-XL | [0,1] | ✓ | ✗ | 0.2733 |
| **Reference-free** (METAMETRICS-MT-QE) | | | | |
| MetricX-23-XXL-QE | [0,25] | ✓ | ✓ | 0.9905 |
| CometKiwi (QE) | [0,1] | ✓ | ✗ | 0.1267 |
| CometKiwi-XL (QE) | [0,1] | ✓ | ✗ | 0.0584 |

Table 1: Metric configuration for METAMETRICS-MT. Metrics not listed in the table have been assigned a weight of zero.

can create a composite metric that improves overall alignment with human evaluations.

### 2.2 Bayesian Optimization

We optimize the weights for each metric using Bayesian optimization with GP as the surrogate model. Bayesian optimization is particularly useful in this context because it efficiently explores and exploits the parameter space when the objective function is expensive to evaluate. By constructing a probabilistic model of the objective function, Bayesian optimization balances exploring new areas with exploiting known promising regions, making it effective even when evaluations are costly.

The GP constructs a joint probability distribution over the variables, assuming a multivariate Gaussian distribution. As the number of observations increases, the posterior distribution becomes more precise, enabling the algorithm to more effectively identify promising regions in the weight space. The Bayesian optimization process involves several iterations. First, the GP model is updated by fitting it to the observed data. Next, the algorithm selects the next set of weights by maximizing the acquisition function, which uses the posterior distribution to choose the next sample from the search space. Finally, the objective function is evaluated at these weights. This iterative process continues until a convergence criterion is met, ensuring that the optimization effectively identifies the optimal weights for the metrics.

### 2.3 METAMETRICS-MT Settings

#### 2.3.1 Hybrid Mode

In the WMT24 shared task dataset, we observe that some samples lack references in the challenge sets, even for reference-based metrics. To address this issue, we implement a hybrid mode that switches from reference-based to reference-free

metrics when reference data is unavailable.

### 2.3.2 Same Language Optimization

During the optimization process, we train a dedicated model for each known language pair in the training set to ensure optimal performance. If a language pair is not present in the training set, we use the entire dataset for tuning.

## 3 Experimental Setup

### 3.1 Training Datasets and Hyper-parameters

We introduce two versions of METAMETRICS-MT to accommodate both reference-based and reference-free evaluations: METAMETRICS-MT, which employs reference-based metrics, and METAMETRICS-MT-QE, which utilizes reference-free metrics. We train METAMETRICS-MT and METAMETRICS-MT-QE using 3 years of MQM datasets from the WMT shared tasks spanning 2020 to 2022 (Mathur et al., 2020; Freitag et al., 2021, 2022). The dataset used for tuning is at the segment level, with Kendall's $\tau$ correlation as the evaluation metric. For the Bayesian optimization, we run GP with a Matérn kernel (Williams and Rasmussen, 2006), a generalization of the RBF kernel, using $\nu = 2.5$. The optimization is performed over 100 steps, starting with 5 initialization points.

### 3.2 Metrics for METAMETRICS-MT

We describe the reference-based metrics utilized for METAMETRICS-MT. During the selection process, we included only metrics that can run on a commercial GPU with 40GB of memory. Consequently, XCOMET-XXL and CometKiwi-XXL were not considered. Additionally, we limited the use of the OpenAI API to GPT4o-mini, which is significantly more cost-effective than other GPT-4 model options.

#### 3.2.1 Reference-based Metric

We utilize nine different metrics in our optimization, including three variations of MetricX-23 and two different BERTScore metrics using precision and F1. The metrics under study are as follows:

**BERTScore (Zhang et al., 2019)** The metric calculates cosine similarity scores for each token in the candidate sentences against each token in the reference sentences, using contextual embeddings derived from pre-trained BERT-based models. From these similarities, BERTScore computes

precision, recall, and F1 scores. In our metrics, we utilize the precision and F1 scores, employing DeBERTa-XL-MNLI (He et al., 2020) as our model, as recommended by the authors.

**YISI-1 (Lo, 2019)** The metric computes the semantic similarity between translations from MT and human references by aggregating lexical semantic similarities, which are weighted by inverse document frequency (IDF) based on the contextual embeddings extracted from pre-trained language model, specifically the last hidden layer of mBERT in our case.

**BLEURT (Sellam et al., 2020)** The metric is fine-tuned using Direct Assessment (DA) dataset. BLEURT jointly encodes the translation and reference using the [CLS] token as an embedding to represent the pair. We employ the BLEURT-20 checkpoint (Pu et al., 2021), which was trained on RemBERT (Chung et al., 2020) using DA data from prior shared tasks between 2015 and 2019 and augmented with synthetic data generated from Wikipedia articles.

**COMET-22 (Rei et al., 2022)** The metric is an ensemble of COMET estimator (Rei et al., 2020) fine-tuned on DA and a Sequence Tagger trained on Multidimensional Quality Metrics (MQM) annotations. We utilize the wmt22-comet-da as our COMET-22 checkpoint, in which the COMET Estimator model and the sequence tagging model are trained on top of XLM-R using DA from 2017 to 2020 and InfoXLM (Chi et al., 2021), respectively.

**XCOMET-XL (Guerreiro et al., 2023)** The metric that performs both sentence-level evaluation and error span detection, making it a more interpretable learned metric. The model utilizes XLM-R XL (3.5B) (Goyal et al., 2021) which is trained in stages, starting with DA annotations and then fine-tuned on MQM data.

**MetricX-23 (Juraska et al., 2023)** The metric uses mT5 encoder-decoder language model. We leverage three different variations of MetricX-23, each fine-tuned from the mT5-Large, mT5-XL, and mT5-XXL respectively. The fine-tuning was performed using DA data from 2015-2020, MQM data from 2020-2021, and synthetic data.

#### 3.2.2 Reference-free Metric

We utilize six different metrics in our optimization, including two variations of CometKiwi and

461

| Model | overall | | en-de | | | | en-es | | | | ja-zh | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $r$ | sys/seg | $r$ | sys | $r$ | seg | $r$ | sys | $r$ | seg | $r$ | sys | $r$ | seg |
| | | avg. corr | | SPA | | acc-t | | SPA | | acc-t | | SPA | | acc-t |
| **Reference-based** | | | | | | | | | | | | | | |
| sentinel-ref-mqm | 10 | 0.513 | 7 | 0.405 | 18 | 0.429 | 4 | 0.581 | 8 | 0.680 | 8 | 0.545 | 17 | 0.435 |
| BLEU | 9 | 0.589 | 4 | 0.736 | 16 | 0.431 | 6 | 0.512 | 8 | 0.680 | 6 | 0.740 | 17 | 0.435 |
| spBLEU | 9 | 0.593 | 4 | 0.741 | 17 | 0.431 | 6 | 0.523 | 7 | 0.680 | 6 | 0.744 | 16 | 0.436 |
| chrfS | 8 | 0.606 | 4 | 0.742 | 14 | 0.434 | 6 | 0.549 | 6 | 0.682 | 4 | 0.788 | 14 | 0.444 |
| chrF | 8 | 0.608 | 4 | 0.750 | 15 | 0.431 | 5 | 0.581 | 8 | 0.680 | 5 | 0.767 | 16 | 0.436 |
| MEE4 | 7 | 0.609 | 5 | 0.731 | 13 | 0.437 | 7 | 0.504 | 4 | 0.683 | 2 | 0.855 | 13 | 0.446 |
| BERTScore | 7 | 0.617 | 4 | 0.749 | 14 | 0.435 | 4 | 0.587 | 6 | 0.682 | 4 | 0.799 | 12 | 0.451 |
| YiSi-1 | 6 | 0.630 | 4 | 0.759 | 13 | 0.436 | 4 | 0.609 | 7 | 0.681 | 3 | 0.835 | 11 | 0.458 |
| PrismRefSmall | 5 | 0.642 | 4 | 0.772 | 14 | 0.433 | 4 | 0.634 | 8 | 0.680 | 2 | 0.875 | 11 | 0.457 |
| PrismRefMedium | 5 | 0.646 | 4 | 0.776 | 14 | 0.434 | 3 | 0.652 | 7 | 0.680 | 2 | 0.872 | 10 | 0.462 |
| BLCOM_1 | 4 | 0.664 | 3 | 0.840 | 10 | 0.455 | 3 | 0.680 | 6 | 0.681 | 3 | 0.843 | 7 | 0.488 |
| BLEURT-20 | 3 | 0.686 | 2 | 0.881 | 7 | 0.486 | 3 | 0.695 | 6 | 0.681 | 1 | 0.887 | 8 | 0.484 |
| COMET | 3 | 0.688 | 2 | 0.879 | 8 | 0.482 | 2 | 0.778 | 5 | 0.683 | 4 | 0.813 | 6 | 0.496 |
| XCOMET | 2 | 0.719 | 1 | **0.906** | 3 | 0.530 | 2 | 0.788 | 1 | **0.688** | 2 | 0.890 | 7 | 0.510 |
| MetricX-24 (Hybrid) | 1 | 0.721 | 2 | 0.874 | 2 | 0.532 | 2 | 0.799 | 3 | 0.685 | 1 | <u>0.897</u> | 2 | 0.539 |
| METAMETRICS-MT | 1 | <u>0.724</u> | 2 | 0.882 | 1 | **0.542** | 2 | **0.804** | 2 | <u>0.686</u> | 3 | 0.871 | 1 | **0.561** |
| METAMETRICS-MT (Same Lang.) | 2 | 0.723 | 1 | <u>0.883</u> | 1 | **0.542** | 2 | <u>0.803</u> | 2 | <u>0.686</u> | 3 | 0.874 | 2 | <u>0.550</u> |
| METAMETRICS-MT (Hybrid) | 1 | **0.725** | 2 | <u>0.883</u> | 1 | **0.542** | 1 | **0.804** | 2 | <u>0.686</u> | 2 | 0.873 | 1 | **0.561** |
| **Reference-free** | | | | | | | | | | | | | | |
| CometKiwi | 5 | 0.640 | 5 | 0.732 | 9 | 0.467 | 3 | 0.693 | 4 | 0.684 | 5 | 0.776 | 7 | 0.490 |
| sentinel-cand-mqm | 5 | 0.650 | 3 | 0.822 | 4 | 0.517 | 2 | 0.785 | 4 | 0.683 | 7 | 0.610 | 8 | 0.481 |
| bright-qe | 4 | 0.681 | 3 | 0.816 | 6 | 0.500 | 2 | 0.792 | 1 | **0.689** | 4 | 0.805 | 8 | 0.484 |
| XCOMET-QE | 3 | 0.695 | 1 | **0.889** | 4 | <u>0.520</u> | 1 | 0.801 | 2 | <u>0.687</u> | 4 | 0.808 | 10 | 0.463 |
| CometKiwi-XXL | 3 | 0.703 | 3 | 0.839 | 9 | 0.481 | 1 | **0.843** | 8 | 0.680 | 2 | <u>0.881</u> | 8 | 0.494 |
| gemba_esa | 2 | <u>0.711</u> | 4 | 0.793 | 5 | 0.507 | 1 | <u>0.838</u> | 5 | 0.683 | 1 | **0.908** | 2 | **0.539** |
| MetricX-24-QE (Hybrid) | 2 | **0.714** | 2 | 0.878 | 3 | **0.526** | 2 | 0.789 | 4 | 0.685 | 2 | 0.875 | 3 | <u>0.530</u> |
| METAMETRICS-MT-QE | 3 | 0.684 | 2 | <u>0.860</u> | 6 | 0.497 | 3 | 0.711 | 2 | 0.686 | 3 | 0.837 | 4 | 0.516 |
| METAMETRICS-MT-QE (Same Lang.) | 4 | 0.688 | 2 | <u>0.860</u> | 7 | 0.497 | 4 | 0.709 | 2 | 0.686 | 4 | 0.853 | 5 | 0.524 |

Table 2: WMT24 results (MQM). **Bold** and <u>underline</u> values indicate the best and second best performance, respectively.

three variations of MetricX-23. We describe the reference-free metrics used for METAMETRICS-MT-QE as follows:

**CometKiwi (Rei et al., 2022)** The metric is a reference-free learned metric fine-tuned on DA on top of RemBERT (Chung et al., 2020) and the same sequence tagger as COMET-22. However, it operates with reference-less inputs during inference. We use two distinct metrics from CometKiwi, each associated with its own separate checkpoint: `wmt22-cometkiwi-da` and `wmt23-cometkiwi-da-xl`. The latter checkpoint replaces InfoXLM with XLM-R XL (3.5B) and is trained on the same dataset, but it also includes newly released DA for Indian languages, which were added as additional training data for the 2023 Quality Estimation (QE) shared task.

**GEMBA-MQM (Kocmi and Federmann, 2023)** The metric is a GPT-based evaluation metric designed for error quality span marking. It employs a three-shot prompting approach using the GPT-4 model, specifically GPT-4o mini in our case.

**MetricX-23-QE (Juraska et al., 2023)** The metric is a reference-free learned metric similar to MetricX-23. We also utilize three different variations, each fine-tuned from the mT5-L, mT5-XL, and mT5-XXL checkpoints, respectively.

## 4 Results and Discussion

### 4.1 Optimized Metric Configuration

Table 1 shows the weight proportion of each metric for METAMETRICS-MT. The optimized configuration is notably sparse. When a metric does not positively contribute to improving performance, the GP assigns it a weight of zero. This is supported by Figure 1, where the GP selects metrics with high Kendall correlation coefficients relative to other provided metrics. In contrast, metrics with low Kendall correlation coefficients are excluded.

| Metric | all | en-de | en-es | ja-zh |
|---|---|---|---|---|
| **Reference-based** | | | | |
| sentinel-ref-mqm | 0.513 | 0.417 | 0.631 | 0.490 |
| BLEU | 0.589 | 0.583 | 0.596 | 0.588 |
| spBLEU | 0.593 | 0.586 | 0.602 | 0.590 |
| chrF | 0.606 | 0.589 | 0.615 | 0.616 |
| chrfS | 0.608 | 0.591 | 0.630 | 0.602 |
| BERTScore | 0.610 | 0.584 | 0.594 | 0.651 |
| MEE4 | 0.617 | 0.592 | 0.635 | 0.625 |
| damonmonli | 0.640 | 0.599 | 0.688 | 0.633 |
| YiSi-1 | 0.643 | 0.603 | 0.657 | 0.666 |
| PrismRefSmall | 0.646 | 0.605 | 0.666 | 0.667 |
| PrismRefMedium | 0.650 | 0.669 | 0.734 | 0.545 |
| BLCOM_1 | 0.684 | 0.679 | 0.698 | 0.676 |
| BLEURT-20 | 0.686 | 0.683 | 0.688 | 0.685 |
| COMET-22 | 0.695 | 0.705 | <u>0.744</u> | 0.636 |
| XCOMET | 0.719 | **0.717** | 0.740 | 0.700 |
| MetricX-24 (Hybrid) | <u>0.721</u> | 0.703 | 0.742 | **0.718** |
| METAMETRICS-MT (Hybrid) | **0.725** | <u>0.713</u> | **0.745** | <u>0.717</u> |
| **Reference-free** | | | | |
| sentinel-src-mqm | 0.513 | 0.418 | 0.630 | 0.491 |
| XLsimMqm | 0.515 | 0.531 | 0.520 | 0.493 |
| sentinel-cand-mqm | 0.630 | 0.597 | 0.645 | 0.647 |
| CometKiwi | 0.635 | 0.569 | 0.644 | 0.691 |
| bright-qe | 0.665 | 0.647 | 0.681 | 0.665 |
| XCOMET-QE | 0.689 | <u>0.680</u> | 0.730 | 0.655 |
| MetricX-24-QE (Hybrid) | **0.714** | **0.702** | 0.737 | <u>0.702</u> |
| gemba_esa | <u>0.711</u> | 0.650 | **0.761** | **0.724** |
| METAMETRICS-MT-QE | 0.681 | 0.658 | <u>0.740</u> | 0.644 |

Table 3: Detailed WMT24 results per language category. **Bold** and <u>underline</u> values indicate the best and second best performance, respectively.

| Metric | all | sys | seg |
|---|---|---|---|
| **Reference-based** | | | |
| sentinel-ref-mqm | 0.513 | 0.510 | 0.515 |
| BLEU | 0.589 | 0.663 | 0.515 |
| spBLEU | 0.593 | 0.669 | 0.516 |
| chrF | 0.606 | 0.693 | 0.520 |
| chrfS | 0.608 | 0.699 | 0.516 |
| BERTScore | 0.609 | 0.697 | 0.522 |
| MEE4 | 0.617 | 0.712 | 0.522 |
| damonmonli | 0.640 | 0.734 | 0.547 |
| YiSi-1 | 0.642 | 0.760 | 0.524 |
| PrismRefSmall | 0.646 | 0.766 | 0.526 |
| PrismRefMedium | 0.650 | 0.739 | 0.560 |
| BLCOM_1 | 0.684 | 0.803 | 0.566 |
| BLEURT-20 | 0.686 | 0.821 | 0.550 |
| COMET-22 | 0.695 | 0.833 | 0.557 |
| XCOMET | 0.719 | **0.862** | 0.576 |
| MetricX-24 (Hybrid) | <u>0.721</u> | <u>0.857</u> | <u>0.586</u> |
| METAMETRICS-MT (Hybrid) | **0.725** | 0.853 | **0.596** |
| **Reference-free** | | | |
| sentinel-src-mqm | 0.513 | 0.511 | 0.515 |
| XLsimMqm | 0.515 | 0.506 | 0.523 |
| sentinel-cand-mqm | 0.630 | 0.734 | 0.525 |
| CometKiwi | 0.635 | 0.738 | 0.532 |
| bright-qe | 0.664 | 0.788 | 0.541 |
| XCOMET-QE | 0.688 | 0.823 | 0.554 |
| gemba_esa | <u>0.711</u> | 0.846 | <u>0.576</u> |
| MetricX-24-QE (Hybrid) | **0.714** | <u>0.847</u> | **0.580** |
| METAMETRICS-MT-QE | 0.681 | 0.804 | 0.557 |

Table 4: Detailed WMT24 results for segment-level and system-level. **Bold** and <u>underline</u> values indicate the best and second best performance, respectively.

Interestingly, in both reference-based and reference-free settings, the optimization process consistently selects only one variant of MetricX-23, specifically MetricX-23-XXL, even though all three variants of MetricX-23 exhibit high Kendall correlation coefficients. The optimization process favors MetricX-23-XXL as the highest-performing metric, leading to the exclusion of the other two variants during the GP assignment. This enhances the efficiency of METAMETRICS-MT as we would only need to use fewer metrics for METAMETRICS-MT. Thus, given a set of metrics, the optimization process would prioritize high-performing metrics, such as the MetricX-23 and COMET variants as shown, leading METAMETRICS-MT and METAMETRICS-MT-QE to construct a better and more robust metric.

### 4.2 Results on WMT24 Shared Task

Table 2 presents the WMT24 shared task results, including system-level soft pairwise ranking accuracy (sys SPA) proposed by Thompson et al. (2024), segment-level pairwise ranking accuracy with tie

calibration (seg acc-t) as described by Deutsch et al. (2023), and system- and segment-level Pearson correlation (avg. corr), as outlined in the WMT23 Metrics Shared Task (Freitag et al., 2023). Based on the overall system and segment average correlation and system accuracy, METAMETRICS-MT outperforms all metrics in the primary submission, with METAMETRICS-MT (Hybrid) achieving the highest performance among its variants.

Table 3 further highlights the performance, where METAMETRICS-MT delivers superior results for en-es, while also maintaining strong performance in en-de and ja-zh, indicating that our methods generalize well across different language pairs. The breakdown in Table 4 shows that METAMETRICS-MT achieves the best segment-level performance, consistent with our optimization approach targeting Kendall correlation at the segment level. Given that our metric optimization

Figure 1: Heatmaps showing Kendall correlation coefficients between human scores and MT metrics over 3 years of MQM datasets from the WMT shared tasks (2020-2022). Panel (a) displays correlations for the metrics used in METAMETRICS-MT, while panel (b) displays correlations for the metrics used in METAMETRICS-MT-QE.

focuses solely on segment-level correlation, incorporating a different weighting method to account for system-level settings could further improve METAMETRICS-MT's alignment with system-level accuracy. While METAMETRICS-MT-QE does not match the performance of gemba_esa, MetricX-24-QE (Hybrid), or CometKiwi-XXL, it remains competitive at the segment level for the en-es language pair. Incorporating better reference-free models such as CometKiwi-XXL

and GEMBA-MQM with GPT-4o instead of GPT-4o mini may help improve the performance of METAMETRICS-MT-QE.

## 4.3 Compute Efficiency

We only run models that can be executed on GPUs with 40GB of memory. We limit our resource usage to GPT-4o mini, a smaller and lower-performing version of GPT-4o, while GEMBA-MQM is a GPT-4 based metric. This constraint restricts our ability to achieve state-of-the-art results or surpass GEMBA-based metrics using GPT-4. However, we demonstrate that even without employing high-memory models like XCOMET-Ensemble in our reference-based setting, we can still outperform other models. Additionally, our QE metric remains competitive and on par with XCOMET-QE.

## 5 Conclusion

In this paper, we propose METAMETRICS-MT, a novel metric designed to evaluate MT tasks by aligning with human preferences through Bayesian optimization with GP. METAMETRICS-MT effectively combines and optimizes existing MT metrics based on human feedback, resulting in a highly flexible and efficient evaluation tool. Our findings show that METAMETRICS-MT surpasses existing baselines for reference-based metrics, establishing a new state-of-the-art, while its reference-free metric performance rivals the best models available. Additionally, METAMETRICS-MT can be tailored to various factors, such as performance and efficiency, making it adaptable to diverse requirements.

## Ethical Considerations

Our research focuses on evaluating MT systems using a newly proposed metric. We are committed to conducting our evaluations with the highest levels of transparency and fairness. By prioritizing these principles, we aim to set a standard for reliability and objectivity in the assessment of the system.

## Limitations

We optimize METAMETRICS-MT using segment-level scores from the MQM dataset. Future work could extend this to other objective functions or system-level optimization and explore non-MQM datasets like DA for further insights. We did not include metrics such as XCOMET-XXL, XCOMET-Ensemble, and XCOMET-QE-Ensemble due to computational constraints.

# References

David Anugraha, Genta Indra Winata, Chenyue Li, Patrick Amadeus Irawan, and En-Shiun Annie Lee. 2024. Proxylm: Predicting language model performance on multilingual tasks via proxy models. *arXiv preprint arXiv:2406.09334*.

Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, He-Yan Huang, and Ming Zhou. 2021. Infoxlm: An information-theoretic framework for cross-lingual language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3576–3588.

Hyung Won Chung, Thibault Fevry, Henry Tsai, Melvin Johnson, and Sebastian Ruder. 2020. Rethinking embedding coupling in pre-trained language models. *arXiv preprint arXiv:2010.12821*.

Daniel Deutsch, George Foster, and Markus Freitag. 2023. Ties matter: Meta-evaluating modern metrics with pairwise accuracy and tie calibration. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12914–12929.

Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are llms breaking mt metrics? results of the wmt24 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frédéric Blain, Daniel Deutsch, Craig Stewart, et al. 2023. Results of wmt23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André FT Martins. 2022. Results of wmt22 metrics shared task: Stop using bleu–neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021. Results of the wmt21 metrics shared task: Evaluating metrics with expert-based human evaluations on ted and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774.

Naman Goyal, Jingfei Du, Myle Ott, Giri Ananthara- man, and Alexis Conneau. 2021. Larger-scale trans- formers for multilingual masked language modeling.

In *Proceedings of the 6th Workshop on Represen- tation Learning for NLP (RepL4NLP-2021)*, pages 29–33.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. xcomet: Transparent machine translation eval- uation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. Metricx-23: The google submission to the wmt 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767.

Tom Kocmi and Christian Federmann. 2023. Gemba- mqm: Detecting translation quality error spans with gpt-4. In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775.

Chi-kiu Lo. 2019. Yisi-a unified semantic mt quality evaluation and estimation metric for languages with different levels of available resources. In *Proceed- ings of the Fourth Conference on Machine Transla- tion (Volume 2: Shared Task Papers, Day 1)*, pages 507–513.

Nitika Mathur, Johnny Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020. Results of the wmt20 metrics shared task. In *Proceedings of the Fifth Con- ference on Machine Translation*, pages 688–725.

Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. Learning compact metrics for mt. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Lan- guage Processing*, pages 751–762.

Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Sev- enth Conference on Machine Translation (WMT)*, pages 578–585.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt eval- uation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. Bleurt: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the As- sociation for Computational Linguistics*, pages 7881– 7892.

Brian Thompson, Nitika Mathur, Daniel Deutsch, and Huda Khayrallah. 2024. Improving statistical significance in human evaluation of automatic metrics via soft pairwise accuracy. *arXiv preprint arXiv:2409.09598*.

Christopher KI Williams and Carl Edward Rasmussen. 2006. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA.

Genta Indra Winata, David Anugraha, Lucky Susanto, Garry Kuwanto, and Derry Tanti Wijaya. 2024a. Metametrics: Calibrating metrics for generation tasks using human preferences. *arXiv preprint arXiv:2410.02381*.

Genta Indra Winata, Hanyang Zhao, Anirban Das, Wenpin Tang, David D Yao, Shi-Xiong Zhang, and Sambit Sahu. 2024b. Preference tuning with human feedback on language, speech, and vision tasks: A survey. *arXiv preprint arXiv:2409.11564*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

## A  Pre-processing

The pre-processing can be defined as follows:

1. **Clipping:** Let the valid range for $y_i$ be defined by $[y_i^{\min}, y_i^{\max}]$. The clipped metric score $y_i'$ can be defined as:

$$y_i' = \begin{cases} y_i^{\min} & \text{if } y_i < y_i^{\min}, \\ y_i & \text{if } y_i^{\min} \leq y_i \leq y_i^{\max}, \\ y_i^{\max} & \text{if } y_i > y_i^{\max}. \end{cases} \quad (4)$$

2. **Normalization:** After clipping, the score is normalized to a common scale of $[0, 1]$:

$$\tilde{y}_i = \frac{y_i' - y_i^{\min}}{y_i^{\max} - y_i^{\min}}. \quad (5)$$

3. **Inversion (if applicable):** If the metric is such that higher scores indicate worse performance, we invert the normalized score:

$$\tilde{y}_i = 1 - \tilde{y}_i. \quad (6)$$

## B  Additional Results

We provide additional details for the results of WMT24 for each task in Tables 5, 6, and 7. Additional results for each domain are also provided in Table 8.

| Domain | | literary | | news | | social | | speech | | literary | | news | | social | | speech |
| | | task1 | | task2 | | task3 | | task4 | | task5 | | task6 | | task7 | | task8 |
| Metric / Level | r | sys SPA | r | sys SPA | r | sys SPA | r | sys SPA | r | seg acc-t | r | seg acc-t | r | seg acc-t | r | seg acc-t |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Reference-based** | | | | | | | | | | | | | | | | |
| sentinel-ref-mqm | 4 | 0.525 | 4 | 0.535 | 6 | 0.439 | 6 | 0.461 | 9 | 0.351 | 9 | 0.421 | 16 | 0.520 | 13 | 0.240 |
| BLEU | 2 | 0.795 | 1 | 0.807 | 5 | 0.691 | 4 | 0.709 | 5 | 0.535 | 9 | 0.421 | 15 | 0.522 | 11 | 0.433 |
| spBLEU | 2 | 0.785 | 1 | 0.810 | 5 | 0.697 | 4 | 0.700 | 4 | 0.540 | 9 | 0.421 | 15 | 0.522 | 10 | 0.446 |
| chrF | 2 | 0.774 | 1 | **0.831** | 4 | 0.728 | 3 | 0.723 | 4 | 0.540 | 9 | 0.421 | 14 | 0.523 | 10 | 0.445 |
| chrfS | 2 | 0.797 | 1 | 0.826 | 4 | 0.712 | 3 | 0.736 | 4 | 0.543 | 9 | 0.421 | 13 | 0.525 | 9 | 0.449 |
| BERTScore | 2 | 0.777 | 1 | 0.821 | 4 | 0.708 | 4 | 0.712 | 4 | 0.550 | 8 | 0.424 | 12 | 0.526 | 11 | 0.436 |
| MEE4 | 2 | 0.792 | 1 | 0.826 | 5 | 0.688 | 4 | 0.712 | 4 | 0.549 | 9 | 0.421 | 10 | 0.531 | 9 | 0.452 |
| damonmonli | 2 | 0.734 | 1 | 0.788 | 5 | 0.695 | 5 | 0.613 | 7 | 0.503 | 7 | 0.427 | 14 | 0.523 | 12 | 0.404 |
| YiSi-1 | 2 | 0.761 | 1 | 0.822 | 4 | 0.719 | 3 | 0.760 | 3 | 0.555 | 9 | 0.421 | 12 | 0.526 | 8 | 0.456 |
| PrismRefSmall | 2 | 0.786 | 1 | _0.829_ | 4 | 0.750 | 3 | 0.736 | 5 | 0.526 | 8 | 0.423 | 13 | 0.524 | 7 | 0.464 |
| PrismRefMedium | 2 | 0.761 | 1 | **0.831** | 4 | 0.756 | 4 | 0.722 | 4 | 0.536 | 8 | 0.424 | 11 | 0.528 | 8 | 0.461 |
| BLCOM_1 | 1 | 0.828 | 1 | 0.812 | 3 | 0.803 | 2 | 0.833 | 3 | 0.562 | 7 | 0.427 | 9 | 0.535 | 5 | 0.487 |
| BLEURT-20 | 1 | 0.827 | 2 | 0.768 | 2 | 0.842 | 3 | 0.784 | 4 | 0.544 | 5 | 0.444 | 7 | 0.554 | 4 | 0.494 |
| COMET-22 | 1 | 0.814 | 1 | 0.804 | 2 | 0.852 | 2 | 0.813 | 2 | 0.571 | 6 | 0.437 | 6 | 0.559 | 3 | 0.503 |
| XCOMET | 1 | _0.830_ | 1 | 0.782 | 1 | _0.889_ | 2 | **0.845** | 2 | 0.573 | 3 | _0.479_ | 2 | 0.575 | 2 | _0.510_ |
| MetricX-24 (Hybrid) | 1 | **0.840** | 1 | 0.774 | 1 | 0.874 | 2 | _0.816_ | 2 | _0.580_ | 3 | 0.478 | 2 | _0.576_ | 1 | **0.520** |
| METAMETRICS-MT (Hybrid) | 1 | 0.822 | 2 | 0.763 | 1 | **0.896** | 3 | 0.788 | 1 | **0.597** | 2 | **0.493** | 1 | **0.588** | 2 | 0.506 |
| **Reference-free** | | | | | | | | | | | | | | | | |
| sentinel-src-mqm | 4 | 0.525 | 4 | 0.534 | 6 | 0.438 | 6 | 0.461 | 9 | 0.351 | 9 | 0.421 | 16 | 0.520 | 13 | 0.240 |
| XLsimMqm | 4 | 0.478 | 4 | 0.497 | 5 | 0.613 | 3 | 0.768 | 8 | 0.474 | 1 | **0.532** | 10 | 0.531 | 12 | 0.410 |
| sentinel-cand-mqm | 2 | 0.776 | 2 | 0.735 | 1 | **0.896** | 3 | 0.760 | 4 | 0.547 | 2 | _0.501_ | 4 | 0.569 | 6 | 0.480 |
| CometKiwi | 3 | 0.722 | 2 | 0.723 | 4 | 0.732 | 4 | 0.685 | 5 | 0.535 | 5 | 0.445 | 9 | 0.532 | 10 | 0.443 |
| bright-qe | 2 | 0.795 | 2 | 0.755 | 3 | 0.760 | 2 | 0.827 | 6 | 0.517 | 4 | 0.457 | 8 | 0.547 | 7 | 0.469 |
| XCOMET-QE | 2 | 0.758 | 1 | **0.790** | 2 | 0.850 | 1 | **0.882** | 4 | 0.541 | 3 | 0.480 | 5 | 0.565 | 3 | _0.498_ |
| gemba_esa | 1 | **0.820** | 2 | 0.755 | 3 | 0.801 | 2 | 0.815 | 3 | _0.562_ | 5 | 0.450 | 3 | _0.569_ | 6 | 0.474 |
| MetricX-24-QE (Hybrid) | 2 | _0.809_ | 1 | _0.783_ | 1 | _0.863_ | 1 | _0.860_ | 2 | **0.575** | 4 | 0.460 | 3 | **0.573** | 1 | **0.518** |
| METAMETRICS-MT-QE | 3 | 0.691 | 3 | 0.690 | 2 | 0.811 | 1 | 0.852 | 6 | 0.520 | 4 | 0.457 | 6 | 0.555 | 7 | 0.471 |

Table 5: Detailed result for language pair en-de. **Bold** and underline values indicate the best and second best performance, respectively.

| Domain | | literary | | news | | social | | speech | | literary | | news | | social | | speech |
| | | task9 | | task10 | | task11 | | task12 | | task13 | | task14 | | task15 | | task16 |
| Metric / Level | r | sys SPA | r | sys SPA | r | sys SPA | r | sys SPA | r | seg acc-t | r | seg acc-t | r | seg acc-t | r | seg acc-t |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Reference-based** | | | | | | | | | | | | | | | | |
| sentinel-ref-mqm | 3 | 0.564 | 4 | 0.460 | 5 | 0.599 | 4 | 0.556 | 5 | 0.615 | 4 | 0.715 | 8 | 0.744 | 6 | 0.535 |
| BLEU | 3 | 0.595 | 4 | 0.557 | 5 | 0.624 | 5 | 0.480 | 5 | 0.615 | 4 | 0.715 | 7 | 0.745 | 5 | 0.536 |
| spBLEU | 3 | 0.602 | 3 | 0.595 | 4 | 0.635 | 5 | 0.486 | 4 | 0.615 | 4 | 0.715 | 7 | 0.745 | 5 | 0.536 |
| chrF | 3 | 0.621 | 3 | 0.593 | 4 | 0.657 | 5 | 0.490 | 4 | 0.615 | 4 | 0.715 | 8 | 0.744 | 4 | 0.537 |
| chrfS | 2 | 0.648 | 3 | 0.604 | 4 | 0.667 | 5 | 0.472 | 3 | 0.617 | 4 | 0.715 | 6 | 0.746 | 5 | 0.537 |
| BERTScore | 2 | 0.665 | 1 | 0.715 | 3 | 0.679 | 5 | 0.488 | 3 | 0.617 | 2 | _0.717_ | 5 | 0.747 | 5 | 0.537 |
| MEE4 | 2 | 0.651 | 2 | 0.628 | 3 | 0.677 | 5 | 0.467 | 3 | 0.617 | 4 | 0.715 | 3 | 0.750 | 4 | 0.539 |
| damonmonli | 1 | 0.720 | 2 | 0.673 | 2 | 0.737 | 4 | 0.555 | 2 | 0.621 | 4 | 0.715 | 5 | 0.747 | 5 | 0.536 |
| YiSi-1 | 1 | 0.706 | 2 | 0.673 | 3 | _0.715_ | 5 | 0.505 | 3 | 0.617 | 4 | 0.715 | 6 | 0.745 | 4 | 0.538 |
| PrismRefSmall | 1 | 0.727 | 2 | 0.624 | 2 | 0.724 | 5 | 0.518 | 5 | 0.615 | 3 | 0.716 | 8 | 0.745 | 5 | 0.537 |
| PrismRefMedium | 1 | 0.733 | 2 | 0.649 | 2 | 0.745 | 5 | 0.518 | 4 | 0.616 | 3 | 0.716 | 7 | 0.745 | 5 | 0.536 |
| BLCOM_1 | 2 | 0.702 | 2 | 0.675 | 2 | 0.773 | 4 | 0.623 | 3 | 0.617 | 4 | 0.715 | 5 | 0.747 | 4 | 0.541 |
| BLEURT-20 | 2 | 0.702 | 2 | 0.648 | 1 | 0.841 | 4 | 0.587 | 2 | 0.620 | 4 | 0.715 | 6 | 0.746 | 6 | 0.535 |
| COMET-22 | 1 | **0.755** | 1 | **0.731** | 1 | **0.865** | 3 | 0.653 | 1 | **0.626** | 4 | 0.715 | 4 | 0.750 | 3 | _0.551_ |
| XCOMET | 1 | 0.733 | 1 | 0.677 | 1 | 0.840 | 2 | _0.685_ | 1 | _0.625_ | 2 | _0.717_ | 1 | **0.756** | 3 | 0.548 |
| MetricX-24 (Hybrid) | 1 | _0.741_ | 1 | 0.683 | 1 | 0.846 | 2 | **0.691** | 2 | 0.621 | 4 | 0.715 | 3 | 0.750 | 2 | **0.559** |
| METAMETRICS-MT (Hybrid) | 1 | 0.734 | 1 | 0.688 | 1 | _0.852_ | 2 | 0.682 | 2 | 0.619 | 1 | **0.720** | 2 | _0.753_ | 3 | 0.550 |
| **Reference-free** | | | | | | | | | | | | | | | | |
| sentinel-src-mqm | 3 | 0.565 | 4 | 0.456 | 5 | 0.598 | 4 | 0.554 | 5 | 0.615 | 4 | 0.715 | 8 | 0.744 | 6 | 0.535 |
| XLsimMqm | 4 | 0.363 | 2 | 0.645 | 6 | 0.410 | 3 | 0.640 | 4 | 0.615 | 4 | 0.715 | 6 | 0.745 | 4 | 0.537 |
| sentinel-cand-mqm | 2 | 0.695 | 1 | 0.678 | 2 | 0.780 | 2 | 0.690 | 2 | 0.620 | 1 | _0.720_ | 4 | 0.749 | 4 | 0.537 |
| CometKiwi | 2 | 0.641 | 2 | 0.661 | 2 | 0.767 | 2 | 0.681 | 2 | 0.620 | 3 | 0.716 | 3 | _0.751_ | 3 | 0.553 |
| bright-qe | 3 | 0.583 | 1 | 0.677 | 2 | 0.764 | 1 | **0.772** | 2 | 0.621 | 2 | 0.718 | 3 | _0.751_ | 1 | **0.571** |
| XCOMET-QE | 1 | _0.731_ | 2 | 0.673 | 2 | 0.779 | 2 | 0.700 | 2 | _0.622_ | 1 | **0.721** | 2 | **0.754** | 3 | 0.547 |
| gemba_esa | 1 | **0.740** | 1 | **0.723** | 1 | **0.820** | 2 | _0.704_ | 2 | 0.621 | 2 | 0.718 | 5 | 0.746 | 3 | 0.549 |
| MetricX-24-QE (Hybrid) | 1 | 0.727 | 1 | 0.694 | 1 | _0.818_ | 2 | 0.703 | 1 | _0.622_ | 4 | 0.715 | 5 | 0.748 | 2 | 0.563 |
| METAMETRICS-MT-QE | 2 | 0.661 | 1 | _0.711_ | 2 | 0.751 | 2 | 0.692 | 1 | **0.624** | 2 | 0.717 | 4 | 0.749 | 2 | _0.565_ |

Table 6: Detailed WMT24 result for language pair en-es. **Bold** and underline values indicate the second best performance, respectively.

| Domain | literary | | news | | speech | | literary | | news | | speech | |
| Metric | task17 | | task18 | | task19 | | task20 | | task21 | | task22 | |
| Level | r | sys SPA | r | sys SPA | r | sys SPA | r | seg acc-t | r | seg acc-t | r | seg acc-t |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Reference-based** | | | | | | | | | | | | |
| sentinel-ref-mqm | 5 | 0.504 | 7 | 0.494 | 8 | 0.569 | 11 | 0.532 | 8 | 0.497 | 12 | 0.197 |
| BLEU | 4 | 0.637 | 3 | 0.762 | 8 | 0.562 | 11 | 0.532 | 8 | 0.497 | 11 | 0.205 |
| spBLEU | 4 | 0.699 | 4 | 0.755 | 7 | 0.743 | 9 | 0.535 | 7 | 0.497 | 7 | 0.506 |
| chrF | 4 | 0.721 | 3 | 0.768 | 6 | 0.766 | 9 | 0.536 | 8 | 0.497 | 6 | 0.513 |
| chrfS | 3 | 0.768 | 3 | 0.773 | 5 | 0.823 | 9 | 0.537 | 7 | 0.497 | 5 | 0.526 |
| BERTScore | 3 | 0.786 | 5 | 0.748 | 5 | 0.833 | 9 | 0.536 | 6 | 0.500 | 5 | 0.524 |
| MEE4 | 2 | 0.816 | 3 | 0.789 | 2 | 0.892 | 8 | 0.538 | 7 | 0.497 | 5 | 0.521 |
| damonmonli | 2 | 0.839 | 1 | **0.857** | 2 | 0.893 | 7 | 0.545 | 5 | 0.504 | 8 | 0.495 |
| YiSi-1 | 2 | 0.813 | 4 | 0.758 | 4 | 0.853 | 8 | 0.539 | 6 | 0.502 | 4 | 0.535 |
| PrismRefSmall | 2 | 0.850 | 3 | 0.786 | 4 | 0.854 | 11 | 0.532 | 7 | 0.498 | 3 | 0.541 |
| PrismRefMedium | 2 | 0.839 | 3 | 0.794 | 3 | 0.875 | 11 | 0.532 | 7 | 0.499 | 3 | 0.544 |
| BLCOM_1 | 2 | 0.827 | 3 | 0.779 | 1 | <u>0.909</u> | 7 | 0.545 | 6 | 0.500 | 1 | <u>0.552</u> |
| BLEURT-20 | 1 | 0.864 | 3 | 0.797 | 2 | 0.904 | 8 | 0.539 | 5 | 0.508 | 4 | 0.535 |
| COMET-22 | 2 | 0.811 | 5 | 0.714 | 2 | 0.906 | 6 | 0.557 | 4 | 0.517 | 1 | <u>0.552</u> |
| XCOMET | 2 | 0.850 | 2 | <u>0.832</u> | 1 | **0.924** | 5 | 0.566 | 3 | 0.527 | 1 | **0.558** |
| MetricX-24 (Hybrid) | 1 | **0.893** | 3 | 0.804 | 3 | 0.890 | 2 | <u>0.607</u> | 2 | <u>0.543</u> | 2 | 0.551 |
| METAMETRICS-MT (Hybrid) | 1 | <u>0.876</u> | 3 | 0.805 | 2 | 0.896 | 1 | **0.643** | 1 | **0.551** | 3 | 0.544 |
| **Reference-free** | | | | | | | | | | | | |
| sentinel-src-mqm | 5 | 0.522 | 7 | 0.491 | 8 | 0.570 | 11 | 0.532 | 8 | 0.497 | 12 | 0.197 |
| XLsimMqm | 5 | 0.592 | 6 | 0.506 | 8 | 0.574 | 9 | 0.535 | 8 | 0.497 | 10 | 0.420 |
| sentinel-cand-mqm | 5 | 0.595 | 6 | 0.581 | 7 | 0.681 | 5 | 0.565 | 5 | 0.505 | 9 | 0.445 |
| CometKiwi | 4 | 0.667 | 3 | 0.797 | 4 | 0.858 | 7 | 0.549 | 4 | 0.519 | 6 | 0.516 |
| bright-qe | 3 | 0.738 | 3 | 0.786 | 6 | 0.759 | 7 | 0.547 | 3 | 0.528 | 9 | 0.438 |
| XCOMET-QE | 3 | 0.740 | 3 | 0.806 | 3 | 0.868 | 10 | 0.534 | 5 | 0.504 | 6 | 0.514 |
| gemba_esa | 2 | <u>0.832</u> | 1 | **0.882** | 1 | **0.930** | 3 | <u>0.592</u> | 2 | **0.545** | 3 | <u>0.538</u> |
| MetricX-24-QE (Hybrid) | 2 | **0.853** | 2 | <u>0.814</u> | 2 | <u>0.907</u> | 3 | **0.597** | 3 | <u>0.529</u> | 2 | **0.548** |
| METAMETRICS-MT-QE | 3 | 0.768 | 4 | 0.749 | 3 | 0.878 | 4 | 0.585 | 4 | 0.522 | 7 | 0.505 |

Table 7: Detailed WMT24 result for language pair ja-zh. **Bold** and <u>underline</u> values indicate the best and second best performance, respectively.

| Metric | all | literary | news | social | speech |
|---|---|---|---|---|---|
| | | Task 1,5,9,13,17,20 | Task 2,6,10,14,18,21 | Task 3,7,11,16 | Task 4,8,12,16,19,22 |
| **Reference-based** | | | | | |
| sentinel-ref-mqm | 0.513 | 0.515 | 0.520 | 0.576 | 0.426 |
| BLEU | 0.589 | 0.618 | 0.626 | 0.645 | 0.488 |
| spBLEU | 0.593 | 0.629 | 0.632 | 0.650 | 0.570 |
| chrF | 0.606 | 0.635 | 0.637 | 0.663 | 0.579 |
| chrfS | 0.608 | 0.652 | 0.640 | 0.662 | 0.590 |
| BERTScore | 0.609 | 0.655 | 0.654 | 0.665 | 0.588 |
| MEE4 | 0.617 | 0.661 | 0.646 | 0.661 | 0.598 |
| damonmonli | 0.640 | 0.660 | 0.661 | 0.676 | 0.583 |
| YiSi-1 | 0.642 | 0.665 | 0.649 | 0.676 | 0.608 |
| PrismRefSmall | 0.646 | 0.672 | 0.646 | 0.686 | 0.608 |
| PrismRefMedium | 0.650 | 0.669 | 0.652 | 0.693 | 0.609 |
| BLCOM_1 | 0.684 | 0.680 | 0.651 | 0.714 | 0.658 |
| BLEURT-20 | 0.686 | 0.683 | 0.647 | 0.746 | 0.640 |
| COMET-22 | 0.695 | 0.689 | 0.653 | 0.757 | 0.663 |
| XCOMET | 0.719 | 0.696 | <u>0.669</u> | 0.765 | **0.678** |
| MetricX-24 (Hybrid) | <u>0.721</u> | <u>0.714</u> | 0.666 | <u>0.761</u> | <u>0.671</u> |
| METAMETRICS-MT (Hybrid) | **0.725** | **0.715** | **0.670** | **0.772** | 0.661 |
| **Reference-free** | | | | | |
| sentinel-src-mqm | 0.513 | 0.518 | 0.519 | 0.575 | 0.426 |
| XLsimMqm | 0.515 | 0.509 | 0.565 | 0.575 | 0.558 |
| sentinel-cand-mqm | 0.630 | 0.633 | 0.620 | <u>0.748</u> | 0.599 |
| CometKiwi | 0.635 | 0.622 | 0.643 | 0.695 | 0.623 |
| bright-qe | 0.664 | 0.634 | 0.653 | 0.706 | 0.639 |
| XCOMET-QE | 0.688 | 0.654 | 0.662 | 0.737 | <u>0.668</u> |
| gemba_esa | <u>0.711</u> | <u>0.694</u> | **0.679** | 0.734 | <u>0.668</u> |
| MetricX-24-Hybrid-QE | **0.714** | **0.697** | <u>0.666</u> | **0.751** | **0.683** |
| METAMETRICS-MT-QE | 0.681 | 0.641 | 0.641 | 0.717 | 0.660 |

Table 8: Detailed WMT24 results per domain category. **Bold** and <u>underline</u> values indicate the best and second best performance, respectively.

# chrF-S: Semantics is All You Need

**Ananya Mukherjee, Manish Shrivastava**
MT-NLP Lab, LTRC, KCIS, IIIT Hyderabad, India
ananya.mukherjee@research.iiit.ac.in
m.shrivastava@iiit.ac.in

## Abstract

Machine translation (MT) evaluation metrics like BLEU and chrF++ are widely used reference-based metrics that do not require training and are language-independent. However, these metrics primarily focus on n-gram matching and often overlook semantic depth and contextual understanding. To address this gap, we introduce chrF-S (Semantic chrF++), an enhanced metric that integrates sentence embeddings to evaluate translation quality more comprehensively. By combining traditional character and word n-gram analysis with semantic information derived from embeddings, chrF-S captures both syntactic accuracy and sentence-level semantics. This paper presents our contributions to the WMT24 shared metrics task, showcasing our participation and the development of chrF-S. We also demonstrate that, according to preliminary results on the leaderboard, our metric performs on par with other supervised and LLM-based metrics. By merging semantic insights with n-gram precision, chrF-S offers a significant enhancement in the assessment of machine-generated translations, advancing the field of MT evaluation. Our code and data will be made available at https://github.com/AnanyaCoder/chrF-S.

## 1 Introduction

In the rapidly advancing field of machine translation (MT), the need for robust and nuanced evaluation metrics has become increasingly critical. The evaluation landscape has expanded significantly in recent years, as evidenced by the WMT Metrics Shared Task, which provides a platform for meta-evaluating these metrics. Notably, in recent iterations of the WMT Metrics Shared Task, apart from learned metrics, lexical-based metrics such as BLEU (Papineni et al., 2002) and chrF (Popović, 2015, 2017) have consistently been regarded as baselines.

These metrics are widely appreciated for their language independence, which require no training and can be applied across diverse languages. However, they primarily address syntactic accuracy and often fall short in capturing the deeper semantic nuances and contextual relevance of translations.

The BLEU metric, with its reliance on modified precision of n-grams, provides a useful measure of how closely a machine-generated translation aligns with reference translations. Similarly, chrF enhances evaluation by incorporating character-level n-grams, offering greater sensitivity to morphological variations.

Despite these advancements, both metrics primarily focus on surface-level features, which can lead to incomplete assessments of translation quality, especially in complex linguistic contexts. To address these limitations, we propose **chrF-S** (Semantic chrF), an extension to the chrF++ metric, which leverages sentence embeddings to provide a more comprehensive evaluation by incorporating semantic analysis alongside traditional n-gram matching. Sentence embeddings (Reimers and Gurevych, 2019, 2020) encode entire sentences, thereby capturing the relationships between words, the structure of the sentence, and the broader contextual meaning. These offer rich semantic representations of sentences, enabling a deeper understanding of meaning and context. By merging these embeddings with chrF++'s character and word n-gram analysis, chrF-S aims to capture both the syntactic and semantic dimensions of translation quality.

This paper details our contributions to the WMT24 shared metrics task, where we have applied chrF-S to evaluate its effectiveness in comparison with existing metrics. We present our methodology of integrating semantic analysis into the chrF framework and discuss the preliminary results from the leaderboard, which indicate that chrF-S performs competitively with other super-

vised and LLM-based metrics. Our findings suggest that chrF-S not only enhances the evaluation of translation quality by incorporating semantic understanding but also provides a significant advancement over traditional metrics.

## 2 chrF-S

The main idea behind chrF-S is to have a combination of character-level match, word-level match and sentence-level match to provide a more comprehensive evaluation of translation quality. While chrF++ (Popović, 2015, 2017) already accounts for character and word-level matches, we enhanced this metric by introducing a sentence-level matching component. **We achieved this by adding a sentence-level component that utilizes sentence embeddings to compute a cosine similarity score, representing the semantic closeness between the reference and translation**. This flow is clearly illustrated in the figure 1. This approach allows chrF-S to assess not only the surface-level accuracy of the translation but also its deeper semantic fidelity, making it a more robust and nuanced evaluation metric.

For our experiments, we employed the LaBSE (Feng et al., 2022) model to generate these sentence embeddings. The ChrF-S score is computed as per equation 2

$$chrF\text{-}S(\text{ref, hyp}) = \alpha \cdot chrF + +(ref, hyp) + (1 - \alpha) \cdot CosSim\left(embed(ref), embed(hyp)\right)$$

In this equation, $ref$ refers to the reference sentence, and $hyp$ is the hypothesis (translation) sentence. $chrF + +(ref, hyp)$ denotes the character- and word-level similarity from ChrF++. The function $embed$ represents the sentence embeddings, which are generated using a sentence embedding model[1]. $CosSim$ computes the cosine similarity between the sentence embeddings. Finally, $\alpha$ is the weighting factor used to balance these two components; in our experiments, we set $\alpha = 2$.

## 3 Experiments

In our experiments, we considered two datasets released by WMT i.e., Direct Assessments (Bojar et al., 2017, 2018; Barrault et al., 2019, 2020; Akhbardeh et al., 2021; Kocmi et al., 2022) from 2017-2022 and MQM (Freitag et al., 2021a,b, 2022) assessments from 2020-2022. As the data is heavily skewed towards west-germanic languages.

---

[1] In this case, we used LaBSE

| lp | #segments | #systems |
|----|-----------|----------|
| en-de | 187 | 1 |
| en-ru | 250 | 1 |
| cs-uk | 3322 | 17 |
| en-es | 240 | 1 |
| en-zh | 120 | 1 |
| en-cs | 202 | 1 |
| en-ja | 242 | 1 |
| en-uk | 226 | 1 |
| en-is | 157 | 1 |
| en-hi | 247 | 1 |
| ja-zh | 243 | 1 |
| zh-en | 239 | 1 |
| de-en | 212 | 1 |
| es-en | 217969 | 356 |
| en-fr | 27730 | 161 |
| fr-en | 17553 | 117 |
| ru-en | 4320 | 18 |
| pt-en | 146508 | 158 |
| en-arz | 27333 | 156 |
| en-twi | 3560 | 16 |
| en-xho | 3285 | 18 |
| en-luo | 1251 | 6 |
| en-hau | 2996 | 14 |
| en-yor | 4774 | 28 |
| en-som | 46340 | 140 |
| yor-en | 25948 | 156 |
| en-kik | 18962 | 114 |
| ary-fr | 19960 | 120 |
| en-swh | 22954 | 138 |
| en-ibo | 19960 | 120 |
| Total | 617290 | 1865 |

Table 1: WMT24 Metrics Shared Task Test Set Statistics

| Test-set | #sentences | BLEU | BERTScore | chrF++ | chrF-S |
|----------|-----------|------|-----------|--------|--------|
| MQM-A | 457 | 0.210 | 0.333 | 0.478 | **0.481** |
| MQM-B | 790 | 0.180 | 0.282 | 0.410 | **0.423** |
| MQM-C | 1399 | 0.140 | 0.222 | 0.329 | **0.355** |
| MQM-D | 2425 | 0.117 | 0.188 | 0.272 | **0.313** |
| MQM-E | 4242 | 0.099 | 0.110 | 0.200 | **0.242** |

Table 2: Pearson Correlation scores on five different test sets curated from WMT-MQM (20-22) data

| Test-set | #sentences | BLEU | BERTScore | chrF++ | chrF-S |
|----------|-----------|------|-----------|--------|--------|
| DA-A | 8903 | 0.186 | 0.208 | 0.290 | **0.328** |
| DA-B | 17663 | 0.183 | 0.209 | 0.290 | **0.336** |
| DA-C | 34715 | 0.180 | 0.191 | 0.288 | **0.333** |
| DA-D | 67487 | 0.180 | 0.179 | 0.285 | **0.333** |
| DA-E | 126957 | 0.191 | 0.188 | 0.291 | **0.336** |

Table 3: Pearson Correlation scores on five different test sets curated from WMT-DA (17-22) data

Figure 1: chrF-S Metric

We created five sub testsets[2] of different sizes having a fair distribution of sentences across all language pairs. We evaluated these five testsets (A, B, C, D, E) using **unsupervised reference based metrics:** BLEU, chrF++, BERTScore and chrF-S and further computed pearson (Kurtz and Mayo, 1979) correlation to compare the metrics in terms of their agreement with human judgements.

## 3.1 Evaluation

Table 2 reports the correlation scores of the metrics with MQM assessments on the testsets built from WMT-MQM (20-22) data. Similarly Table 3 displays the correlation scores of the metrics with direct assessments on the testsets created from WMT-DA (17-22) data. In both the tables, it is clearly evident that chrF-S has performed better. We notice that the correlation scores of chrF-S in WMT-DA testset is slightly less (<0.4), however when compared to other metrics it still stands as winner.

By incorporating sentence-level embeddings, chrF-S enhances its ability to evaluate the semantic closeness between the reference and translation, leading to better alignment with human judgments that prioritize meaning and context. This semantic dimension improves correlation scores with human assessments, making chrF-S a more accurate and reliable metric, especially when translations **differ lexically but are semantically equivalent**.

## 4 WMT24 Metrics Shared Task Participation

We have participated in the WMT24 Metric Shared Task by submitting the translation scores for official evaluation (en-de, en-es, ja-zh) and secondary evaluation (for all langauge pairs from the generalMT task). The test-set statistics are reported in Table 1.

The preliminary leaderboard for the official language pairs is released by the shared task is reported at Table 4, displaying the system-level Pearson correlations and segment-level Kendall Tau correlations of of en-de, en-es and ja-zh language

---

[2]code and testsets will be released

| Rank | Participant | En-De sys-level Pearson | En-De seg-level Kendall | En-Es sys-level Pearson | En-Es seg-level Kendall | Ja-Zh sys-level Pearson | Ja-Zh seg-level Kendall |
|---|---|---|---|---|---|---|---|
| 1 | mengyao | 1.0 | 0.85 | 1.0 | 0.82 | 1.0 | 0.98 |
| 2 | jjuraska | 1.0 | 0.57 | 0.99 | 0.59 | 0.99 | 0.55 |
| 3 | gentaiscool | 1.0 | 0.67 | 0.98 | 0.69 | 0.99 | 0.61 |
| 7 | GEMBA-ESA | 0.98 | 0.53 | 0.99 | 0.51 | 0.94 | 0.49 |
| 8 | chrF-S | **0.97** | **0.51** | **0.99** | **0.5** | **0.97** | **0.56** |
| 12 | MetricsTaskBaseline | 0.95 | 0.45 | 0.92 | 0.46 | 0.5 | 0.17 |

Table 4: WMT24 Prelimnary Leaderboard reporting system-level and segment-level correlations. Our metric correlations are highlighted in bold.

pairs. It is noteworthy that chrF-S has not only surpassed the baseline but also demonstrated performance on par with GEMBA, an LLM-based metric. When compared to other preceding supervised metrics, chrF-S, an unsupervised metric proves to be competitive, standing alongside other top performers in the field.

## 5 Conclusion

This paper contributes to the WMT24 metrics shared task by introducing chrF-S, an enhanced version of chrF++ that incorporates sentence-level semantics for more accurate MT evaluation. Our metric effectively captures both surface accuracy and deeper semantic meaning by integrating character-level, word-level, and sentence-level matching. The use of sentence embeddings enables chrF-S to better assess semantic closeness between translations and references, leading to improved correlation with human judgments such as MQM and direct assessments. Preliminary leaderboard results indicate that chrF-S is competitive with other leading metrics, underscoring its potential as a reliable and nuanced tool for evaluating translation quality.

## Limitations

One significant limitation of this approach is its dependency on embedding models for sentence embeddings. The effectiveness of this method is restricted to languages for which appropriate sentence embedding models are available.

## References

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88, Online. Association for Computational Linguistics.

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 conference on machine translation (WMT20). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55, Online. Association for Computational Linguistics.

Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. Findings of the 2019 conference on machine translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy. Association for Computational Linguistics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. Findings of the 2017 conference on machine translation (WMT17). In *Proceedings of the Second Conference on Machine Translation*, pages 169–

214, Copenhagen, Denmark. Association for Computational Linguistics.

Ondřej Bojar, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, and Christof Monz. 2018. Findings of the 2018 conference on machine translation (WMT18). In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 272–303, Belgium, Brussels. Association for Computational Linguistics.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021b. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774, Online. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Albert K. Kurtz and Samuel T. Mayo. 1979. *Pearson Product Moment Coefficient of Correlation*, pages 192–277. Springer New York, New York, NY.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

# MSLC24: Further Challenges for Metrics
# on a Wide Landscape of Translation Quality

**Rebecca Knowles**          **Samuel Larkin**          **Chi-kiu Lo** 羅致翹

Digital Technologies Research Centre
National Research Council Canada (NRC-CNRC)
{rebecca.knowles,samuel.larkin,chikiu.lo}@nrc-cnrc.gc.ca

## Abstract

In this second edition of the Metric Score Landscape Challenge (MSLC), we examine how automatic metrics for machine translation perform on a wide variety of machine translation output, ranging from very low quality systems to the types of high-quality systems submitted to the General MT shared task at WMT. We also explore metric results on specific types of data, such as empty strings, wrong- or mixed-language text, and more. We raise several alarms about inconsistencies in metric scores, some of which can be resolved by increasingly explicit instructions for metric use, while others highlight technical flaws.

## 1 Introduction

This work builds on Lo et al. (2023), which introduced the Metric Score Landscape Challenge (MSLC).[1] At the Conference on Machine Translation (WMT), the Metrics Shared Task typically focuses on high-performing machine translation (MT) systems, in order to determine which new and improved metrics provide the most accurate and reliable scores (via comparison to human evaluation). However, the goal is for these metrics to go on to be used more broadly, which will likely result in their use on a wider range of systems. Since the Metrics Task primarily focuses on high-performing MT systems and their human evaluations, there is a risk that the new knowledge generated by the task about metrics may not generalize to lower-quality MT. For this reason, we submit a challenge set that covers a wider range of MT quality, in order to give potential users as well as metrics researchers a view of a broader range of performance. We also consider specific phenomena that may result in unexpected results from some metrics. We focus on three language pairs: English→Spanish

(eng→spa), English→German (eng→deu), and Japanese→Chinese (jpn→zho).

## 2 Data

We divide this MSLC into two subsets: the first challenge set (MSLC-A) follows the approach set out in MSLC23, merging together our low- to mid-quality systems with the systems submitted to the General MT shared task, while the second challenge set focuses on specific phenomena (MSLC-B; developed based on notable results from 2023 and new aspects of this year's General MT Task).

### 2.1 MSLC-A: News Data

We focus only on the "news" subset of the WMT General Task test set, as this better matches the domain of our trained MSLC systems and because of concerns with some of the other domains. All figures and values for MSLC-A will be shown over the subset of the "news" data that was manually evaluated with MQM (Multidimensional Quality Metrics; Lommel et al., 2013) by the Metrics Shared Task unless otherwise noted.

The MSLC-A systems we evaluate are a range of low- to medium-quality sets of MT output for the three identified language pairs.

The MT models we build for MSLC24 are all constrained (as per the WMT General Task rules) models, built using standard WMT training data (or subsets thereof), without the application of common additional techniques like backtranslation or tagging. We train all NMT models using Sockeye version 3.1.31 (Hieber et al., 2022), commit 13c63be5 with PyTorch 1.13.1 (Paszke et al., 2019).

The English→German systems are the same ones described in Lo et al. (2023); we direct the reader to that work for more details. The English→Spanish and Japanese→Chinese systems are described in more detail in Larkin et al. (2024). We use checkpoints from training the systems as

---

[1]MSLC data and additional figures can be found at https://github.com/nrc-cnrc/MSLC.

representative of varying levels of quality. The levels of quality are manually checked by authors familiar with the relevant target languages on a small sample of the data. We list the checkpoints used for the systems in Appendix A. The lowest-quality systems are indicated with the letter A, and the quality approximately increases as the system labels proceed alphabetically.

## 2.2 MSLC-B: Specific Phenomena

We target three specific phenomena in the MSLC-B challenge set: empty strings, mixed- and wrong-language text, and language variants. In addition to this, across these, we consider an overarching theme of consistency. We begin by describing and justifying our study of these phenomena and the topic of consistency.

Lo et al. (2023) observed unusual performance around empty strings (which appeared due to a submitted system's output in 2023). This may, at first glance, seem like a trivial and uninteresting issue. We argue that it is worth exploring, for three primary reasons: it is a real scenario that we observe in the WMT submissions and in more general MT (empty strings *do* appear in output and sometimes even input or references), it is important to know how metrics handle the empty string (as different metrics take different approaches to handling empty strings), and because of the question of consistency (some metrics may score empty strings in internally-inconsistent or surprising ways). It would be simple for all implementations of metrics to treat empty strings (in the source, reference, or hypothesis) as an edge case to be handled separately; in practice this is not what we observe, so it is important for users of metrics to be aware of how metrics may perform in these cases.

We also consider questions of how metrics perform when the MT output is mixed-language or wrong-language text. This is a situation that can arise, for example, due to noise in training data.

In a similar vein, since the General MT Task specified translation into Latin American Spanish, we build a very small test set of terms that differ between variants of Spanish spoken in Latin America and in Spain. For example, the word *computer* may commonly be translated as *ordenador* in Spain but *computadora* in Latin America. We use this to examine how metrics, particularly reference-free metrics, score translations from different language variants. This is a very small-scale study, but our results indicate that this is an area that should be considered for future work.

We now describe how we build this portion of the challenge set in order to study these issues.

### 2.2.1 English→German and English→Spanish

Here we produce a small data set to explore these issues more closely. We begin by selecting data that will be used repeatedly:

- 10 segments (paragraphs and sentences) from the English language source (WMT news data) with their Spanish and German[2] reference translations

- 10 short phrases in English with reference translations (confirmed via wikipedia, Linguee, and WordReference)[3]

- 10 words in English with reference translations

- 10 punctuation marks or other characters

Taking all of these, we consider the following situations: empty source and reference paired with the reference segments described above (simulating an MT system generating fluent text after empty string input) and empty string hypothesis paired with the known source and reference (simulating an MT system outputting the empty string).

Using only the segment (paragraph or sentence length) portion, we also consider the situation where the output is fluent but in the wrong-language by pairing the source with the correct reference but the opposite language hypothesis (e.g., English source, Spanish reference, German reference used as hypothesis). We also consider a mixed-language hypothesis, manually produced by substituting substrings of the Spanish reference with substrings from the German refB reference.[4] For German, because we have access to refB, we also submit a version with English source, refA as the reference, and refB as the hypothesis; this permits a full range from incorrect language to

---

[2] For the German references we use refA.

[3] https://en.wikipedia.org/, https://www. linguee.com/, https://www.wordreference.com/

[4] This was done in such a way to maintain (approximate) fluency and adequacy, such that a reader familiar with both German and Spanish should still be able to understand the text. We would urge caution, however, in assuming that any results from this part of the test set would extend to more natural code-mixing.

| Source | Last year, the World Economic Forum forecast that it would take five generations to achieve gender equality in every nation. Now the World Bank wants to rapidly accelerate that time frame. |
|---|---|
| Reference | Im vergangenen Jahr hat das Weltwirtschaftsforum vorausgesagt, es würde fünf Generationen dauern, bis in allen Staaten Geschlechtergleichstellung herrsche. Die Weltbank hat sich nun zum Ziel gesetzt, diesen Zeitraum deutlich zu verkürzen. |
| refB | Im vergangenen Jahr prognostizierte das Weltwirtschaftsforum, dass es fünf Generationen dauern werde, die Gleichstellung der Geschlechter in jeder Nation zu erzielen. Jetzt möchte die Weltbank diesen Zeitrahmen erheblich verkürzen. |
| Mixed-Lang. | *El año pasado,* prognostizierte das Weltwirtschaftsforum *que harían falta cinco* Generationen *para lograr la igualdad de género* in jeder Nation. Jetzt möchte die Weltbank *acelerar ese plazo rápidamente.* |
| Wrong Lang. | El año pasado, el Foro Económico Mundial pronosticó que harían falta cinco generaciones para lograr la igualdad de género en todas las naciones. Ahora, el Banco Mundial quiere acelerar ese plazo rápidamente. |

Table 1: Example of wrong-language, mixed-language (Spanish shown in italics), and refB (correct language alternate human reference) as hypotheses in the English→German MSLC-B dataset.

mixed-language to matched language (but different human translation). We show an example of this in Table 1.

For Spanish, since the WMT General MT Task explicitly describes this translation task as "EN to Spanish (Latin America)", we provide a very small sample (8 words) of words that tend to have differing translations between varieties of Spanish spoken in Latin America and varieties of Spanish spoken in Spain. This has very limited coverage but may permit us to begin asking questions about whether quality estimation systems have tendencies or biases towards certain language varieties.

### 2.2.2 Japanese→Chinese

For Japanese→Chinese, we examine metrics' performance around empty strings by first selecting data that will be used repeatedly:

- 5 segments (paragraphs and sentences) from the Japanese language source with their Chinese reference translations

- 5 short phrases in Japanese with reference translations

- 5 words in Japanese with reference translations

- 5 punctuation marks in Japanese with reference translations

We consider the same two types of empty string situations as in the other language pairs. The empty strings challenge examples make up 40 items in the MSLC-B Japanese→Chinese test set.

Similarly to the other language pairs, we consider wrong-language output (an English translation of the Japanese source, produced as a human translation from the Chinese reference by one of the authors) and mixed-language output (substituting words or phrases in the Chinese reference with corresponding Japanese and English words or phrases); these make up 10 items in the MSLC-B Japanese→Chinese test set.

## 3 Metrics

There are dozens of metrics submitted by the task organizers and participants to the WMT24 Metrics Shared Task. Given time and space limitations, we only examine the baseline metrics submitted by the task organizers and the primary metrics submitted by the participants. We describe the metrics included in this work in Appendix B.

## 4 Results and Plots

We divide our examination of the results into the two parts of the challenge set: MSLC-A and MSLC-B.

### 4.1 MSLC-A

Here we present preliminary results for the MSLC-A subset of the challenge set. We begin with the segment level and then consider system-level results. We make use of the MQM results provided by the Metrics Task organizers.

### 4.1.1 Segment Level

The histograms along the diagonal of Figure 1 show the distributions of segment-level scores produced by a subset of the baselines and submitted primary metrics. We can see that different metrics exhibit very different score distributions. Some show a somewhat bimodal distribution of scores, others are closer to normally distributed. For the metrics that are closer to normally distributed, we also see different skews. Most metrics are left skewed (i.e., they more frequently give segment scores in the higher-end of their possible score range), while BLEU is right skewed and more frequently gives segment scores in the lower end of its possible score range.

Metrics also differ in whether they exhibit a strong separation between the segments produced by the low-quality systems from our challenge set and the segments produced by the WMT submissions or whether they assign a range of low to high scores to most systems (i.e., having clear overlap in score range across all systems). This variation in characteristics suggests that metrics may have different strengths and weaknesses across the translation quality landscape; not all metrics are equally appropriate for scoring high-quality and low-quality MT.

*XCOMET* gives very low scores to segments from the very low-quality systems, but uses much more of the score space for the mid-quality systems. On the low-quality side, this is somewhat similar to the distribution of *BLEU* scores, but the high-quality systems have *XCOMET* scores that are much higher due to *XCOMET*'s bimodal distribution. Meanwhile, *chrF* shows a fairly normal distribution, but with a clear distinction between the various MSLC systems. We can also see this reflected when we examine system-level scores.

There are also metrics that use an approximation of a discrete score space, such as *GEMBA-ESA*. Lo et al. (2023) noted several metrics that did this in 2023; *GEMBA-ESA* is the only one in this year's set that does.

### 4.1.2 System Level

To analyze system-level scores, we compute an average over all of the segment-level scores in the news domain for a given MT system. There are two reasons why we are using this segment average instead of the submitted system-level score: 1) not all metrics submitted system-level scores and 2) using averaged segment-level scores allow us to show a representation of uncertainty (computed with bootstrap resampling, 1000 times, for $p < 0.05$) for the metrics. These system-level scores can also be used in order to gain a better understanding of the overall range of a metric's scores, as well as what kind of scores are assigned to very low quality machine translation (e.g., the A and B systems from the challenge set).

Figure 2 shows the system average scores for a subset of English→German (see Appendix D.1 for other translation directions). We observe that metrics show different patterns of scores at the system level. Both *PrismRefMedium* and *PrismRefSmall* appear to have serious difficulties in accurately scoring the lowest-quality system and give it a score higher than some of the better (still low-quality) systems.[5] Some metrics, such as *GEMBA-ESA*, *XCOMET* and *XCOMET-QE*, give very close scores to all of the low-scoring systems. For a use case (e.g., a low-resource language) where one expects to have low- to medium-quality systems at least initially, one may want to choose a metric that provides clearer distinctions between various systems on the lower range of quality.

For the high-quality systems the string-based metrics, such as *BLEU* and *chrF*, show wider error bars and thus may not distinguish between them. We leave analysis of the high-quality systems to the Metrics Shared Task.

By having our top MSLC system evaluated alongside the submitted WMT systems, we are able to observe that for Japanese→Chinese our systems combined with the high-performing submitted WMT systems do cover the wide range of quality. For English→German and English→Spanish, however, there may be a "missing middle" gap in quality that is not covered, an issue we aim to address in future work.

### 4.1.3 Conclusions: MSLC-A

As we saw in Lo et al. (2023), metrics differ in how they use their available score space. Some make fairly full use of their score range, others discretize the score space, and yet others display bimodal distributions of scores. All of these impact how individual segments are scored as well as how the system-level scores are distributed (i.e., whether the system-level scores are distributed more uniformly over the score space from low quality to

---

[5]Though we only do small-scale informal human evaluation, we expect, e.g., system E should not be ranked below system A.

Figure 1: Matrix of segment-level scores for English→German. Along the diagonal are stacked histograms of segment scores across the challenge set (cool colours/bottom) and submitted WMT systems (warm colours/top). The off-diagonal entries are scatterplots where each point is a single segment positioned according to the score assigned to it by row and column metrics; each point is coloured according to the same colours as the histogram. Note: for a full, scalable version of this figure, see https://github.com/nrc-cnrc/MSLC; all other figures in this paper are scalable.

Figure 2: System average scores for English→German. MSLC systems (cool colours, left) are ordered by BLEU score and brief manual examination; WMT submitted systems are ranked by average MQM score.

high quality, or whether most systems are clustered near the low and high ends of the score space). This year we noted fewer extremely unusual distributions; we did not see a repeat of the "universal scores" results observed in Lo et al. (2023).

The MQM evaluation of our top-performing system helps us to get a better idea of how to interpret these scores, though we note the issue of the missing middle range of human scores in two out of our three language pairs. We also note a weakness of error-based evaluations: they may not always capture non-errors (e.g., ways of translating that are not incorrect, but may be dispreferred by translators or end-users).

In future work, we may wish to apply more formal human evaluation to our lower-scoring systems, to better clarify the full range, but this year's introduction of human scores for one system per language pair takes a step towards that goal.

## 4.2 MSLC-B: Empty Strings

In MSLC23 (Lo et al., 2023), we observed a variety of system scores on empty strings produced by one of the participating systems in the WMT task. Here, we expand on that in a controlled fashion, examining the scores that metrics output when scoring empty strings.

### 4.2.1 Empty Source and Reference

We begin with empty source and reference, paired with four different types of output: single punc-

tuation characters (*punct*), single words (*word*), short phrases (*phrase*), and full sentences/short paragraphs (*sent*). All of the hypothesis text is in the target language, with the full sentences drawn from the WMT news data reference (refA, in the case of German). If these had been produced by an MT system taking an empty source and generating text, this might be considered a "hallucination"— generating fluent text that is not conditioned on any relevant source text. As such, we would expect that MT metrics should give these low scores. While some metrics (*BERTScore*, *BLEU*, *YiSi-1*, *chrF*, *spBLEU*, *mmm_qe*) do consistently give their lowest score (0) to all of these test segments, others show a greater variety of results.

Figure 3 shows a subset of the remaining metrics for English→German, covering a range of the variations in scores. Each subfigure shows the scores assigned to the 10 items in each category, with the vertical red lines indicating the lowest and highest scores assigned by this metric to any of the WMT news test data for any submitted MT system. *COMET-22* demonstrates the most common pattern: assigning a range of scores, with a tendency to have slightly higher scores for the shorter categories (e.g., punctuation—a single character has a very small edit distance to the empty string, perhaps making it more similar to the empty string than longer text) and lower scores to items in the longer categories (i.e., penalizing generating a full sentence out of nothing). *PrismRefMedium* and

Figure 3: English→German scores assigned to text when paired with empty source and reference. Red vertical lines indicate the minimum and maximum scores assigned over all WMT News primary submission data.

*chrfS* show another common pattern, by assigning low scores to all items; this is more in line with the desired and anticipated performance on this set of data. We note that the scores from *chrfS* are quite clustered around the lowest scores assigned to the WMT news data, while *PrismRefMedium* has scores expanding into a much lower range than the range of scores it assigned to the news data in the WMT test set. *MetricX-24-Hybrid* shows concerning results on this test set, assigning scores *higher* than any assigned to the WMT news test data to some of the samples, particularly the punctuation (perhaps not entirely unreasonably, as MT systems may need to occasionally generate additional punctuation in the target language), but also in some word and phrase examples. Finally, *GEMBA-ESA* assigns its lowest score most of the time, but occasionally assigns a top score or a score exactly in the middle of the range, an unexpected inconsistency.

### 4.2.2 Empty Hypothesis

Next, we flip the empty strings to the output side and pair them with real sources and references for the four different types of text mentioned in the

previous subsection. This is simulating the extreme case of omission where the complete output is missing. We understand that MT users may find it acceptable to omit translation for a single punctuation. As such, we again may expect that MT metrics would give gradually lower scores to the empty string output as the length of the source and reference increase. Similarly to the empty source and reference test cases, some metrics (*BERTScore*, *BLEU*, *YiSi-1*, *chrF*, *spBLEU*, *mmm_qe*) do consistently give their lowest score (0) to all of these test segments.

Figure 4 shows a subset of metrics for English→German, covering a range of the variations in scores. As we observe, *PrismRefMedium* and *chrfS* also give low scores (although not their lowest possible score) to empty string output. Some metrics (e.g. *COMET-22*) indeed give gradually lower scores to the empty string output according to the length of the input, with the items in the *sent* category receiving the low end of scores. We find that this is still relatively unsurprising behavior for metrics scoring empty string output. However, *MetricX-24-Hybrid* and *XCOMET* show concerning results on this test set, assigning mid-range to high scores to empty string output. Finally, as was the case in the empty source and reference test set, *GEMBA-ESA* assigns its lowest score most of the time, but occasionally assigns a top score to the empty string output.

### 4.2.3 Conclusions: Empty Strings

These empty string test cases (both empty source and reference and empty output) reveal undesirable metric results: giving high scores to extreme hallucination and omission. This leads us to be particularly concerned about the decision by the WMT General MT Task to use *MetricX-23-XL* and *CometKiwi-DA-XL* to decide which participating systems would receive human annotations, because related metrics (*MetricX-24-Hybrid* and *CometKiwi*) are two of the metrics showing these undesirable phenomena.

This may be an opening for a wider discussion about whether it is better for an MT system to fail to generate output than to generate output that is incorrect; nevertheless this would be a departure from past expectations (where, e.g., in human evaluation, "no translation" is typically given as a prototypical example of something that should receive a low score). In any case, we can likely find common ground in agreeing that metrics should not give

Figure 4: English→German scores assigned to the empty string paired with real source and reference.

high scores to non-empty output when given an empty input and empty reference. We would encourage a broader conversation about this, and in the meantime would encourage those presenting new metrics to be sure to specify how their metrics handle empty strings.

We encourage both metric builders and metric users to be aware of how metrics treat these edge cases. They do occur in practice and a user anticipating one type of performance on empty strings (e.g., low or 0 scores) may come to erroneous conclusions if they unknowingly use a metric that treats empty strings in another way (e.g., as high-scoring). We were also somewhat surprised to encounter the level of variation across empty string scores, and expect that users who are most familiar with string-matching metrics like BLEU may also not expect this variety of results.

### 4.3 MSLC-B: Mixed- or Wrong-Language

In this section, we explore what kind of output the metrics produce when they are applied to mixed-language and wrong-language translation hypotheses. We focus on English→German, because access to a second human-translated reference (refB),

allows us to explore a range of translation hypotheses, from a good human translation (refB, in German), a mixed-language translation (composed of a mix of the text of refB and the human-translated Spanish reference), and a wrong-language translation (the Spanish language reference). Our small test set for this is composed of 10 English source sentences, along with the various translations described above. We would expect a well-performing and usable metric to assign high scores to the refB German translation, lower scores to the mixed-language translation (portions of which are correct German translations of words and phrases), and even lower scores for the Spanish translation (a fluent and accurate translation, but in the wrong language).

However, this is not precisely what we observe in Table 2. Most metrics do give a score to refB that is greater than or equal to the score given to the mixed-language text in all 10 examples, while others score it at or above the mixed-language the majority of the time (*mmm_qe* (9), *CometKiwi* (8), and *damonmonli* (6)). Only *XLsimMqm* scores the mixed-language text higher than refB in 8 out of 10 examples. When it comes to the wrong-language text, most metrics again score refB equal or higher all of the time, but others at least occasionally rank the wrong-language text above refB— reference-free metrics in particular tend to make some errors (with the exceptions of *GEMBA-ESA* and *MetricX-24-Hybrid-QE*), but the two *PrismRef*\* metrics also make these errors. When comparing the scores given to the mixed-language text and the wrong-language text, we see even more of a mix. Some systems (both *PrismRef*\* systems, *XLsimMqm*, *mmm_qe*, and *MetricX-24-Hybrid-QE*) never score the mixed-language text above the wrong-language text.

#### 4.3.1 Conclusions: Mixed- and Wrong-Language

This varies somewhat between language pairs (see Appendix C), but string-based metrics like *BLEU* and *chrF* consistently score the mixed-language text above the wrong-language text. The weaknesses of string-based methods, such as their reliance on exact matches and lack of partial credit for synonyms (especially when evaluated against a single reference), have resulted in a shift towards embedding-based metrics that can provide more flexible semantic representations. However, given these results, it raises the question: are all modern

| Metric | refB≥Mix | refB≥Wrong | Mix≥Wrong |
|---|---|---|---|
| *BERTScore* | 10 | 10 | 8 |
| *BLEU* | 10 | 10 | 10 |
| *BLEURT-20* | 10 | 10 | 2 |
| *COMET-22* | 10 | 10 | 3 |
| *CometKiwi* | 8 | 4 | 1 |
| *PrismRefMedium* | 10 | 6 | 0 |
| *PrismRefSmall* | 10 | 7 | 0 |
| *YiSi-1* | 10 | 10 | 3 |
| *chrF* | 10 | 10 | 10 |
| *spBLEU* | 10 | 10 | 9 |
| *chrfS* | 10 | 10 | 10 |
| *MEE4* | 10 | 10 | 10 |
| *XLsimMqm* | 2 | 1 | 0 |
| *mmm_qe* | 9 | 6 | 0 |
| *mmm_hybrid* | 10 | 10 | 1 |
| *MetricX-24-Hybrid* | 10 | 10 | 2 |
| *MetricX-24-Hybrid-QE* | 10 | 10 | 0 |
| *GEMBA-ESA* | 10 | 10 | 10 |
| *XCOMET* | 10 | 10 | 1 |
| *XCOMET-QE* | 10 | 6 | 2 |
| *damonmonli* | 6 | 7 | 7 |

Table 2: eng→deu: Number of times (out of 10) that the metric scored `refB` higher than or equal to its mixed-language pair (refB≥Mix), higher than or equal to its wrong-language pair (refB≥Wrong), and a mixed-language hypothesis higher than or equal to its wrong-language pair (Mix≥Wrong).

metrics suitable for providing information about whether a text is a good translation *into the target language*, or simply whether it is a good translation (into some language(s))? We argue that these preliminary, small-scale results suggest the importance of additional analysis of this question. While this is unlikely to be a problem in many cases (especially when, e.g., language ID could also be performed), this may be particularly risky in low-resource settings where high-quality language ID is not available (cf. issues described in Kreutzer et al., 2022). In concurrent work, Zouhar et al. (2024) propose incorporating language ID to handle this issue as they explore it specifically in the context of *COMET*.

### 4.4 MSLC-B: Language Variants

The WMT General Task specifically called on researchers to build MT systems for English→Spanish using *Latin American* Spanish. We choose a small selection of terms that exhibit some of the vocabulary differences between the language variants of Spanish spoken in Latin America and Spain. We note several limitations to this: this is a very small set of terms, the terms are evaluated in isolation, and they are certainly not fully representative of all Spanish language variants spoken in Latin America or Spain.[6]

Due to the structure of the challenge set submission process, each source term was submitted four times: once for each language variant with the matching reference and once for each language

---

[6]In several of the cases presented here, there exist a number of other translations that we could have selected.

variant with the opposite reference. Considering only the reference-free metrics (those that do not use the provided reference in order to compute their score), we observe results in Table 3. A checkmark(✓) indicates that in all cases for that term, the Latin American term chosen was scored higher than the term used more commonly in Spain; an ✗ indicates that the term used in Spain scored equal to or higher than the term used in Latin America. This somewhat arbitrary choice to include repeated versions was fortuitous, because it highlights a concern with one of the metrics: a question mark (**?**) in a cell indicates that the rankings computed were mixed. This means that, on repeated scoring, the variations within the metric scores returned were great enough (in the case of *MetricX-24-Hybrid-QE*) to result in different rankings at least once. This should be alarming to potential users of metrics, who would expect consistent results on repeated strings. That is: a user may reasonably expect that if they submit the same input to a metric twice in a row, they will get the same output twice in a row; here we observe that not all metrics have this as a guarantee. We discuss this more in Section 5. We note that other metrics also exhibited some variation in their scores, but the rest did not vary enough to change which of the two term variants received the higher score.

We now observe that *XLsimMqm* is the only metric to prefer the term used more commonly in Latin America more than half the time (5/8). We note that *GEMBA-ESA* only prefers the term from Latin America in 1/8 terms, but for the remaining 7 terms, both variants are given identical scores of 100 (*GEMBA-ESA* is the only metric whose counts would change, were we to score it as correct if a term used in Latin America scores *equal or greater than* the term used more commonly in Spain).

This raises similar questions to those we considered in the wrong-language and mixed-language experiments, albeit at a finer-grained level. Metrics may not be equally appropriate for use across all language variants, and may in fact demonstrate a scoring preference to one or the other. This will require considerably more experimentation, with larger test sets, in the future.

### 5 Consistency

Our experiments in MSLC-B highlighted some issues in metric score consistency: repeated instances of scoring the same string resulting in dif-

| Metric | English source: computer / Latin America: computadora / Spain: ordenador | sandwich / sándwich / bocadillo | potato / papa / patata | juice / jugo / zumo | waiter / mesero / camarero | tires / llantas / neumáticos | peanut butter / mantequilla de maní / crema de cacahuete | drive / manejar / conducir | Counts |
|---|---|---|---|---|---|---|---|---|---|
| CometKiwi | ✓ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | 4/8 |
| XLsimMqm | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ | ✓ | 5/8 |
| mmm_qe | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | 4/8 |
| MetricX-24-Hybrid-QE | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ? | ✓ | 4/8 |
| GEMBA-ESA | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | 1/8 |
| XCOMET-QE | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | 2/8 |

Table 3: Metric preferences for terms that are more common in Spanish language variants spoke in Latin America (✓), for terms more common to language variants spoken in Spain (✗), or inconsistent preferences (?).

ferent scores. The MSLC-B dataset provided only a small set of examples on which to test this, and only for reference-free metrics. However, because we submitted our highest-scoring MSLC MT systems to the WMT General MT task, we actually do have a larger set of data with which to explore repeated scores. This consists of the intersection between the MSLC-A data (news domain only) and the full General MT test set (149 segments for eng→deu and eng→spa and 269 for jpn→zho), each of which was scored by each metric as part of the MSLC-A challenge set and as a WMT submitted system.

It is important to highlight that, while the Metrics Shared Task calls for metrics that provide scores at the segment level and the system level, it does not currently specify how or when metrics may make use of extrasentential information (e.g., information from other parts of the test set or document) in order to produce segment level scores. This could include approaches that compute some statistics from the full test set (like *YiSi* does for TF-IDF) or that operate on the batch level (like *PrismRef\**). This is to say, some of these apparent inconsistencies may be intentional (i.e., giving a segment a different score depending on the context in which it appears).

Since metrics have different score ranges, we first calculate the lowest and highest scores assigned by each metric to any of the MSLC-A or MSLC WMT submission segments within the news subset. This gives us a range of metric scores. Then, for each source-reference-hypothesis in the news subset, we compute the absolute difference in the score that it was assigned as part of the MSLC-A dataset and the score it was assigned as part of the MSLC WMT submission dataset and express this as a percentage of the metric's score range described above.

For many metrics (*BERTScore*, *BLEU*, *chrF*, *spBLEU*, *chrfS*, *XLsimMqm*, *mmm_qe*, *mmm_hybrid*, *GEMBA-ESA*), there is never any difference in these two scores. For other metrics, like

*CometKiwi*, there are some small differences (never greater than 0.1% of the metric's score range); these seem likely attributable to rounding/floating point errors. In other cases, it is possible that other even larger differences may be accounted for due to differences in batch size and hardware used, such as the case of *MetricX-24-Hybrid-QE*, which sees its largest score difference as 7.3% of the score range for eng→deu.

For *YiSi-1*, there is a known reason for the observed differences (up to 2.3% eng→deu, 2.8% eng→spa, 4.4% jpn→zho): the *YiSi* score is computed taking into account TF-IDF statistics from the full test set; since MSLC-A included only the news data while the full WMT General Task submission for MSLC included other domains, the scores assigned to individual segments may differ, as the segment-level scores are conditioned on the full test set. The largest difference we observe is for *PrismRefMedium*, with one score difference of 98.9% of the full score range; this is likely also due to the model operating at the level of the document or document chunk. The MSLC-A challenge set did not include document boundaries, which could account for the differences we observe. In future tasks, we would suggest incorporating document information in the challenge set submission in order to avoid these issues, and it would also be helpful to clarify which metrics incorporate extrasentential information (specifically from other parts of the challenge set data). We know that there are at least three different levels on which metrics are operating: the single-segment level (i.e., each segment is scored individually, so repeated segments should be scored identically), incorporating information from the full test set (in which case repeated segments within the same test set may receive identical scores), and incorporating document/batch/multi-segment input (in which case, scores may depend on how the batching is performed). It could also be possible to have more complex interactions (e.g., taking into account where in a document a segment occurs in order to score it); metrics users and challenge set

484

builders need to be aware of these in order to ensure that they are measuring what they think they are measuring.

As we can see, there are at least two different reasons for these apparent inconsistencies: 1) purposeful differences that arise from metrics that use contextual information for computing sentence-level scores (as in the case of *YiSi*) and 2) errors and noise resulting from computational or implementational factors. In the case of these purposeful differences, the primary thing for metric users to be aware of is the scope of the context that is used, in order to be able to reproduce scores. The latter issue is a larger problem, especially when we see score differences that cover substantial portions of the metric's score range. If a metric is unstable or produces different scores based on the hardware used to compute it, we face an issue at least as concerning as the preprocessing one identified in Post (2018). We propose two main (but not entirely satisfactory) solutions to this: 1) it may be best to report such metrics as an average over multiple runs and 2) metrics should adopt the proposals outlined in Zouhar et al. (2024) to include metric signatures for better reproducibility.

## 6   Conclusion

We once again show the diversity of ways that metrics perform on a wide range of system quality. We also observe quite a bit of variation in terms of how systems handle empty strings, which may influence how they are used (e.g., when comparing a system that frequently generates empty strings to one that never does). We also consider questions of wrong-language text and mixed-language text as well as language variants, and argue that metrics researchers should consider whether their metrics are overgeneralizing (i.e., whether they give high scores to good translations regardless of whether the translation is in the desired target language or not) or are biased towards particular language variants. Many of our results support the conclusions that Zouhar et al. (2024) describe in their concurrent analysis of *COMET*, such as the need to better handle empty strings, questions of target language, biases, and the importance of metric signatures when metric variations may introduce score differences. In concert with that work, we raise the concern that as new metrics are introduced, we are not learning the lessons from our field's past errors. We argue for the importance of examining real-world corner cases and issues of reproducibility in order to more responsibly introduce new metrics to the research community. Both metrics researchers and users should be alarmed by the levels of inconsistency that we observe. One of the benefits of using automatic metrics should be to make fair comparisons (for repeated scoring, across papers, and so forth)—inconsistent metrics cannot serve this purpose. When there are intentional sources of differences in scores for repeated segments (i.e., due to the context in which they appear), users need to be aware of the scope and approaches used to incorporate context, in order to ensure that they are using metrics as intended in order to measure what they intend to measure. This will become increasingly important as we see a shift to document-level

## Limitations

We focus only on three language pairs (English→German, English→Spanish, and Japanese→Chinese) in the News domain in this work, due to the availability of human-annotated scores for this set. Several of our additional experiments use extremely small sets of data (e.g., 5-10 examples); in most cases these are designed to help us establish whether additional future study would be helpful, rather than to make definitive claims about the results. In time for the camera-ready submission, we had access to MQM scores, but not to the General MT Task ESA annotations.

## Acknowledgements

## References

Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are llms breaking mt metrics? results of the wmt24 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Felix Hieber, Michael Denkowski, Tobias Domhan, Barbara Darques Barros, Celina Dong Ye, Xing Niu, Cuong Hoang, Ke Tran, Benjamin Hsu, Maria Nadejde, Surafel Lakew, Prashant Mathur, Anna Currey, and Marcello Federico. 2022. Sockeye 3: Fast neural machine translation with pytorch. *arXiv*, abs/2207.05851.

J. D. Hunter. 2007. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.

Samuel Larkin, Chi-kiu Lo, and Rebecca Knowles. 2024. MSLC24 submissions to the general machine translation task. In *Proceedings of the Ninth Conference on Machine Translation*. Association for Computational Linguistics.

Chi-kiu Lo, Samuel Larkin, and Rebecca Knowles. 2023. Metric score landscape challenge (MSLC23): Understanding metrics' performance on a wider landscape of translation quality. In *Proceedings of the Eighth Conference on Machine Translation*, pages 776–799, Singapore. Association for Computational Linguistics.

Arle Richard Lommel, Aljoscha Burchardt, and Hans Uszkoreit. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Vilém Zouhar, Pinzhen Chen, Tsz Kin Lam, Nikita Moghe, and Barry Haddow. 2024. Pitfalls and outlooks in using comet.

# A MSLC24 MT Systems

In Table 4 we see the checkpoint IDs for systems included in the challenge set for eng→deu. Table 5 and 6 show the same for eng→spa and jpn→zho.

| **System** | Checkpoints | *BLEU* |
|---|---|---|
| A | 54 | 0.50 |
| B (50k) | 1 | 1.85 |
| C | 79 | 3.13 |
| D (50k) | 7 | 4.19 |
| E (200k) | 2 | 4.54 |
| F | 91 | 6.88 |
| G (200k) | 27 | 7.87 |
| H (400k) | 4 | 8.73 |
| I (400k) | 43 | 9.64 |
| J | 102 | 9.24 |
| K | 129 | 13.91 |
| L | 313 | 22.79 |
| M (MSLC) | 311 | 22.65 |

Table 4: Checkpoint IDs and *BLEU* scores (nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1) on MSLC-A for systems included in challenge set (eng→deu); parenthetical numbers indicate one of the pseudo-low-resource systems rather than the full training data system.

| **System** | Checkpoints | *BLEU* |
|---|---|---|
| A | 52 | 0.75 |
| B | 65 | 4.94 |
| C | 74 | 8.55 |
| D | 84 | 13.14 |
| E | 98 | 19.91 |
| F | 123 | 25.61 |
| G | 207 | 31.47 |
| H (MSLC) | 800 | 37.97 |

Table 5: Checkpoint IDs and *BLEU* scores (nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1) on MSLC-A for systems included in challenge set (eng→spa).

| System | Checkpoints | BLEU |
|---|---|---|
| A | 37 | 0.05 |
| B | 70 | 4.97 |
| C | 80 | 10.74 |
| D | 97 | 15.79 |
| E | 133 | 19.48 |
| F (MSLC) | 358 | 23.12 |

Table 6: Checkpoint IDs and *BLEU* scores (`nrefs:1|case:mixed|eff:no|tok:zh|smooth:exp|version:2.3.1`) on MSLC-A for systems included in challenge set (`jpn→zho`).

| Metric Name | Reference-based |
|---|---|
| *Human annotation* | |
| MQM | |
| *Metrics* | |
| *BERTScore* | ✓ |
| *BLEU* | ✓ |
| *BLEURT-20* | ✓ |
| *chrF* | ✓ |
| *chrfS* | ✓ |
| *COMET-22* | ✓ |
| *CometKiwi* | |
| *damonmonli* | ✓ |
| *GEMBA-ESA* | |
| *MetaMetrics-MT* | ✓ |
| *MetaMetrics-MT-QE* | |
| *MEE4* | ✓ |
| *MetricX-24-Hybrid* | ✓ |
| *MetricX-24-Hybrid-QE* | |
| *prismRefMedium* | ✓ |
| *prismRefSmall* | ✓ |
| *sentinel-cand-mqm* | |
| *sentinel-ref-mqm* | ✓ |
| *sentinel-src-mqm* | |
| *spBLEU* (flores-200) | ✓ |
| *XCOMET* | ✓ |
| *XCOMET-QE* | |
| *XLsimMqm* | |
| *YiSi-1* | ✓ |

Table 7: Human annotation and metrics included in this work, with their coverage of language pairs. Metrics that are not marked as reference-based are reference-free (a.k.a quality estimation) metrics.

## B  Metrics

Table 7 shows a summary of the human annotations and metrics included in this work and the translation directions they participated in. For detail descriptions of the metrics, please refer to the Metrics Task overview paper (Freitag et al., 2024).

Note: in the main body of the text, for space reasons, we abbreviate the *MetaMetrics-MT-QE* and *MetaMetrics-MT* names as *mmm_qe* and *mmm_hybrid*, respectively.

| Metric | Mix≥Wrong |
|---|---|
| *BERTScore* | 10 |
| *BLEU* | 10 |
| *BLEURT-20* | 5 |
| *COMET-22* | 4 |
| *CometKiwi* | 3 |
| *PrismRefMedium* | 9 |
| *PrismRefSmall* | 8 |
| *YiSi-1* | 10 |
| *chrF* | 10 |
| *spBLEU* | 10 |
| *chrfS* | 10 |
| *MEE4* | 10 |
| *XLsimMqm* | 7 |
| *MetaMetrics-MT-QE* | 2 |
| *MetaMetrics-MT* | 1 |
| *MetricX-24-Hybrid* | 0 |
| *MetricX-24-Hybrid-QE* | 0 |
| *GEMBA-ESA* | 10 |
| *XCOMET* | 1 |
| *XCOMET-QE* | 0 |
| *damonmonli* | 6 |

Table 8: `eng→spa`: Number of times (out of 10) that the metric scored a mixed-language hypothesis higher than or equal to its wrong-language pair.

## C  Additional Mixed/wrong-language Tables

Tables 8 and 9 show the how the metrics scores mixed-language and wrong-language data for English→Spanish and Japanese→Chinese. For English→Spanish, the wrong-language text was German and the mixed-language was a mix of German and Spanish. For Japanese→Chinese, the mixed-language was a mix of Chinese, English and Japanese, while the wrong-language was English. Note that because the Chinese text in the mixed-language hypotheses is drawn directly from the reference, this should be a particularly easy task for string-based metrics.

## D  Additional Figures

Figures in this paper are produced using Matplotlib (Hunter, 2007), version 3.7.1.

### D.1  MSLC-A System-Level

Figures 5, 6, and 7 show the system average scores for English→German, English→Spanish, and Japanese→Chinese across all metrics.

Figure 5: System average scores for English→German.

Figure 6: System average scores for English→Spanish.

Figure 7: System average scores for Japanese→Chinese.

| Metric | Mix≥Wrong |
|---|---|
| *BERTScore* | 5 |
| *BLEU* | 5 |
| *BLEURT-20* | 5 |
| *COMET-22* | 5 |
| *CometKiwi* | 0 |
| *PrismRefMedium* | 0 |
| *PrismRefSmall* | 0 |
| *YiSi-1* | 5 |
| *chrF* | 5 |
| *spBLEU* | 5 |
| *chrfS* | 5 |
| *MEE4* | 5 |
| *XLsimMqm* | 5 |
| *MetaMetrics-MT-QE* | 0 |
| *MetaMetrics-MT* | 4 |
| *MetricX-24-Hybrid* | 1 |
| *MetricX-24-Hybrid-QE* | 0 |
| *GEMBA-ESA* | 0 |
| *XCOMET* | 0 |
| *XCOMET-QE* | 0 |
| *damonmonli* | 3 |

Table 9: jpn→zho: Number of times (out of 5) that the metric scored a mixed-language hypothesis higher than or equal its wrong-language pair.

## D.2 Remaining Additional Plots

For other examples of of the empty string plots, as well as for additional plots showing histograms and scatterplots, see https://github.com/nrc-cnrc/MSLC.

# MetricX-24: The Google Submission to the WMT 2024 Metrics Shared Task

**Juraj Juraska, Daniel Deutsch, Mara Finkelstein** and **Markus Freitag**
Google
{jjuraska,dandeutsch,marafin,freitag}@google.com

## Abstract

In this paper, we present the MetricX-24 submissions to the WMT24 Metrics Shared Task and provide details on the improvements we made over the previous version of MetricX. Our primary submission is a hybrid reference-based/-free metric, which can score a translation irrespective of whether it is given the source segment, the reference, or both. The metric is trained on previous WMT data in a two-stage fashion, first on the DA ratings only, then on a mixture of MQM and DA ratings. The training set in both stages is augmented with synthetic examples that we created to make the metric more robust to several common failure modes, such as fluent but unrelated translation, or undertranslation. We demonstrate the benefits of the individual modifications via an ablation study, and show a significant performance increase over MetricX-23 on the WMT23 MQM ratings, as well as our new synthetic challenge set.[1]

## 1 Introduction

Automatic evaluation metrics are critical to the development of machine translation (MT) systems. In recent years, the landscape of MT evaluation has changed dramatically since the use of lexical metrics, like BLEU (Papineni et al., 2002) and ChrF (Popović, 2015), that compared the tokens or characters of the candidate translation to a reference translation to predict a scalar score that represents the quality of the translation. Evaluation metrics based on neural networks opened up the door for more experimentation, and metrics now vary based on what type of output they produce, what they require as input for prediction, and whether they use a dedicated evaluation model or a general-purpose large language model.

This paper provides details on MetricX-24, the successor to MetricX-23. MetricX is a learned

regression-based metric trained to predict a floating point score representing the quality of a translation. This year, we made four submissions to the WMT24 Metrics Shared Task, all based on the mT5 language model (Xue et al., 2021), which is further fine-tuned on direct assessment (DA) ratings, MQM ratings (Lommel et al., 2014; Freitag et al., 2021), and newly constructed synthetic data. The primary submission, denoted MetricX-24-Hybrid, is a hybrid reference-based/-free metric, which can score a translation irrespective of whether it is given the source segment, the reference, or both. The same model is thus the primary submission for both the reference-based evaluation and the quality estimation (QE) task, having predicted the scores once with and once without the references provided in the input. Our contrasting submissions, MetricX-24(-QE), are standalone reference-based/QE models, trained only for their specific task.

The key takeaways from our experiments, detailed in this report, include:

1. Learned metrics cannot reliably detect under-translation, duplication, missing punctuation, and fluent but unrelated translation;
2. Adding a relatively small amount of synthetic data to the training set can boost the metric's performance, especially on lower-quality translations with the above issues;
3. It is possible to effectively train a metric on a mixture of MQM and DA ratings, thus maintaining high performance on a larger set of language pairs;
4. Training a metric in the hybrid input mode, i.e., with and without the reference included in the input, allows it to learn to rely less on the reference when it is of poor quality.

## 2 Data

Developing MetricX-24, we relied solely on publicly available data from the WMT Metrics shared

---

[1] Our code and models can be found at https://github.com/google-research/metricx.

tasks between 2015 and 2023. The translation ratings from these years come in two different flavors: (1) direct assessment (DA) scores on a scale from 0 to 100, collected in general from non-expert raters, and (2) MQM scores (Lommel et al., 2014; Freitag et al., 2021) on a scale from 0 to 25 (with 0 being the best), which are grounded on error spans and their corresponding severity levels, annotated by professional raters. MQM ratings have been collected as part of the WMT campaign only since 2020 and, because the annotations are considerably more time-consuming and expensive to obtain, they are only available for a few language pairs. The DA scores, on the other hand, offer a broader language coverage of nearly 50 language pairs, but the raw ratings are noisy (due to different rating strategies) and generally of lower quality. Therefore, it is often beneficial to $z$-normalize DA ratings per rater before training models on them, so as to make the ratings more comparable across different annotators. In contrast, models do not benefit from MQM scores being $z$-normalized because the scores come from a rather small group of annotators and they adhere to a rubric.

In the rest of this section, we provide details on which data we use for training and evaluation, as well as how the different datasets are preprocessed. Furthermore, we describe new synthetic data we created from the WMT datasets, with the goal of addressing some of MetricX's known failure modes.

## 2.1 Training Data

**DA.** We utilize most of the DA data from the 2015–2022 period for training, with the following exceptions. As we observed during the development of the previous version of MetricX (Juraska et al., 2023), the into-English portion of the WMT21 DA ratings drags the model performance down. We confirmed this observation again this year and excluded these language pairs from the training data. With the gradually declining quality of DA ratings collected for WMT using the MTurk platform, we also exclude all into-English language pairs from WMT22.[2] Additionally, we exclude the en-zh language pair from WMT22, as we use the equivalent slice of data, but with MQM ratings, for evaluation. We use $z$-normalized ratings when training models on DA data only, but raw ratings

when training on a mixture of MQM and DA data.

**MQM.** Besides the DA ratings, we also take advantage of the higher-quality MQM ratings from the years up to 2022 for training. These include four language pairs: en-de, en-ru, en-zh and zh-en.[3] We only use the conversation, e-commerce and social domains from WMT22 en-zh for training. In our experiments with different subsets of MQM ratings, we observed a consistent boost in performance with the 2020 data excluded, hence, our final models are only trained on MQM ratings from 2021 and 2022. We always train models on raw MQM ratings, i.e., using the 0–25 scale.

## 2.2 Evaluation Data

**MQM.** Our primary evaluation set consists of the WMT23 MQM ratings, which includes three language pairs: en-de, he-en and zh-en. Since the zh-en language pair is known to have low-quality references (Kocmi et al., 2023), we replace them with newly collected references. Note that this has no effect on the MQM ratings, as those were collected in a source-based fashion. Additionally, given the fact that one of the official WMT24 test language pairs is ja-zh, we reserve the news domain subset of the WMT22 en-zh ratings for evaluation, allowing us to assess our models' performance on a language pair with Chinese as the target language.

**DA.** We use the WMT23 DA ratings as a secondary evaluation set, taking advantage of its better language coverage (8 language pairs). Nevertheless, with DA ratings generally following a significantly different distribution than MQM ratings, a higher correlation of the metric scores with these DA ratings does not necessarily imply better performance. For example, fine-tuning a model on zh-en MQM ratings results in lower performance than fine-tuning it on DA ratings, according to the zh-en DA evaluation set (but not the MQM one). Therefore, we only consider the WMT23 DA evaluation set in experiments where we mix MQM and DA training data together.

## 2.3 Synthetic Data

After seeing the initial benefits from the simple synthetic data used for training MetricX-23, we decided to construct a more comprehensive collection

---

[2] One exception is zh-en, for which DA ratings were collected in two different ways, including using the same method and framework as the out-of-English language pairs (Kocmi et al., 2022).

[3] The en-zh MQM ratings, available at `https://github.com/google/wmt-mqm-human-evaluation`, were collected post-WMT22.

of synthetic training examples. They cover additional, less trivial failure modes of MetricX, i.e., translation issues commonly unrecognized by the metric. The DEMETR challenge set (Karpinska et al., 2022), which we relied on last year, does not cover several of the failure modes we created the new synthetic training examples for, hence we also constructed a set of test examples for each of them. Next, we describe how we designed both the training and the test synthetic datasets.

### 2.3.1 Training Sets

In order for the MetricX models to learn to identify certain types of bad translations that are not sufficiently (or at all) represented in the regular WMT training data, we generated synthetic examples that we augment the training data with. They were created by modifying examples from the DA datasets ranging from WMT15 to WMT22, comprising 49 language pairs. Table 1 provides an overview of the various failure modes that we considered, including brief descriptions of how we prepared the synthetic data to address them. Additional details regarding the creation process can be found in Appendix A.

### 2.3.2 Test Set

We constructed a new DEMETR-style test set based on the WMT23 DA dataset, with examples generated analogously to our synthetic training examples, as described in Table 1. Each synthetic example is paired with its original counterpart (although using the reference instead of the candidate translation whenever the synthetic translation was created from the reference), which allows for a metric to be evaluated on how frequently it ranks the pairs correctly.

## 3 Metric Descriptions

The MetricX-24 submissions to the WMT24 Metrics Shared Task build on top of the successful MetricX-23 (Juraska et al., 2023; Kocmi et al., 2023), with several major improvements. We start this section by summarizing the aspects this year's submissions have in common with MetricX-23, then provide an overview of the modifications, and finally describe the differences between the individual submissions.

### 3.1 MetricX Model

MetricX is a learned metric, powered by a regression model trained to predict a floating point number that represents the quality of a given transla-

tion. The reference-based variant takes the candidate translation (hypothesis) and reference segments as input, and concatenates them, along with corresponding prefixes ("candidate:" and "reference:", respectively). In contrast to the previous versions, MetricX-24 also prepends the source segment (along with the prefix "source:") to the input, offering the model additional context to make a better prediction in the reference-based setting, which may be beneficial especially in cases where the reference is inadequate. The model then encodes this combined input and uses it to predict the translation quality score. The QE variant works in an analogous way, but taking only the source segment and the hypothesis as the input.

With MetricX-24, we continue to rely on mT5 (Xue et al., 2021) as the pretrained language model that we fine-tune on translation evaluation data. We refer the reader to Juraska et al. (2023) for details on how we adapted this encoder-decoder model to the regression task. Similar to MetricX-23, we fine-tune the model in two stages: first on DA ratings ($z$-normalized, aggregated per segment, negated, and finally clipped to the $[-1.0, 1.0]$ range) and then further on raw MQM ratings. As a result, the metric produces scores in the [0, 25] range. The model is trained with a mean squared error (MSE) loss function. Further implementation details can be found in §4.

### 3.2 Design Improvements

We achieve some initial improvement in performance by simply including the WMT22 data in the training set – both the DA and the MQM ratings, which we previously used as the evaluation set when developing MetricX-23. The additional MQM ratings (including en-ru, a language pair not present in the older MQM data) are especially valuable, considering the scarcity of MQM data. Besides that, we introduce three major modifications to the training procedure and data in order to further improve MetricX's performance, described throughout the rest of this section.

### 3.2.1 Training With Synthetic Data

Although we used synthetic training data alongside the DA and MQM ratings already for training MetricX-23, the synthetic examples covered only the two trivial cases of empty and reference-matching translations. As described in §2.3, we prepared a significantly more comprehensive synthetic training set for MetricX-24, which we combine

| Failure mode | Synthetic candidate translation description | MQM score |
|---|---|---|
| Empty translation | Empty string. | 25 |
| Gibberish | Text of a similar length as the reference, generated by sampling words from the vocabulary built from all references in the data with a matching target language. | 25 |
| Fluent but unrelated translation | Arbitrary reference from the dataset of a similar length and in the same language. | 25 |
| Undertranslation | Candidate translation with an arbitrary sentence removed, if a multi-sentence segment, otherwise, candidate translation with 20–80% words removed from the end. | 5–25 |
| Duplication | Candidate translation duplicated, with a space in between. | 25 |
| Missing punctuation | Reference translation with the end punctuation removed (11 punctuation symbols considered, such as period, question mark, closing parenthesis or quotation mark). | 1 |
| Reference-matching translation | Reference translation itself (unlike the rest, these synthetic examples are meant to train the metric to predict a perfect score for translations matching the reference). | 0 |

Table 1: Failure mode categories we prepared synthetic data for, along with brief descriptions of how we created the synthetic examples from the WMT data, and the MQM scores we label the training examples with.

with the DA and MQM data in both fine-tuning stages. We experimented with various ratios, and settled on 1:100 for each synthetic example category in the first stage and 1:5000 in the second stage. We evaluate the effects of adding the synthetic training data by measuring accuracy and average score differences on the synthetic test set, also described in §2.3.

### 3.2.2 Mixing DA and MQM Data

Next, we attempt to address the inevitable decline in MetricX performance on other languages after fine-tuning the model on MQM data, which only covers a few language pairs. The performance, as measured by the WMT23 DA evaluation set with 8 language pairs, quickly declines after starting to fine-tune on MQM ratings. While it is expected that the change in the general score distribution – caused by the switch from DA to MQM ratings – results in the Pearson correlations with the ground-truth scores dropping, we believe the model should be able to retain its system- and segment-level pairwise accuracy from the first stage of fine-tuning on DA data. Moreover, we observe a significant drop in system-level performance on the zh-en language pair of the MQM evaluation set, despite zh-en being present in the MQM training data.

In order to remedy these behaviors, we mix in a smaller proportion of DA ratings in the second-stage fine-tuning. That way the model is trained primarily on MQM ratings, but has a continued exposure to the additional 40+ language pairs from the first stage of fine-tuning. We experimented with different combinations of DA and MQM rating formats (e.g., raw vs. $z$-normalized, transformed to the MQM scale or not, etc.), and the one yielding

the best results was raw MQM ratings combined with raw DA ratings linearly transformed to the MQM scale of $[0, 25]$. Finally, we determined that a DA:MQM ratio of 1:4 works well for boosting the performance on the DA evaluation set back to the levels from the first stage of fine-tuning, without a significant negative impact on the model's performance on the MQM evaluation set.[4]

### 3.2.3 Hybrid Input Mode

The third major modification we make to the training procedure when developing MetricX-24, is mixing training examples in three different formats: (1) source + hypothesis, (2) hypothesis + reference, and (3) source + hypothesis + reference. This allows the model to operate in both a QE and a reference-based mode (and the latter either with or without the source included). But perhaps more importantly, it gives the model an opportunity to learn how much weight to put on the source and the reference in different scenarios, or possibly to completely ignore the reference when it is of low quality. Such a hybrid model is then evaluated as a QE model by only passing it the source segment and the hypothesis as input, and as a reference-based model by additionally passing it the reference.

### 3.3 MetricX-24 Variants

There are four variants of MetricX-24 that we submitted to the WMT24 Metrics Shared Task:

- MetricX-24-Hybrid (primary)
- MetricX-24-Hybrid-QE (primary)
- MetricX-24

---

[4]After a more extensive post-submission experimentation, we determined the optimal ratio to be 1:10.

- MetricX-24-QE

Our primary reference-based and QE submissions are actually the same hybrid model, with the scores predicted with and without the references provided as part of the input. The secondary submissions are the standalone reference-based and QE counterparts of the hybrid model, i.e., only trained on examples with the references (as well as the source segments) included and on examples with the references omitted, respectively. Other than that, all of the submission models are identical in terms of training data mixtures, as described in §3.2.1 and §3.2.2, as well as training hyperparameters.

## 4 Experimental Setup

### 4.1 Meta-Evaluation

As mentioned in §2.2, our primary evaluation set consists of the MQM ratings from WMT23, as well as the news domain subset of the en-zh language pair from WMT22. Considering there is no into-English language pair among the official test sets this year, we focus primarily on en-de and en-zh when evaluating our models, but also keeping zh-en (the dataset with alternate references) in the mix, in order to ensure that we do not overfit the models to out-of-English language pairs. To evaluate our models, we calculate the agreements between their predicted scores and the human judgments of translation quality using the four different correlations from the WMT23 Metrics Shared Task (Freitag et al., 2023), detailed in Appendix B.

### 4.2 Checkpoint Selection

In both the first and the second stage of fine-tuning, we pick the best checkpoint $c_{\text{best}}$ based on the following linear combination of segment- and system-level pairwise accuracy:

$$\arg\max_c \frac{3}{4} \sum_l \text{acc}_l^{\text{seg}}(c) + \frac{1}{4} \sum_l \text{acc}_l^{\text{sys}}(c) \, ,$$

where $l \in \{\text{en-de, en-zh, zh-en}\}$, and $\text{acc}_l^{\text{seg}}(c)$ and $\text{acc}_l^{\text{sys}}(c)$ are the segment- and the system-level pairwise accuracy calculated for checkpoint $c$ on the language pair $l$ of the evaluation set. We down-weight the system-level component to account for its greater variance and to thus avoid a checkpoint being picked due to a rare spike in system-level accuracy if segment-level accuracy is low.

### 4.3 Implementation Details

MetricX-24, similar to its predecessor, is implemented with TensorFlow (Abadi et al., 2015) and the T5X library (Roberts et al., 2023). All of the metric variants are based on mT5-XXL with 13B parameters. We defer further implementation details to Appendix C. We are publicly releasing our submissions, converted from TensorFlow to PyTorch (Paszke et al., 2019) checkpoints.[5]

## 5 Results and Discussion

Here we present the results of our experiments, focusing solely on fully trained models (i.e., those that went through both stages of fine-tuning) and modifications in the second stage. Since the ablation studies performed with reference-based and QE models show similar trends, we discuss the reference-based experiments in depth in this section, and provide the QE results in Appendix D.2. Due to limited resource availability, we were only able to run each experiment with one random seed.

### 5.1 Training With Synthetic Data

We start by examining the benefits of including synthetic training examples, as described in 2.3. In Table 2, the bottom four rows – corresponding to the hybrid model – demonstrate the effects of progressively adding DA data only, synthetic data only, and finally both, to the training set in the second stage of fine-tuning.[6] We ended up not using the duplication synthetic training set, as we observed that the models learn to correctly identify such cases even without it.

The first thing to notice is that mixing in DA ratings actually improves the metric's performance on the synthetic test set over fine-tuning on MQM ratings alone, especially in the unrelated, undertranslation and duplication failure modes. Adding synthetic data instead is, however, significantly more effective in general, boosting the accuracy to the 94–100% range in most categories. Finally, augmenting the training set with both the DA and the synthetic data results in an overall similar performance as with the synthetic data only.

Missing punctuation is one of two categories in which our metrics score not so close to perfect. In fact, the synthetic training examples appear not to be helpful in improving the performance at all. Our

---

[5] https://github.com/google-research/metricx
[6] The models that did not include synthetic training data in the second stage, did not use it in the first stage either.

| MetricX variant | +DA | +Synth | Empty transl. | Gibberish | Unrelated | Undertransl. | Duplication | Missing punct. | Refmatch |
|---|---|---|---|---|---|---|---|---|---|
| 23 | – | ∼ | **100.00** | **100.00** | 88.14 | 57.75 | 38.14 | 66.01 | **94.00** |
| 24 | ✓ | ✓ | 99.29 | 99.86 | **99.29** | **98.75** | 99.14 | 83.01 | 78.14 |
| 24-Hybrid | – | – | 51.43 | 99.86 | 81.00 | 68.75 | 87.57 | 83.66 | 76.00 |
| | ✓ | – | 53.57 | 99.71 | 92.14 | 82.25 | **99.57** | **85.62** | 72.86 |
| | – | ✓ | 94.14 | 99.71 | 99.14 | 96.25 | 94.43 | 84.97 | 79.86 |
| | ✓ | ✓ | 97.29 | 99.71 | 98.71 | 96.25 | 99.43 | 82.35 | 75.14 |

Table 2: Accuracy of reference-based MetricX variants in all 7 categories of our synthetic test set. "23" is the baseline, the last row of "24-Hybrid" corresponds to our primary submission, and "24" is our secondary submission.

| MetricX variant | +DA | +Synth | Segment-level pairwise accuracy | | | | System-level pairwise accuracy | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | en-de | zh-en | zh-en† | en-zh | en-de | zh-en | zh-en† | en-zh |
| 23 | – | ∼ | 60.20 | 53.12 | 54.06 | 55.73 | 90.91 | 89.52 | 86.67 | 74.36 |
| 24 | ✓ | ✓ | 60.71 | 54.50 | 55.78 | 56.16 | 96.97 | **92.38** | **95.00** | **88.46** |
| 24-Hybrid | – | – | 61.17 | 54.63 | 55.52 | 57.43 | **100.00** | 89.52 | 91.67 | 85.90 |
| | ✓ | – | 60.75 | 54.89 | 55.58 | 57.65 | 98.48 | **92.38** | 92.50 | 84.62 |
| | – | ✓ | **61.75** | 54.38 | 55.43 | **57.73** | 98.48 | 90.48 | 91.67 | **88.46** |
| | ✓ | ✓ | 61.11 | **55.00** | **55.82** | 57.02 | 98.48 | **92.38** | 94.17 | 85.90 |

Table 3: Meta-evaluation scores of reference-based MetricX variants on the WMT23 MQM evaluation set. "23" is the baseline, the last row of "24-Hybrid" corresponds to our primary submission, and "24" is our secondary submission. †Alternate references.

hypothesis is that using references to create this category of synthetic examples results in a significant proportion of misleading examples because we assume references to be perfect, but that is not always the case. That, combined with the fact that the removal of the punctuation symbol from the end of the segment warrants just a minor score change, means that some of the synthetic examples might have an unreasonably high ground-truth score associated with them, thus giving the model the opposite signal to what is desired.

The reference-matching translation synthetic training set appears not to be effective either, however, its benefits are somewhat concealed by the fact that mixing in DA data drags the performance in this category down. With the non-hybrid model, we observed a significantly bigger drop with DA data included ($77\% \rightarrow 64\%$) and a greater increase with synthetic data included instead ($77\% \rightarrow 83\%$). Granted, that is still far from perfect, however, expecting a 100% accuracy in this category equates to expecting that the candidate translation is never better than the reference, which, as we pointed out earlier, is not always true when judging the translation quality based on the source segment.

Overall, thanks to the new synthetic training data, MetricX-24 (hybrid or not) is clearly more robust to the failure modes than MetricX-23 (see first row

in the table), with the reference-matching translations being an exception. That might have to do with the absence of WMT22 data in the training set of MetricX-23, or the only synthetic examples present therein being those of empty and reference-matching translations.

## 5.2 Mixing DA and MQM Data

We already discussed the effects of adding DA data to the training set in the second stage of fine-tuning in terms of the synthetic test set performance; let us now have a look at the correlations with human MQM scores. Comparing the first two rows of the "24-Hybrid" section in Table 3, we see that there are just relatively minor changes in either direction across all language pairs, with score differences within the expected variance between runs.

What the table does not show, however, is the huge jump in all correlations across all language pairs on the WMT23 DA evaluation set, typically back to the levels from the first stage of fine-tuning on DA data only, or above. Segment- and system-level pairwise accuracy increases by up to 2 and 5 points, respectively, and Pearson's $r$ sees improvements of up to 10 points. These are valuable gains, considering we achieved them without sacrificing the performance on the MQM evaluation set. An overview of the results and a more detailed

analysis on the DA evaluation set can be found in Appendix D.1.

## 5.3 Hybrid Input Mode

To wrap up the evaluation, we discuss the performance difference between MetricX-24 and MetricX-24-Hybrid (rows 2 and 6 in Table 3). At the system level, the hybrid variant lags slightly behind in zh-en and en-zh, but it makes up for it by outperforming the non-hybrid across the board at the segment level. Notably, the hybrid metric achieves an almost 1% higher segment-level accuracy on en-zh, and the 0.5% boost on zh-en (with original references) may be evidence for the hybrid model handling examples with poor-quality references better, especially considering the accuracy difference on the zh-en set with alternate references is only 0.04%. The other performance differences between the two models are largely insignificant.

Finally, comparing our primary submission with MetricX-23 (row 1 in the table), we can see consistent gains of 1–2 points in segment-level accuracy, and substantially bigger gains at the system level, with the accuracy on en-zh improving by a whopping 11.5 points. We conclude that this a significant improvement over our last year's submission, ranked overall second in the WMT23 Metrics Shared Task.

## 6 Related Work

Traditionally, evaluation metrics predict a scalar quality score for the translation. This type of metric includes BLEU, ChrF, MetricX (Juraska et al., 2023), BLEURT (Sellam et al., 2020; Pu et al., 2021), COMET (Rei et al., 2020, 2022a), COMETKiwi (Rei et al., 2022b), Prism (Thompson and Post, 2020), and more. While these metrics have historically been the dominant category of metric, newly proposed methods provide structured (Perrella et al., 2022; Fernandes et al., 2023; Kocmi and Federmann, 2023; Guerreiro et al., 2023) or natural language explanations (Xu et al., 2023) for the predicted scores.

Then, evaluation metrics are considered to be reference-based or reference-free (also known as "quality estimation") depending on whether or not they require a reference to evaluate a translation. Metric developers usually train separate models for each type of metric (e.g., COMET and COMETKiwi, or MetricX-23 and MetricX-23-QE), but some opt for combining both tasks into a single

model (Wan et al., 2022; Guerreiro et al., 2023), which is the approach we took in this work with our hybrid model.

Finally, while most metrics like MetricX-24 use a dedicated model for scoring translations, some recent works have begun to leverage general-purpose large language models instead (Fernandes et al., 2023; Kocmi and Federmann, 2023; Xu et al., 2023; Leiter et al., 2023; Leiter and Eger, 2024). While LLM-based metrics have achieved strong system-level performance, using a learned dedicated model was the best approach at the segment-level in last year's Metrics Shared Task (Freitag et al., 2023).

## 7 Conclusion

We presented in detail our approach to training MetricX-24, a regression-based MT evaluation metric. We submitted four versions of MetricX-24 to the WMT24 Metrics Shared Task, including a reference-based and a QE variant, as well as a new hybrid variant evaluated with and without the references. By evaluating on the WMT23 MQM dataset, we showed all of them to significantly outperform our last year's submission, MetricX-23. In addition, we made MetricX-24 more robust to various types of bad translations, which do not frequently occur in the WMT data, such as undertranslation, or fluent but unrelated translation. Finally, by combining DA and MQM ratings together in the final stage of fine-tuning, we were able to dramatically increase the performance on the WMT23 DA dataset covering 8 language pairs, while maintaining the high correlations with the MQM ratings at the same time.

## References

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.

Daniel Deutsch, George Foster, and Markus Freitag. 2023. Ties matter: Meta-evaluating modern metrics

with pairwise accuracy and tie calibration. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12914–12929, Singapore. Association for Computational Linguistics.

Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André Martins, Graham Neubig, Ankush Garg, Jonathan Clark, Markus Freitag, and Orhan Firat. 2023. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1066–1083, Singapore. Association for Computational Linguistics.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. XCOMET: Transparent machine translation evaluation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Marzena Karpinska, Nishant Raj, Katherine Thai, Yixiao Song, Ankita Gupta, and Mohit Iyyer. 2022. DEMETR: Diagnosing Evaluation Metrics for Translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9540–9561, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.

Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To Ship or Not to Ship: An Extensive Evaluation of Automatic Metrics for Machine Translation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 478–494, Online. Association for Computational Linguistics.

Christoph Leiter and Steffen Eger. 2024. PrExMe! Large Scale Prompt Exploration of Open Source LLMs for Machine Translation and Summarization Evaluation. *arXiv preprint arXiv:2406.18528*.

Christoph Leiter, Juri Opitz, Daniel Deutsch, Yang Gao, Rotem Dror, and Steffen Eger. 2023. The Eval4NLP 2023 Shared Task on Prompting Large Language Models as Explainable Metrics. In *Proceedings of the 4th Workshop on Evaluation and Comparison of NLP Systems*, pages 117–138, Bali, Indonesia. Association for Computational Linguistics.

Arle Lommel, Hans Uszkoreit, and Aljoscha Burchardt. 2014. Multidimensional quality metrics (mqm): A framework for declaring and describing translation quality metrics. *Tradumàtica*, (12):0455–463.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor

Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library.

Stefano Perrella, Lorenzo Proietti, Alessandro Scirè, Niccolò Campolungo, and Roberto Navigli. 2022. MaTESe: Machine translation evaluation as a sequence tagging problem. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 569–577, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. Learning compact metrics for MT. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 751–762, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Adam Roberts, Hyung Won Chung, Gaurav Mishra, Anselm Levskaya, James Bradbury, Daniel Andor, Sharan Narang, Brian Lester, Colin Gaffney, Afroz Mohiuddin, et al. 2023. Scaling up models and data with t5x and seqio. *Journal of Machine Learning Research*, 24(377):1–8.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning Robust Metrics for Text Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR.

Brian Thompson and Matt Post. 2020. Automatic Machine Translation Evaluation in Many Languages via Zero-Shot Paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 90–121, Online. Association for Computational Linguistics.

Yu Wan, Keqin Bao, Dayiheng Liu, Baosong Yang, Derek F. Wong, Lidia S. Chao, Wenqiang Lei, and Jun Xie. 2022. Alibaba-translate China's submission for WMT2022 metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 586–592, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Wenda Xu, Danqing Wang, Liangming Pan, Zhenqiao Song, Markus Freitag, William Wang, and Lei Li. 2023. INSTRUCTSCORE: Towards Explainable Text Generation Evaluation with Automatic Feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5967–5994, Singapore. Association for Computational Linguistics.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

500

## A  Synthetic Data Creation

We sample 500 examples from each language pair, whose candidate translations (hypotheses) we then manipulate in different ways to create the synthetic examples for each failure mode category. The missing punctuation category is an exception, with a stratified sample across the 11 end-punctuation symbols, rather than language pairs, and 250 examples each.

In general, the synthetic examples have the candidate translation manipulated, turning it into a worse, or an outright bad, translation. One exception is the reference-matching category, whose purpose is to actually teach the metric to score translations that match the reference highly, which it does not learn to do reliably when only trained on the WMT data. Table 4 shows a few concrete examples from the synthetic training set.

## B  Meta-Evaluation Details

**System-Level.** At the system level, we measure pairwise ranking accuracy (Kocmi et al., 2021), as well as Pearson's $r$. Pairwise accuracy assesses how well a metric ranks MT systems by calculating the proportion of all possible pairs of systems that are ranked the same by the metric and human scores. Pearson's $r$, on the other hand, captures how strong the linear relationship is between the metric and human scores for MT systems. We obtain the system-level scores (both metric and human) as the mean segment-level score for each system.

**Segment-Level.** At the segment level, we use the group-by-item pairwise accuracy with tie calibration, as described by Deutsch et al. (2023), and the no-grouping Pearson's $r$. The pairwise accuracy calculates the proportion of all possible pairs of translations for the same source segment that are ranked the same by the metric and human, then averages the accuracies over all input segments. At the same time, it rewards correct tie predictions by introducing ties for any two translations with a metric score difference below an automatically determined threshold. The no-grouping Pearson's $r$ quantifies the linear relationship between the metric and human scores across all translations from every system and document.

## C  Implementation Details

Having increased the maximum segment length from 256 to 512 SPM tokens, and including up to three segments (source, hypothesis and reference) in the model's input, each training run requires 256 TPUs. Using a batch size of 256, we train our models for 16K steps in the first stage, using a learning rate of 0.001 with an inverse square root decay after the first 2K steps. We then fine-tune the best checkpoint for another 8K steps in the second stage, lowering the learning rate to 0.0002 and decaying it after 1K steps. The models are trained using the Adafactor optimizer (Shazeer and Stern, 2018).

## D  Additional Results

### D.1  Mixing DA and MQM Data

Table 5 compiles the results of the meta-evaluation of a group of reference-based models on the WMT23 DA evaluation set. All of the models are standalone reference-based models. In the table, we contrast four variants of the model fine-tuned in two stages (DA then MQM data) with a model fine-tuned on DA data only (i.e., the first stage only). We present the results on a subset of four language pairs, two of which are present in our MQM training data (en-de and zh-en) and two which are not (en-cs and de-en).

The experiments with mixing DA and MQM data in the second stage of fine-tuning were motivated by the large differences in performance on the WMT23 DA evaluation set observed between a model trained on DA ratings only (row 1 in Table 5) and the same model further fine-tuned on MQM ratings (row 2). As already discussed in §2.2, this can be partly explained by the discrepancy in DA and MQM rating distributions. This discrepancy understandably affects Pearson correlations, however, it should not have a significant effect on how the metric ranks segments or systems. Nevertheless, while we observed large drops in Pearson's $r$, the pairwise accuracy also dropped substantially for most of the language pairs, both at the segment and the system level. For example, on en-cs the segment-level accuracy drops from $59.54$ to $57.43$, and the system-level accuracy from $87.62$ to $82.86$.

Considering the fact that the performance difference between the models in rows 1 and 2 on en-de and zh-en (i.e., the language pairs with a good amount of MQM training data), are relatively

**Gibberish (zh-en example)**

*Created from: corpus hypothesis vocabulary*

src    我希望你们能准时，不是想要你们的优惠券！！

hyp    filter two that to also in allegations train 800 city, continuous the

ref    I hope you can be on time, and it's not that I want your coupons! !

label   25

---

**Fluent but unrelated translation (de-en example)**

*Created from: corpus references*

src    Damit können doppelt so viele Studierende ausgebildet werden wie bisher.

hyp    She booked a return flight and went home the next day.

ref    In that way, twice as many students can be educated as before.

label   25

---

**Undertranslation (cs-en example)**

*Created from: hypothesis*

src    Dlouhodobě napjaté vztahy mezi oběma zeměmi se vyostřily v roce 2018 poté, co Washington odstoupil od jaderné dohody z roku 2015 mezi Íránem a světovými mocnostmi a zavedl vůči Íránu sankce, které mají tvrdý dopad na jeho ekonomiku.

hyp    Long-tense relations between the two countries sharpened in 2018 after Washington withdrew from the 2015 nuclear deal between Iran and world powers and imposed sanctions.

ref    Long-term tense relations between both countries escalated in 2018 after that Washington withdrew from the nuclear deal closed in 2015 between Iran and the world powers and imposed sanctions against Iran, which have had hard impacts on its economy.

label   12.75

---

**Duplication (fi-en example)**

*Created from: hypothesis*

src    Ensi vuoden vaje on yli 2,4 prosenttia kansantuotteesta.

hyp    Next year's deficit will be over 2.4 per cent of national product. Next year's deficit will be over 2.4 per cent of national product.

ref    Next year's deficit is over 2.4 per cent of GDP.

label   15

---

**Missing punctuation (ru-en example)**

*Created from: reference*

src    Последний альбом Ace вышел в 2016 году.

hyp    Their last album, "Ace", came out in 2016

ref    Their last album, "Ace", came out in 2016.

label   1

---

**Reference-matching translation (ja-en example)**

*Created from: reference*

src    グレタさんは、27日の金曜日にも行うことを呼びかけていた。

hyp    Now, Greta is calling for further strikes to be held on Friday the 27th.

ref    Now, Greta is calling for further strikes to be held on Friday the 27th.

label   0

---

Table 4: Synthetic examples for the different failure mode categories (except for the trivial empty translation case), along with the MQM scores we label the training examples with. Each category also has an indication of how the hypothesis was created/generated in order to produce a synthetic example (e.g., by modifying the original hypothesis or reference).

| MetricX variant | +DA | +Synth | Segment-level pairwise accuracy | | | | System-level pairwise accuracy | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | en-de | zh-en | en-cs | de-en | en-de | zh-en | en-cs | de-en |
| DA only | N/A | N/A | 61.77 | 56.33 | 59.54 | 61.14 | **95.45** | 79.05 | **87.62** | 92.31 |
| DA then MQM | – | – | 61.59 | 55.99 | 57.43 | 61.65 | 93.94 | 81.90 | 82.86 | 85.90 |
| | ✓ | – | 61.88 | **56.67** | **60.16** | 62.29 | **95.45** | 80.95 | 86.67 | 88.46 |
| | – | ✓ | **62.60** | 56.35 | 59.02 | 61.92 | **95.45** | **84.76** | 81.90 | **93.59** |
| | ✓ | ✓ | 61.89 | 56.64 | 60.04 | **62.32** | 93.94 | 83.81 | 86.67 | **93.59** |

Table 5: Meta-evaluation scores of reference-based MetricX variants on a subset of the language pairs of the WMT23 DA evaluation set. "DA only" is a model after just the first stage of fine-tuning (i.e., on DA data only), whereas the "DA then MQM" section contains models fine-tuned in full two stages. The last row thus corresponds to the "24" row in Tables 2 and 3, i.e., our secondary submission "MetricX-24".

| MetricX variant | +DA | +Synth | Segment-level Pearson's $r$ | | | | System-level Pearson's $r$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | en-de | zh-en | en-cs | de-en | en-de | zh-en | en-cs | de-en |
| DA only | N/A | N/A | **60.10** | **41.52** | **43.47** | 52.59 | 98.48 | **89.21** | **92.49** | 97.20 |
| DA then MQM | – | – | 48.18 | 34.66 | 39.77 | 44.09 | 93.41 | 87.58 | 90.60 | 87.40 |
| | ✓ | – | 53.67 | 36.29 | 43.11 | 52.79 | 93.67 | 87.75 | 90.40 | 91.53 |
| | – | ✓ | 56.03 | 35.21 | 37.24 | 45.26 | **99.48** | 88.40 | 89.32 | 96.79 |
| | ✓ | ✓ | 57.92 | 37.17 | 42.55 | **55.23** | 98.56 | 88.94 | 91.58 | **97.97** |

Table 6: Same as Table 5, but showing Pearson correlations instead of pairwise accuracies.

small, we conjecture that further fine-tuning on MQM data alone causes the model to partially "forget" other languages from the first stage of fine-tuning. We attempt to prevent the model from this sort of forgetting by mixing some DA ratings into the training set in the second stage.

As the scores of the model in row 3 in the table demonstrate, we are able, for the most part, to restore the performance observed in the first stage of fine-tuning by adding a small proportion of DA training data in the second stage too. Adding not only the DA data, but also the synthetic data, in the second stage (row 5) sometimes boosts the performance further, significantly improving even over the first-stage performance (row 1). Most importantly, the gains over fine-tuning on MQM data alone (row 2) are achieved not at the expense of the model's performance on the MQM or the synthetic test set, as evidenced by the results in Tables 2 and 3.

Finally, Table 6 shows the expected big drops in Pearson correlation with the DA ratings after fine-tuning on MQM data (see rows 1 and 2), especially at the segment level. Adding DA data in the second stage helps recover most of the performance (compare rows 3 and 5 with row 1), but as expected, the correlations remain lower particularly for language pairs present in the MQM data the model is fine-tuned on in the second stage (en-de and zh-en).

### D.2 QE Models

In Tables 7 and 8, we present the meta-evaluation results for our QE models. These are analogous to those presented in §5, only the hybrid model is evaluated in a reference-free mode, and the non-hybrid models are ones trained on the source and hypothesis segments only. Note that the hybrid model is the same checkpoint as the one for which we reported the reference-based results in Tables 2 and 3, i.e., not one optimized for QE performance.

Examining first the results on the synthetic test set, summarized in Table 7, we see similar trends to those observed with reference-based models (Table 2). The main difference is that the QE models achieve significantly lower performance in the missing punctuation and the reference-matching translation categories. This, however, is expected because both the types of synthetic examples were created from references. In case of the missing punctuation examples, the synthetic translation is simply the reference with the end punctuation removed. Comparing such a hypothesis with the corresponding reference is arguably a significantly easier task than comparing it to the source segment and identifying a missing punctuation symbol. Moreover, there may be a mismatch in the presence of punctuation between the source and the reference in the training examples, making it even more difficult for a QE model to reliably identify missing punctuation. As for the reference-matching

| MetricX variant | +DA | +Synth | Empty transl. | Gibberish | Unrelated | Under-transl. | Duplication | Missing punct. | Ref-match |
|---|---|---|---|---|---|---|---|---|---|
| 23 | – | ~ | **100.00** | **99.86** | 96.43 | 63.25 | 88.29 | 69.93 | 63.00 |
| 24 | ✓ | ✓ | 97.86 | **99.86** | **99.43** | **98.50** | **98.14** | 65.36 | 63.43 |
| 24-Hybrid | – | – | 69.86 | **99.86** | 82.43 | 81.25 | 63.00 | **77.78** | 63.00 |
| 24-Hybrid | ✓ | – | 66.14 | 99.57 | 95.29 | 93.50 | 97.86 | 73.86 | 62.57 |
| 24-Hybrid | – | ✓ | 93.57 | 99.71 | 99.29 | 96.50 | 84.43 | 69.28 | 62.14 |
| 24-Hybrid | ✓ | ✓ | 93.71 | **99.86** | **99.43** | 97.25 | **98.14** | 69.28 | **64.14** |

Table 7: Accuracy of reference-free (QE) MetricX variants in all 7 categories of our synthetic test set. "23" is the baseline, the last row of "24-Hybrid" corresponds to our primary submission, and "24" is our secondary submission. The hybrid model is the same as in Table 2, only evaluated without references provided as input.

| MetricX variant | +DA | +Synth | Segment-level pairwise accuracy | | | | System-level pairwise accuracy | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | en-de | zh-en | zh-en† | en-zh | en-de | zh-en | zh-en† | en-zh |
| 23 | – | ~ | 59.57 | 52.64 | 52.89 | 54.47 | 92.42 | 86.67 | 85.83 | 74.36 |
| 24 | ✓ | ✓ | 59.70 | **54.30** | **54.48** | 56.00 | 98.48 | **92.38** | 90.83 | **87.18** |
| 24-Hybrid | – | – | 60.11 | 53.80 | 54.00 | **56.27** | **100.00** | 89.52 | 89.17 | 84.62 |
| 24-Hybrid | ✓ | – | 59.18 | 54.08 | 54.30 | 56.14 | **100.00** | **92.38** | 90.00 | 84.62 |
| 24-Hybrid | – | ✓ | **60.27** | 53.76 | 53.99 | 55.88 | 98.48 | 89.52 | 90.00 | 83.33 |
| 24-Hybrid | ✓ | ✓ | 59.52 | 54.15 | 54.41 | 55.94 | 98.48 | 90.48 | **91.67** | 83.33 |

Table 8: Meta-evaluation scores of reference-free (QE) MetricX variants on the WMT23 MQM evaluation set. "23" is the baseline, the last row of "24-Hybrid" corresponds to our primary submission, and "24" is our secondary submission. The hybrid model is the same as in Table 3, only evaluated without references provided as input. †Alternate references.

translation category, a QE model does not have access to the reference, so it makes perfect sense for it to score a candidate translation better than the reference translation if the reference is of low quality.

Switching over to Table 8, which shows the pairwise accuracy of the QE model scores, the trends are also in line with those of the reference-based models in Table 3. In contrast to the reference-based results, however, the hybrid model (row 6) does not outperform the standalone model (row 2), although most of the differences are within the expected variance. An astute reader might notice that the accuracy scores on the zh-en test set with the original references and the one with the alternate references do not match (despite the QE models not using the references), and that is because the latter has the original references included as an additional "human system".

Finally, we note that our QE models do not fall far behind their reference-based counterparts. In fact, both our primary and secondary QE submissions of MetricX-24 outperform our reference-based MetricX-23 submission from last year, according to the WMT23 MQM evaluation set.

# Evaluating WMT 2024 Metrics Shared Task Submissions on AfriMTE (the African Challenge Set)

**Jiayi Wang[1], David Ifeoluwa Adelani[2,3,4], Pontus Stenetorp[1]**

[1]University College London, UK, [2]Mila - Quebec AI Institute, Canada
[3]McGill University, Canada, [4]Canada CIFAR AI Chair

ucabj45@ucl.ac.uk, david.adelani@mila.quebec, p.stenetorp@cs.ucl.ac.uk

## Abstract

The AFRIMTE challenge set from WMT 2024 Metrics Shared Task aims to evaluate the capabilities of evaluation metrics for machine translation on low-resource African languages, which primarily assesses cross-lingual transfer learning and generalization of machine translation metrics across a wide range of under-resourced languages. In this paper, we analyze the submissions to WMT 2024 Metrics Shared Task. Our findings indicate that language-specific adaptation, cross-lingual transfer learning, and larger language model sizes contribute significantly to improved metric performance. Moreover, supervised models with relatively moderate sizes demonstrate robust performance, when augmented with specific language adaptation for low-resource African languages. Finally, submissions show promising results for language pairs including Darija-French, English-Egyptian Arabic, and English-Swahili. However, significant challenges persist for extremely low-resource languages such as English-Luo and English-Twi, highlighting areas for future research and improvement in machine translation metrics for African languages.

## 1 Introduction

Recent machine translation (MT) research has scaled dramatically, encompassing hundreds of languages, including many under-resourced ones (Fan et al., 2021a; NLLB-Team et al., 2022; Bapna et al., 2022; Kudugunta et al., 2023). However, accurately measuring MT quality in low-resource languages remains challenging. Conventional metrics like BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and chrF (Popović, 2015), which rely on n-gram matching, often fail to capture deeper semantic similarities (Zhang et al., 2020; Rei et al., 2020; Sai B et al., 2023).

Newer approaches include embedding-based metrics like BERTScore (Zhang et al., 2020) and learned metrics such as COMET (Rei et al., 2020), which have shown promise in more accurately evaluating translations across diverse languages. However, the application of these neural-based metrics to under-resourced languages continues to face significant challenges (Wang et al., 2024), highlighting ongoing areas of research in multilingual MT evaluation. These challenges include: (1) data scarcity impeding metric development, (2) complexity in annotation guidelines challenging non-expert evaluators, and (3) limited language model coverage restricting applicability, which underscore the need for continued innovation in MT evaluation methods, particularly for under-resourced African languages.

In response to these challenges, Wang et al. (2024) have introduced AFRIMTE, a human evaluation dataset focusing on MT adequacy and fluency for 13 typologically diverse African languages. This dataset addresses the data scarcity issue and employs simplified MQM evaluation guidelines tailored for non-expert translators, thus tackling two of the primary challenges in this field. Moreover, the authors establish benchmark systems for MT Evaluation and reference-free Quality Estimation (QE) by leveraging transfer learning techniques. These techniques draw from existing, well-resourced Direct Assessment (Graham et al., 2013) (DA) data and utilize an African-centric multilingual pre-trained language model, thereby addressing the challenge of limited language model coverage for African languages.

Building on this work, the WMT 2024 Metrics Shared task incorporates the translation adequacy test set from AFRIMTE as a challenge set. This inclusion aims to evaluate the capabilities of metric systems for machine translation on low-resource African languages, primarily assessing the cross-lingual transfer learning ability and generalization of these systems across a wide range of under-resourced African languages.

505

Our examination of task submissions has yielded several key findings in the development of machine translation metrics for African languages. We observed that language-specific adaptation, cross-lingual transfer learning, and increased language model sizes contribute to significant improvements in metric performance. Even supervised models of relatively modest scale can achieve robust results when augmented with language adaptation techniques. In addition, our analysis reveals promising outcomes for certain language pairs, such as Darija-French, English-Egyptian Arabic, and English-Swahili. However, persistent challenges remain evident in extremely low-resource languages like English-Luo and English-Twi. These disparities underscore critical areas requiring further investigation and highlight the need for targeted research in developing effective metrics across the diverse linguistic landscape of Africa.

## 2  AFRIMTE

AFRIMTE (Wang et al., 2024) focuses on the dev and devtest subsets of the FLORES-200 dataset (NLLB-Team et al., 2022). It covers 13 language pairs (LPs), primarily focusing on African languages with English, plus Darija-French and a control pair of English-French. In details, there are Darija-French (ary-fr), English-Egyptian Arabic (en-arz), English-French (en-fr), English-Hausa (en-hau), English-Igbo (en-ibo), English-Kikuyu (en-kik), English-Luo (en-luo), English-Somali (en-som), English-Swahili (en-swh), English-Twi (en-twi), English-isiXhosa (en-xho), English-Yoruba (en-yor), and Yoruba-English (yor-en). The annotations were also extended on domain-specific translations for English-Yoruba.

The MT outputs were generated using two open-source MT engines: NLLB-200 (600M) (NLLB-Team et al., 2022) and M2M-100 (418M) (Fan et al., 2021b). Most language pairs use NLLB-200, except for English-French and English-Swahili, which use M2M-100 due to their exceptionally high translation quality based on NLLB-200. The authors noted that while some language pairs like English-isiXhosa showed high overall quality, minor errors at the word level were still present.

AFRIMTE initially provides both fine-grained word-level error annotations and sentence-level Direct Assessment scoring for translation adequacy and fluency. For the WMT 2024 Metrics Shared Task, we utilize the adequacy test set from

| LP | Test # | LP | Test # |
|---|---|---|---|
| ary-fr | 187 | en-som | 226 |
| en-arz | 250 | en-swh | 157 |
| en-fr | 250 | en-twi | 247 |
| en-hau | 240 | en-xho | 243 |
| en-ibo | 120 | en-yor | 239 |
| en-kik | 202 | yor-en | 212 |
| en-luo | 242 | | |
| Total: 2815 annotations | | | |

Table 1: Counts of adequacy annotations for each language pair (LP) in the test set of AFRIMTE.

AFRIMTE as the African Challenge set to evaluate the sentence-level scoring performance of submitted metrics, focusing specifically on the FLORES-200 subsets within the dataset. Table 1 presents the counts of translation annotations in this challenge set. Due to the limited sizes of annotations for individual language pairs, we merge test data from all LPs into a single African-centric dataset to enhance evaluation significance for MT evaluation and reference-free quality estimation (QE) metrics. However, recognizing that different LPs may have varying score ranges, potentially favoring metrics that correlate with these distributions more than actual quality, we also report metric performance on each LP separately. This approach balances the need for statistical robustness with LP-specific insights.

## 3  Metrics

The WMT 2024 Metrics Shared Task received various metric submissions from both task organizers and participants. Our analysis will concentrate on the baseline metrics provided by the task organizers and the primary and contrastive metrics submitted by the participants.

### 3.1  Baselines

The baseline metrics for MT evaluation include BLEU (Papineni et al., 2002), chrF (Popović, 2015), spBLEU (Fan et al., 2021a), prism-Ref (Thompson and Post, 2020), YiSi-1 (Lo, 2019), COMET-22 (Rei et al., 2022a), BLUERT-20 (Sellam et al., 2020), and BertScore (Zhang et al., 2019). For reference-free quality estimation, the baseline metric is CometKiwi (Rei et al., 2022b). Additionally, we include AfriCOMET and AfriCOMET-QE for comparison, which are the African extensions of COMET-22 (Rei et al., 2022a) and CometKiwi (Rei et al., 2022b) pro-

posed by Wang et al. (2024). They employ transfer learning from well-resourced DA data and utilize an African-centric multilingual pre-trained encoder, AfroXLMR (Alabi et al., 2022), to build MT evaluation and QE models for African languages.

## 3.2 Submissions from Participants

The metrics submitted by participants for MT evaluation include XCOMET (Guerreiro et al., 2023), METRICX-24 and METRICX-24-HYBRID (Juraska et al., 2024)[1], chrF-S (Mukherjee and Shrivastava, 2024), METAMETRICS-MT (Anugraha et al., 2024), damonmonli, and monmonli[2]. For reference-free QE, the submitted metrics are XCOMET-QE (Guerreiro et al., 2023), METRICX-24-QE and METRICX-24-HYBRID-QE (Juraska et al., 2024)[3], QE model of METAMETRICS-MT (Anugraha et al., 2024), GEMBA-ESA (Kocmi and Federmann, 2023), and XLsimMQM (Mukherjee and Shrivastava, 2023). Details of all metrics can be found in Freitag et al. (2024).

## 3.3 AfriCOMET-1.1 and AfriCOMET-QE-1.1

In the ongoing efforts to enhance performance on African languages, we explore the use of a more advanced African pre-trained encoder. Specifically, we re-train AfriCOMET and AfriCOMET-QE using AfroXLMR-76 (Adelani et al., 2024) and conduct the training in single-task learning mode (Wang et al., 2024).

AfroXLMR-76 (Adelani et al., 2024) is an enhanced version of AfroXLMR (Alabi et al., 2022), which itself was a multilingual adaptation of the XLM-R-large model for 20 widely spoken African languages (each with at least 50MB of data). AfroXLMR-76 scales the language coverage up to 76 languages, including 61 languages with at least 10MB of data and an additional 15 languages with less than 10MB. To address the scarcity of monolingual data for some African languages, Adelani et al. (2024) proposed to generate synthetic parallel sentences by translating an English news commentary dataset (Kocmi et al., 2022) using the NLLB (600M) model.This expanded language coverage and increased training data volume have resulted in AfroXLMR-76 outperforming its predecessor, AfroXLMR, on the SIB-200 topic classification

| Metrics | Pearson | Spearman | Kendall |
|---|---|---|---|
| METRICX-24* | 0.5188 | 0.3949 | 0.2714 |
| AfriCOMET-1.1* | 0.5117 | 0.4129 | 0.2865 |
| AfriCOMET-1.0 | 0.4821 | 0.3857 | 0.2675 |
| METRICX-24-HYBRID | 0.4764 | 0.3844 | 0.2640 |
| METAMETRICS-MT | 0.3934 | 0.3429 | 0.2360 |
| COMET-22 | 0.3674 | 0.2835 | 0.1943 |
| YiSi-1 | 0.3058 | 0.2453 | 0.1666 |
| chrF-S | 0.3121 | 0.2332 | 0.1584 |
| chrF | 0.2833 | 0.2193 | 0.1492 |
| BERTScore | 0.2959 | 0.1834 | 0.1248 |
| BLEURT-20 | 0.2284 | 0.2225 | 0.1492 |
| XCOMET | 0.2224 | 0.2119 | 0.1451 |
| spBLEU | 0.2159 | 0.2052 | 0.1388 |
| monmonli | 0.2022 | 0.1713 | 0.1152 |
| damonmonli | 0.2007 | 0.1690 | 0.1138 |
| BLEU | 0.1863 | 0.1897 | 0.1282 |
| PrismRefMedium | 0.1149 | 0.1799 | 0.1202 |
| PrismRefSmall | 0.1058 | 0.1642 | 0.1099 |

Table 2: Segment-level correlation coefficients of **MT evaluation** metrics on the entire AFRIMTE. Metrics marked with * are ranked first based on the Perm-Input hypothesis test (Deutsch et al., 2021).

dataset for African languages (Adelani et al., 2024).

We refer to the original models using AfroXLMR as AfriCOMET-1.0[4] and AfriCOMET-QE-1.0[5], while the new versions leveraging AfroXLMR-76 are called AfriCOMET-1.1[6] and AfriCOMET-QE-1.1[7], respectively.

## 4 Analysis

This section presents a comprehensive analysis of the metrics outlined in Section 3. Our evaluation framework is structured around two primary components. First, we assess segment-level performance by examining the correlation between metric scores and human DA scores. This assessment involves analyzing correlation coefficients on the entire mixed African Challenge set and calculating weighted average correlation coefficients across various language pairs. Second, we conduct a language-specific analysis by computing average correlation coefficients for each individual language pair across all metric systems.

---

[1]METRICX-24 is the contrastive system to METRICX-24-HYBRID

[2]The monmonli is the contrastive system to damonmonli.

[3]METRICX-24-QE is the contrastive system to METRICX-24-HYBRID-QE

[4]https://huggingface.co/masakhane/africomet-stl

[5]https://huggingface.co/masakhane/africomet-qe-stl

[6]https://huggingface.co/masakhane/africomet-stl-1.1

[7]https://huggingface.co/masakhane/africomet-qe-stl-1.1

| Metrics | Pearson | Spearman | Kendall |
|---|---|---|---|
| METRICX-24* | 0.6269 | 0.4833 | 0.3455 |
| METRICX-24-HYBRID | 0.5972 | 0.4695 | 0.3351 |
| METAMETRICS-MT | 0.5295 | 0.4726 | 0.3368 |
| AfriCOMET-1.1 | 0.5399 | 0.4363 | 0.3097 |
| AfriCOMET-1.0 | 0.5260 | 0.4261 | 0.3027 |
| XCOMET | 0.4108 | 0.4045 | 0.2874 |
| COMET-22 | 0.4513 | 0.3432 | 0.2430 |
| YiSi-1 | 0.4233 | 0.3125 | 0.2182 |
| BLEURT-20 | 0.3604 | 0.3428 | 0.2396 |
| BERTScore | 0.3997 | 0.2933 | 0.2049 |
| chrF-S | 0.3763 | 0.3025 | 0.2106 |
| damonmonli | 0.3627 | 0.3013 | 0.2100 |
| chrF | 0.3593 | 0.2955 | 0.2053 |
| monmonli | 0.3215 | 0.2877 | 0.1991 |
| PrismRefMedium | 0.2389 | 0.2978 | 0.2053 |
| PrismRefSmall | 0.2250 | 0.2868 | 0.1984 |
| spBLEU | 0.2585 | 0.2515 | 0.1733 |
| BLEU | 0.2394 | 0.2457 | 0.1691 |

Table 3: Segment-level weighted average correlation coefficients of **MT evaluation** metrics, averaged across language pairs on AFRIMTE, with weights based on the size of each language pair group. The metric marked with * ranks first based on the average of Pearson, Spearman, and Kendall correlation coefficients.

## 4.1 Segment-level Averaged Correlation

For both MT evaluation and reference-free QE tasks, we assess the metric performance using three widely adopted correlation coefficients: Pearson, Spearman-rank, and Kendall-rank. These coefficients measure the correlation between metric scores and human DA scores, each capturing different aspects of the relationship (Deutsch et al., 2023). To validate the statistical significance of our results, we additionaly employ the Perm-Input hypothesis test (Deutsch et al., 2021), which is conducted with 200 re-sampling runs and a significance level of $p = 0.05$, producing rankings of the various automatic metrics based on their performance.

### 4.1.1 MT Evaluation Metric

We present the segment-level correlation coefficients of MT evaluation metrics on the entire AFRIMTE test set in Table 2 and the weighted average correlation coefficients across various language pairs in Table 3. Detailed Pearson, Spearman-rank, and Kendall-rank correlations of baseline metrics and primary submissions for each language pair are shown in Figures 3, 4, and 5 of Appendix A.

For MT evaluation, Table 2 provides valuable insights into evaluation metrics' performance on the African Challenge Set. Generally, Pearson correlations are generally higher

than Spearman and Kendall, with rankings remaining largely consistent across correlation types. The top-performing metrics—METRICX-24, AfriCOMET-1.1, AfriCOMET-1.0, and METRICX-24-HYBRID—are all based on pretrained multilingual large language models (LLMs) and utilize supervised learning. These metrics consistently outperform other types across all correlation coefficients. METRICX-24 and AfriCOMET-1.1 emerge as the best performers, statistically indistinguishable from each other. The improved performance of AfriCOMET-1.1 over its predecessor suggests ongoing enhancements in these LLM-based metrics. It is evident that African-centric LLM-based metrics (AfriCOMET variants) perform exceptionally well, highlighting the importance of language-specific fine-tuning for low-resource African languages.

Moreover, the weighted average correlation results presented in Table 3 offer additional valuable insights. METRICX-24 still emerges as the top-performing metric, achieving the highest correlation with human judgments across all three correlation coefficients (Pearson: 0.6269, Spearman: 0.4833, Kendall: 0.3455). Its hybrid variant, METRICX-24-HYBRID, follows closely, suggesting the robustness of this metric family. METAMETRICS-MT shows strong performance, ranking third overall with high correlation coefficients. As an ensemble method, it selectively combines complementary metrics, proves effective for African languages despite these metrics being trained on general WMT data. In addition, AfriCOMET-1.1 and its predecessor AfriCOMET-1.0 show robust performance indicating their effectiveness for African language pairs.

Traditional metrics like BLEU and its variant spBLEU demonstrate relatively weak correlations, reinforcing the need for more advanced metrics in evaluating MT quality for African languages. Interestingly, some widely-used metrics such as BERTScore and BLEURT-20 show moderate performance, outperforming traditional metrics but falling behind the top-performing ones. The consistent ranking across different correlation coefficients suggests a reliable performance hierarchy among these metrics. However, the overall moderate correlation values (mostly below 0.5 for Spearman and Kendall) highlight the ongoing challenges in accurately evaluating MT quality for African languages.

Figure 1: Average correlations across MT evaluation metrics for each language pair.

| Metrics | Pearson | Spearman | Kendall |
|---|---|---|---|
| METRICX-24-QE* | 0.4857 | 0.3810 | 0.2616 |
| AfriCOMET-QE-1.1* | 0.4760 | 0.3961 | 0.2747 |
| METRICX-24-HYBRID-QE | 0.4337 | 0.3594 | 0.2464 |
| GEMBA-ESA | 0.4033 | 0.3300 | 0.2427 |
| METAMETRICS-MT | 0.3781 | 0.3004 | 0.2050 |
| AfriCOMET-QE-1.0 | 0.3496 | 0.2524 | 0.1729 |
| CometKiwi-XXL | 0.2149 | 0.1814 | 0.1254 |
| XCOMET-QE | 0.1717 | 0.1528 | 0.1042 |
| CometKiwi | 0.1685 | 0.1259 | 0.0838 |
| XLsimMQM | 0.0886 | 0.0925 | 0.0619 |

Table 4: Segment-level correlation coefficients of **QE** metrics on AFRIMTE. Metrics marked with * are ranked first based on the Perm-Input hypothesis test ([Deutsch et al., 2021](#)).

| Metrics | Pearson | Spearman | Kendall |
|---|---|---|---|
| METRICX-24-QE* | 0.5790 | 0.4383 | 0.3117 |
| METRICX-24-HYBRID-QE | 0.5530 | 0.4289 | 0.3048 |
| AfriCOMET-QE-1.1 | 0.4905 | 0.4117 | 0.2900 |
| GEMBA-ESA | 0.4624 | 0.3793 | 0.2900 |
| METAMETRICS-MT | 0.5010 | 0.3610 | 0.2528 |
| AfriCOMET-QE-1.0 | 0.4774 | 0.3743 | 0.2628 |
| CometKiwi-XXL | 0.3709 | 0.3428 | 0.2417 |
| XCOMET-QE | 0.3087 | 0.3290 | 0.2317 |
| CometKiwi | 0.3301 | 0.2914 | 0.2046 |
| XLsimMQM | 0.1548 | 0.1817 | 0.1256 |

Table 5: Segment-level weighted average correlation coefficients of **QE** metrics, averaged across language pairs on AFRIMTE, with weights based on the size of each language pair group. The metric marked with * ranks first based on the average of Pearson, Spearman, and Kendall correlation coefficients.

### 4.1.2 Quality Estimation as a Metric

QE presents a more challenging and purely cross-lingual task, making its investigation essential. Ta-

bles 4 and 5 presents the segment-level correlation coefficients of QE metrics on the entire AFRIMTE and weighted average correlations across language pairs. Detailed Pearson, Spearman-rank, and Kendall-rank correlations of baseline metrics and primary submissions for each language pair are shown in Figures 6, 7, and 8 of Appendix A.

Comparing results in Tables 2 and 4, and results in Tables 3 and 5, we have observed significant performance gaps between MT evaluation models and their QE counterparts. This is evident when comparing specific versions, such as the differences between METRICX-24 and METRICX-24-QE, XCOMET and XCOMET-QE, as well as AfriCOMET-1.1 and AfriCOMET-QE-1.1. These disparities underscore the increased complexity of the QE task, which requires assessing translation quality without access to reference translations.

Tables 4 and 5 reveal the superior performance of LLM-based supervised-learning metrics in the QE task. Specifically, METRICX-24-QE and AfriCOMET-QE-1.1 emerge as the top-performing metrics on the entire AFRIMTE test set (Table 4). These metrics demonstrate statistically indistinguishable performance, as confirmed by the Perm-Input hypothesis test. Furthermore, in the weighted average correlation across different language pairs (Table 5), METRICX-24-QE consistently outperforms other approaches. This trend in QE metrics mirrors the pattern observed in MT evaluation metrics, underscoring the effectiveness of LLM-based supervised-learning approaches in both contexts for African languages. Additionally,

METAMETRICS-MT, as a meta-metric, continues to show strong performance, further validating the effectiveness of ensemble methods in addressing the complexities of African language evaluation. Another LLM-based metric, GEMBA-ESA, which employs a two-step approach: first collecting MQM error spans, and then assigning the final score also demonstrates robust performance, further highlighting the potential of LLM-based techniques in QE tasks for African languages. However, supervised QE metrics such as CometKiwi, CometKiwi-XXL, and XCOMET-QE show relatively poor performance, suggesting they might not be well-suited for African languages without specific language adaptation.

### 4.1.3 Language Adaptation, Cross-lingual Transfer, and Model Size as Key Factors in Metric Performance

Our analysis on the baseline and task submissions reveals that language-specific tuning, cross-lingual transfer learning, and model size are crucial factors in MT evaluation and Quality Estimation.

The top-performing systems demonstrate these principles in various ways. METRICX-24 systems, based on mT5-XXL (Xue et al., 2020), cover a wide range of languages, including several African languages such as Hausa, Igbo, Somali, Swahili, Xhosa, Yoruba, and Zulu. In contrast, AfriCOMET models use African-enhanced masked language models (AfroXLMR and AfroXLMR-76) with well-resourced DA training data, showcasing the benefits of language-specific adaptation. Both METRICX-24 and AfriCOMET variants employ supervised training and cross-lingual transfer learning, proving effective for low-resource language scenarios. The impact of model size is evident, with AfriCOMET variants (560 million parameters) and METRICX-24 (13 billion parameters) both achieving strong results. While METRICX-24's larger size contributes to its superior performance, AfriCOMET's performance demonstrates that well-adapted smaller models can also yield robust results.

Moreover, the excellent performance of METAMETRICS-MT underscores the potential of ensembling robust metrics to create effective meta-metrics. The promising results of GEMBA-ESA further highlight the effectiveness of LLM-based prompting techniques in this domain. These findings collectively emphasize the potentials of model ensemble and innovative

LLM prompting strategies in developing effective MT evaluation and QE metrics, particularly for low-resource languages.

### 4.2 Language-Specific Performance: Average Correlations across Metrics

To investigate model performance on specific language pairs, we calculate the average correlation coefficients for each individual language pair across all metric systems, providing insights into how well metrics perform for specific language pairs. Results are shown in Figure 1 and 2.

#### 4.2.1 Performance on MT Evaluation

Figure 1 depicting the average correlation coefficients across metric submissions for MT evaluation reveals significant variations in metric performance on different language pairs. Consistently across all pairs, Pearson correlation shows the highest values, followed by Spearman and then Kendall, suggesting stronger linear relationships between human and metric scores compared to monotonic or ordinal relationships. English-Swahili (en-swh) and Darija-French (ary-fr) demonstrate the highest correlations across all three metrics, likely due to their status as more resource-rich or commonly studied pairs. In contrast, English-Luo (en-luo), English-Twi (en-twi), and English-isiXhosa (eng-xho) exhibit the lowest correlations, indicating particular challenges for MT evaluation in these language pairs.

#### 4.2.2 Performance on QE

A consistent pattern emerges in the QE task (Figure 2) where Pearson correlations generally show the highest values. Language pair performance is notably similar across both figures, with resource-rich pairs like English-Swahili (en-swh) consistently demonstrating higher correlations, while extremely low-resource pairs such as English-Luo (en-luo) and English-Twi (en-twi) show persistently lower correlations. Interestingly, some language pairs show improved relative performance in QE compared to MT Evaluation. For example, English-Egyptian Arabic (en-arz) and English-Hausa (en-hau) demonstrate better results in QE, possibly indicating their suitability for reference-free evaluation methods.

#### 4.2.3 Some Special Cases

Contrary to expectations, English-French (en-fr) does not emerge as the top-performing language

**Average Correlation Coefficients across Systems for QE**

Figure 2: Average correlations across QE metrics for each language pair.

pair in either the MT evaluation or the QE task. This surprising result might be attributed to two factors. First, as illustrated in Table 7 of Wang et al. (2024), there is a scarcity of supervised DA training datasets for English-French. Second, the performance may be affected by the "curse of multilinguality" (Pfeiffer et al., 2022), a phenomenon where model performance on high-resource languages can degrade when the pre-trained model is fine-tuned and enhanced with data from multiple low-resource languages, in this case, African languages.

Another noteworthy case is English-isiXhosa (en-xho). As previously observed, English-isiXhosa translations demonstrated high overall sentence-level quality (median DA: 100 according to Wang et al. (2024)) , with only minor errors at the word level. This characteristic makes it particularly challenging to differentiate and rank translation quality. Consequently, the relatively lower performance of Spearman and Kendall for English-isiXhosa is expected.

## 5 Conclusion

In conclusion, our analysis on submissions to the AFRIMTE challenge set of WMT 2024 Metrics Shared Task for African languages reveals that LLM-based supervised-learning metrics, especially those with African-centric tuning, consistently outperform traditional and other neural-based approaches in both MT evaluation and Quality Estimation tasks. Language-specific adaptation, cross-lingual transfer learning, and larger model sizes contribute significantly to improved metric performance. However, challenges persist for extremely low-resource languages such as Luo and Twi. Our analysis also highlights unexpected performance patterns in certain language pairs, including English-French and English-isiXhosa, demonstrating the complexities of evaluating machine translation across diverse African languages.

## References

David Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba Alabi, Yanke Mao, Haonan Gao, and En-Shiun Lee. 2024. SIB-200: A simple, inclusive, and big evaluation dataset for topic classification in 200+ languages and dialects. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian's, Malta. Association for Computational Linguistics.

Jesujoba O. Alabi, David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. 2022. Adapting pre-trained language models to African languages via multilingual adaptive fine-tuning. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4336–4349, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

David Anugraha, Garry Kuwanto, Lucky Susanto, Derry Tanti Wijaya, and Genta Indra Winata. 2024. Metametrics-MT: Tuning machine translation metametrics via human preference calibration. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR:

An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Ankur Bapna, Isaac Caswell, Julia Kreutzer, Orhan Firat, Daan van Esch, Aditya Siddhant, Mengmeng Niu, Pallavi N. Baljekar, Xavier García, Wolfgang Macherey, Theresa Breiner, Vera Axelrod, Jason Riesa, Yuan Cao, Mia Xu Chen, Klaus Macherey, Maxim Krikun, Pidong Wang, Alexander Gutkin, Apurva Shah, Yanping Huang, Z. Chen, Yonghui Wu, and Macduff Hughes. 2022. Building machine translation systems for the next thousand languages. *ArXiv*, abs/2205.03983.

Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. A statistical analysis of summarization evaluation metrics using resampling methods. *Transactions of the Association for Computational Linguistics*, 9:1132–1146.

Daniel Deutsch, George Foster, and Markus Freitag. 2023. Ties matter: Meta-evaluating modern metrics with pairwise accuracy and tie calibration. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12914–12929.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2021a. Beyond english-centric multilingual machine translation. *J. Mach. Learn. Res.*, 22(1).

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021b. Beyond english-centric multilingual machine translation. *The Journal of Machine Learning Research*, 22(1):4839–4886.

Markus Freitag, Nitika Mathur, Daniel Deutsch, Chikiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are llms breaking mt metrics? results of the wmt24 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria. Association for Computational Linguistics.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.

Juraj Juraska, Daniel Deutsch, Mara Finkelstein, and Markus Freitag. 2024. Metricx-24: The google submission to the WMT 2024 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. Gembamqm: Detecting translation quality error spans with gpt-4. *arXiv preprint arXiv:2310.13988*.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier García, Christopher A. Choquette-Choo, Katherine Lee, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2023. Madlad-400: A multilingual and document-level large audited dataset. *ArXiv*, abs/2309.04662.

Chi-kiu Lo. 2019. YiSi - a unified semantic MT quality evaluation and estimation metric for languages with different levels of available resources. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 507–513, Florence, Italy. Association for Computational Linguistics.

Ananya Mukherjee and Manish Shrivastava. 2023. MEE4 and XLsim: IIIT HYD's Submissions for WMT23 Metrics Shared Task. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.

Ananya Mukherjee and Manish Shrivastava. 2024. chrf-s: Semantics is all you need. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

NLLB-Team, Marta Ruiz Costa-jussà, James Cross, Onur cCelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Alison Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon L.

Spruit, C. Tran, Pierre Yves Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzm'an, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *ArXiv*, abs/2207.04672.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Jonas Pfeiffer, Naman Goyal, Xi Victoria Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. 2022. Lifting the curse of multilinguality by pre-training modular transformers. *arXiv preprint arXiv:2205.06266*.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022a. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ananya Sai B, Tanay Dixit, Vignesh Nagarajan, Anoop Kunchukuttan, Pratyush Kumar, Mitesh M. Khapra, and Raj Dabre. 2023. IndicMT eval: A dataset to meta-evaluate machine translation metrics for Indian languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14210–14228, Toronto, Canada. Association for Computational Linguistics.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Brian Thompson and Matt Post. 2020. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 90–121, Online. Association for Computational Linguistics.

Jiayi Wang, David Adelani, Sweta Agrawal, Marek Masiak, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Xuanli He, Sofia Bourhim, Andiswa Bukula, Muhidin Mohamed, Temitayo Olatoye, Tosin Adewumi, Hamam Mokayed, Christine Mwase, Wangui Kimotho, Foutse Yuehgoh, Anuoluwapo Aremu, Jessica Ojo, Shamsuddeen Muhammad, Salomey Osei, Abdul-Hakeem Omotayo, Chiamaka Chukwuneke, Perez Ogayo, Oumaima Hourrane, Salma El Anigri, Lolwethu Ndolela, Thabiso Mangwana, Shafie Mohamed, Hassan Ayinde, Oluwabusayo Awoyomi, Lama Alkhaled, Sana Al-azzawi, Naome Etori, Millicent Ochieng, Clemencia Siro, Njoroge Kiragu, Eric Muchiri, Wangari Kimotho, Toadoum Sari Sakayo, Lyse Naomi Wamba, Daud Abolade, Simbiat Ajao, Iyanuoluwa Shode, Ricky Macharm, Ruqayya Iro, Saheed Abdullahi, Stephen Moore, Bernard Opoku, Zainab Akinjobi, Abeeb Afolabi, Nnaemeka Obiefuna, Onyekachi Ogbu, Sam Ochieng', Verrah Otiende, Chinedu Mbonu, Yao Lu, and Pontus Stenetorp. 2024. AfriMTE and AfriCOMET: Enhancing COMET to embrace under-resourced African languages. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5997–6023, Mexico City, Mexico. Association for Computational Linguistics.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2020. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

# A  Appendix

Detailed Pearson, Spearman-rank, and Kendall correlation coefficients of MT evaluation and QE metrics for each language pair are shown in Figures 3, 4, 5, 6, 7, and 8 accordingly.

Figure 3: Pearson Correlations of MT Evaluation Metrics for each language pair.



Figure 4: Spearman-rank Correlations of MT Evaluation Metrics for each language pair.

Figure 5: Kendall-rank Correlations of MT Evaluation Metrics for each language pair.



Figure 6: Pearson Correlations of QE Metrics for each language pair.

Figure 7: Spearman-rank Correlations of QE Metrics for each language pair.



Figure 8: Kendall-rank Correlations of QE Metrics for each language pair.

# Machine Translation Metrics are better in evaluating Linguistic Errors on LLMs than on Encoder-Decoder Systems

**Eleftherios Avramidis, Shushen Manakhimova, Vivien Macketanz and Sebastian Möller**
German Research Center for Artificial Intelligence (DFKI),
Berlin, Germany
`firstname.lastname@dfki.de`

## Abstract

This year's MT metrics challenge set submission by DFKI expands previous years' linguistically motivated challenge set. It includes 137,000 items extracted from 100 MT systems for the two language directions (en→de, en→ru), covering more than 100 linguistically motivated phenomena organized in 14 linguistic categories. The metrics with the statistically significant best performance with regard to our linguistically motivated analysis are METRICX-24-HYBRID and METRICX-24 for en→de and METRICX-24 for en→ru, whereas METAMETRICS and XCOMET are in the next ranking positions in both language pairs. Metrics are more accurate in detecting linguistic errors among LLM translations than in translations based on the encoder-decoder NMT architecture. Some of the most difficult phenomena for the metrics to score are the transitive past progressive, the multiple connectors, and the ditransitive simple future I for en→de and the pseudogapping, the contact clause and the cleft sentences for en→ru. Despite its overall low performance, the LLM-based metric GEMBA performs best in scoring German negation errors.

## 1 Introduction

For almost two decades, the development and evaluation of machine translation (MT) have relied on automatic metrics. MT metrics aim to digest and automate various aspects of human judgment of MT output into numerical scores. Over the years, these metrics have undergone several technological changes (from measuring overlap to grammatical features and neural models). Still, at the same time, they have had to follow the technological evolution of MT systems, moving from phrase-based statistical systems to NMT encoder-decoder models and, more recently, to large language models (LLMs). As we witness the first efforts to use and evaluate LLMs in the task of MT, it is of great interest to see to what extent pre-existing MT methodologies

can adapt to the needs of the new technologies. An obvious question is to what extent MT metrics developed and tested for NMT can be applied to evaluating LLMs.

This year's Metrics Task (WMT24; Freitag et al., 2024) provides a very good opportunity to evaluate the metrics under these particular circumstances, as the evaluated MT outputs have for the first time been produced by numerous LLMs (Kocmi et al., 2024). Meanwhile, the ability of LLMs to act as judges for translations is being explored through the participation of an LLM-based metric.

Given this perspective, this paper extends previous work on linguistically motivated challenge sets for MT metrics to investigate whether LLMs can influence MT evaluation. As part of this year's submission to the challenge set subtask of the WMT24 Metrics Task, we repeat the methodology of previous years to evaluate the metrics on a controlled test set that can rank them with regard to their ability to detect linguistic errors by providing fine-grained statistics for each linguistic phenomenon. We then analyze whether the metrics perform differently on MT output from LLMs as opposed to output from encoder-decoder systems. In addition, we see in which linguistic aspects the LLM-based metric performs better or worse than the specialized metrics.

The rest of the paper is structured as following: Section 2 describes briefly the generation of the challenge set. Section 3 presents and discusses the results, whereas the conclusion is given in section 4

## 2 Method

This year's linguistically-motivated challenge set is an extension of the challenge sets that were submitted the previous years (Avramidis and Macketanz, 2022; Avramidis et al., 2023).

The source sentences $s$ originate from an MT evaluation test suite (Macketanz et al., 2022a). Each sentence has been carefully constructed to test one particular phenomenon. Every phenomenon is

tested by more sentences (with a minimum of 20 sentences), whereas the phenomena are aggregated in a few categories. At the moment, there are more than 100 phenomena and 14 categories.

As part of the WMT shared tasks of the previous years, these source sentences have been given to a large amount of MT systems, and their output has been evaluated by combining regular expressions and annotations by linguists, labeling every output as correct ($t \in T$) or incorrect ($\hat{t} \in T'$).

In order to use this test set to evaluate the MT metrics, we create examples in the form of $(s, \hat{t}, t, r) \in S$, where each example contains one source sentence $s$, one incorrect translation hypothesis $\hat{t}$, one correct translation hypothesis $t$ and one reference translation $r$. The correct translation hypotheses $t$ and the reference translations $r$ are sampled with permutations from the same set of correct translations $T$. Then, we decompose the set of examples $S$ into a blind test set $S'$, where each example includes either an incorrect translation $(s, \hat{t}, r)$ or a correct translation $(s, t, r)$ along with the source and the reference. The separated contrastive examples are shuffled, and we set aside a file that contains the golden truth, indicating which samples are correct or incorrect.

As part of the Metrics Task, every shuffled translation $t$ and $hatt$ is scored by every $M$, given the reference $r$ in the given blind test set $S'$, without knowing if it is correct or incorrect. A contrastive pair scoring is considered correct if the metric delivers a score for the incorrect translation hypothesis, which is lower than the one of the correct translation hypothesis $M(s, \hat{t}, r) < M(s, t, r)$. Finally, for every phenomenon and category and for every metric, the respective accuracy is calculated by dividing the number of correctly scored contrastive pairs by the total amount of examples.

$$\text{acc}_M = \frac{|M(s, \hat{t}, r) < M(s, t, r)|}{|(s, \hat{t}, t, r)|}$$

$$(s, \hat{t}, r) \cup (s, t, r) \in S' \quad (s, \hat{t}, t, r) \in S$$

Lastly, we provide three types of score averaging:

i) **Micro-average:** This approach treats all items equally, aggregating all test items to compute the average percentages.

ii) **Category macro-average:** Here, all categories are treated equally, with the percent-ages being computed independently for each category and then averaged.

iii) **Phenomenon macro-average:** This average treats all phenomena equally, with the percent-ages being computed independently for each phenomenon and then averaged.

The current version of the challenge set contains MT outputs from the WMT Shared Tasks of the years 2019-2024 (Avramidis et al., 2019, 2020; Macketanz et al., 2021, 2022b; Manakhimova et al., 2023, 2024). The English to German version contains 39,463 contrastive pairs, while the English to Russian version contains 30,108 pairs.

## 3 Results

### 3.1 English-German

The comparison of the metrics based on the accuracies per category for English-German can be seen in table 2, whereas the detailed phenomena in table 4. One can see that the metrics which have the highest accuracy with statistical significance are METRICX24-HYBRID and METRICX24 (Juraska et al., 2024), with more than 80.7 % macro-average. Both metrics are very good at multi-word expressions (mostly verbal MWEs). The former is the best of all metrics at coordination/ellipsis and non-verbal agreement (genitive and personal pronoun coreference). In contrast, the latter performs best at verb valency (resultative and passive voice). The metrics "METAMETRICS" (Anugraha et al., 2024) and XCOMET (Guerreiro et al., 2023) follow in the ranking, with more than 80% macro-averaged accuracy.

The LLM-based metric GEMBA (Kocmi and Fe-dermann, 2023) performs relatively low, with an average accuracy of 69.7%, even below the base-line non-tuned metric CHRF (Popović, 2015). It is nevertheless remarkable that this metric has the best score on negation, among all metrics (97.4%, 4.5% higher than the best system). The fact that most of the metrics will miss 10% of the nega-tions is rather noteworthy, given the implications of such a mistake on the meaning of the sentence. It is also remarkable that a reference-less metric, METRICX24-HYBRID-QE, achieves the highest accuracy on long-distance dependencies and inter-rogatives, mainly on the phenomenon of negative inversion.

Some of the most difficult phenomena for the

| | METRICX24 | METRICX24-HYB | METAMETRICS | XCOMET |
|---|---|---|---|---|
| encdec vs. encdec | 73.2 | 72.3 | 70.8 | 69.7 |
| LLM vs. encdec | 77.3 | 76.9 | 79.9 | 77.6 |
| LLM vs. LLM | 79.9 | 78.1 | 80.0 | 79.1 |

Table 1: Accuracy of the metrics when they evaluate contrastive pairs containing (a) MT output only by encoder/decoder systems, (b) one encoder/decoder output and one LLM output, (c) only LLM output

metrics to score are transitive past progressive, multiple connectors, and ditransitive simple future I.

### 3.2 English-Russian

The comparison of the metrics based on the accuracies per category for English-Russian can be seen in table 3, whereas the detailed phenomena in table 5. MetricX-24 is the clear winner in this language direction, achieving a macro-averaged accuracy of 82.5% MetricX-24 excels in ambiguity, false friends, non-verbal agreement (coreference & genitive), verb semantics, and verb valency. The ranking of the metrics is similar to the one for English-German, with METAMETRICS, METRICX24-HYBRID and XCOMET having the next position, with more than 79.6% accuracy in macro-average.

If one focuses again on the phenomenon of negation, they would notice that in English-Russian, the highest accuracy is achieved by the baseline metric CHRF, whereas most metrics perform here very low (61% on average) Some of the most difficult phenomena for this language direction are the pseudogapping, the contract clause, and the cleft sentences for en→ru.

### 3.3 Comparing performance of metrics over LLM vs. encoder-decoder systems

Table 1 presents the accuracies of the 4 best performing metrics on three subsets of the challenge sets. Here every subset contains contrastive pairs which consist of

(a) two MT outputs, both by encoder/decoder NMT systems
(b) one encoder/decoder and one LLM output
(c) two LLM outputs

One can see that all four metrics exhibit higher accuracy when scoring contrastive translations originating from LLMs. This indicates that despite the fact that LLM translations achieve very good performance (Kocmi et al., 2024), their fewer errors are easier to be distinguished by the automatic metrics. Whether there is a systematic reason for this phenomenon remains to be investigated.

### 4 Conclusion

We presented the MT metrics challenge set of DFKI for two language directions (en-de, en-ru). This year, we have expanded the set to include outputs from encoder-decoder NMT systems and LLMs. The number of test items (total of 137,000) allows for producing fine-grained scores for every linguistic phenomenon and statistically significant comparisons among the MT metrics. We also identified the best-performing metric, METRICX-24, for both language directions.

### Acknowledgements

### References

David Anugraha, Garry Kuwanto, Lucky Susanto, Derry Tanti Wijaya, and Genta Indra Winata. 2024. Metametrics-MT: Tuning machine translation metametrics via human preference calibration. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Eleftherios Avramidis and Vivien Macketanz. 2022. Linguistically motivated evaluation of machine translation metrics based on a challenge set. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 514–529, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, Aljoscha Burchardt, and Sebastian Möller. 2020. Fine-grained linguistic evaluation for state-of-the-art machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 346–356, Online. Association for Computational Linguistics.

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, and Hans Uszkoreit. 2019. Linguistic evaluation of German-English machine translation using a test suite. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 445–454, Florence, Italy. Association for Computational Linguistics.

Eleftherios Avramidis, Shushen Manakhimova, Vivien Macketanz, and Sebastian Möller. 2023. Challenging the state-of-the-art machine translation metrics from a linguistic perspective. In *Proceedings of the Eighth Conference on Machine Translation*, pages 713–729, Singapore. Association for Computational Linguistics.

Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are llms breaking mt metrics? results of the wmt24 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection.

Juraj Juraska, Daniel Deutsch, Mara Finkelstein, and Markus Freitag. 2024. Metricx-24: The google submission to the WMT 2024 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Findings of the WMT24 general machine translation shared task: the LLM era is here but mt is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. GEMBA-MQM: Detecting translation quality error spans with GPT-4. In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775, Singapore. Association for Computational Linguistics.

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, and Hans Uszkoreit. 2022a. A linguistically motivated test suite to semi-automatically evaluate German–English machine translation output. In *Proceedings*

*of the Thirteenth Language Resources and Evaluation Conference*, pages 936–947, Marseille, France. European Language Resources Association.

Vivien Macketanz, Eleftherios Avramidis, Shushen Manakhimova, and Sebastian Möller. 2021. Linguistic evaluation for the 2021 state-of-the-art machine translation systems for German to English and English to German. In *Proceedings of the Sixth Conference on Machine Translation*, pages 1059–1073, Online. Association for Computational Linguistics.

Vivien Macketanz, Shushen Manakhimova, Eleftherios Avramidis, Ekaterina Lapshinova-koltunski, Sergei Bagdasarov, and Sebastian Möller. 2022b. Linguistically motivated evaluation of the 2022 state-of-the-art machine translation systems for three language directions. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 432–449, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Shushen Manakhimova, Eleftherios Avramidis, Vivien Macketanz, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2023. Linguistically motivated evaluation of the 2023 state-of-the-art machine translation: Can ChatGPT outperform NMT? In *Proceedings of the Eighth Conference on Machine Translation*, pages 224–245, Singapore. Association for Computational Linguistics.

Shushen Manakhimova, Vivien Macketanz, Eleftherios Avramidis, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2024. Investigating the linguistic performance of large language models in machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.

# A  Accuracies per category

Table 2: Accuracy of the metrics(%) with regards to the linguistically-motivated categories for English-German

| ling. category | # | Metrix-24-Hybrid | Metrix-24 | metricx | XCOMET | BLEURT-20 | COMET-22 | CometKiwi-XXL | Metrix-24-QE | Metrix-24-Hybrid-QE | XCOMET-QE | YiSi-1 | sentinel-cand-mqm | CometKiwi | MEE4 | chrfS | BERTScore | chrF | gemba | spBLEU | damonmonti | monmonti | BLEU | XLsimMqm | XLsimDA | PrismRefSmall | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | 4614.0 | 85.1 | 85.9 | 89.9 | 80.0 | 89.7 | 89.5 | 60.8 | 74.6 | 70.6 | 61.9 | 88.6 | 77.7 | 48.2 | 82.1 | 83.8 | 78.2 | 85.2 | 70.0 | 80.0 | 83.8 | 83.0 | 64.1 | 55.2 | 55.2 | 68.1 | 60.6 | 75.1 |
| Coordination & ellipsis | 4373.0 | 81.3 | 74.2 | 74.4 | 77.4 | 76.5 | 76.7 | 80.2 | 78.2 | 78.8 | 74.4 | 69.2 | 76.7 | 71.1 | 63.8 | 62.9 | 67.3 | 62.2 | 66.5 | 61.8 | 61.0 | 62.9 | 60.6 | 49.5 | 49.5 | 51.1 | 49.1 | 67.6 |
| False friends | 1389.0 | 79.9 | 78.2 | 78.3 | 73.9 | 72.7 | 85.9 | 85.2 | 69.8 | 74.3 | 74.4 | 73.1 | 77.0 | 72.1 | 80.4 | 77.1 | 74.4 | 74.9 | 81.9 | 65.9 | 48.6 | 38.2 | 64.1 | 78.3 | 78.3 | 74.3 | 58.8 | 72.4 |
| Function word | 1900.0 | 78.1 | 80.6 | 82.2 | 86.0 | 81.9 | 87.3 | 83.0 | 81.7 | 78.6 | 82.2 | 72.9 | 85.8 | 86.7 | 76.9 | 77.2 | 86.9 | 74.4 | 70.9 | 74.2 | 64.9 | 60.1 | 78.6 | 55.7 | 55.7 | 52.2 | 53.8 | 74.9 |
| LDD & interrogatives | 1002.0 | 83.4 | 80.1 | 80.8 | 80.6 | 80.3 | 74.5 | 78.7 | 81.8 | 84.7 | 81.7 | 59.1 | 68.6 | 78.4 | 64.2 | 59.3 | 64.1 | 57.5 | 58.6 | 61.5 | 66.7 | 64.5 | 60.5 | 62.0 | 62.0 | 49.6 | 47.7 | 68.9 |
| MWE | 5816.0 | 87.0 | 87.3 | 85.9 | 86.2 | 84.1 | 82.9 | 80.0 | 82.5 | 81.2 | 80.5 | 80.3 | 84.0 | 76.4 | 75.1 | 75.9 | 76.1 | 73.3 | 82.0 | 70.7 | 77.3 | 76.6 | 71.5 | 67.0 | 67.0 | 59.4 | 55.6 | 77.1 |
| Named entity & terminology | 22891.0 | 71.5 | 74.2 | 74.2 | 68.8 | 71.7 | 73.6 | 58.0 | 55.3 | 60.9 | 56.9 | 74.7 | 52.1 | 50.2 | 72.2 | 70.5 | 67.1 | 68.9 | 48.1 | 70.0 | 75.4 | 73.1 | 62.0 | 48.5 | 48.5 | 49.8 | 50.1 | 63.3 |
| Negation | 506.0 | 92.9 | 89.5 | 88.5 | 91.1 | 92.7 | 92.9 | 93.3 | 93.9 | 91.3 | 90.9 | 87.9 | 74.5 | 95.3 | 90.7 | 82.8 | 86.0 | 76.7 | 97.4 | 73.9 | 86.6 | 88.3 | 73.7 | 58.3 | 58.3 | 58.5 | 58.1 | 83.2 |
| Non-verbal agreement | 15497.0 | 83.6 | 80.6 | 77.4 | 80.9 | 78.2 | 73.3 | 80.2 | 82.3 | 82.4 | 79.2 | 65.7 | 76.2 | 72.9 | 65.6 | 66.1 | 63.7 | 65.9 | 72.7 | 64.3 | 59.6 | 59.5 | 62.2 | 57.8 | 57.8 | 51.0 | 49.0 | 69.5 |
| Punctuation | 2435.0 | 62.2 | 64.4 | 64.9 | 63.2 | 71.9 | 72.4 | 70.1 | 70.4 | 65.9 | 64.9 | 71.6 | 80.1 | 71.3 | 69.9 | 72.1 | 66.0 | 68.5 | 44.3 | 67.3 | 66.9 | 50.7 | 50.3 | 50.3 | 50.3 | 50.6 | 50.8 | 64.2 |
| Subordination | 4698.0 | 89.1 | 87.5 | 86.3 | 89.3 | 84.1 | 83.9 | 89.2 | 89.5 | 89.4 | 86.9 | 78.9 | 80.8 | 89.8 | 76.6 | 76.1 | 76.4 | 74.1 | 72.6 | 72.3 | 66.1 | 70.9 | 73.9 | 44.4 | 44.4 | 57.5 | 54.3 | 76.3 |
| Verb tense/aspect/mood | 10120.0 | 78.6 | 81.8 | 79.2 | 83.4 | 73.0 | 68.2 | 80.6 | 77.4 | 77.3 | 81.8 | 67.5 | 52.2 | 72.1 | 67.8 | 67.8 | 65.9 | 66.7 | 73.3 | 63.0 | 63.6 | 66.0 | 62.6 | 51.8 | 51.8 | 59.1 | 52.7 | 68.7 |
| Verb valency | 3486.0 | 80.8 | 84.6 | 84.5 | 81.7 | 81.7 | 77.0 | 83.8 | 82.9 | 80.3 | 80.9 | 73.7 | 75.5 | 73.2 | 67.4 | 71.1 | 67.9 | 71.6 | 67.4 | 67.2 | 70.5 | 71.3 | 62.3 | 61.2 | 61.2 | 55.0 | 53.9 | 72.6 |
| macro avg. | 78727.0 | 81.0 | 80.7 | 80.5 | 80.2 | 79.9 | 79.9 | 78.7 | 78.5 | 78.1 | 76.4 | 74.1 | 73.9 | 73.7 | 73.3 | 72.5 | 71.9 | 70.8 | 69.7 | 68.6 | 68.5 | 66.5 | 66.5 | 56.9 | 56.9 | 56.6 | 53.4 | 71.8 |
| micro avg. | 78727.0 | 79.1 | 79.4 | 78.6 | 77.9 | 77.1 | 76.0 | 73.3 | 73.1 | 74.2 | 72.0 | 72.8 | 67.3 | 66.4 | 70.9 | 70.7 | 68.8 | 69.5 | 64.8 | 68.0 | 67.9 | 67.9 | 64.3 | 53.9 | 54.4 | 54.4 | 51.9 | 69.0 |

Table 3: Accuracy of the metrics(%) with regards to the linguistically-motivated categories for English-Russian

| ling. category | # | MetricX-24 | metametrics | MetricX-24-Hybrid | XCOMET | BLEURT-20 | COMET-22 | CometKiwi-XXL | MetricX-24-QE | XCOMET-QE | MetricX-24-Hybrid-QE | Yisi-1 | CometKiwi | BERTScore | sentinel-cand-mqm | chrFS | chrF | spBLEU | BLEU | damonmonli | monmonli | gemba | XLsimMqm | XLsimDA | PrismRefSmall | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | 3788.0 | **96.9** | 96.4 | 95.1 | 93.2 | 89.8 | 87.4 | 80.9 | 96.3 | 83.8 | 91.4 | 82.6 | 77.2 | 75.3 | 87.1 | 74.7 | 73.1 | 70.6 | 68.9 | 80.5 | 78.4 | 89.4 | 43.9 | 43.9 | 48.1 | 45.3 | 78.0 |
| Coordination & ellipsis | 2273.0 | 80.6 | 79.3 | **81.5** | 80.4 | 74.9 | 76.6 | 81.0 | **81.4** | 77.2 | **81.8** | 68.6 | 78.6 | 68.1 | 75.5 | 63.5 | 61.5 | 62.7 | 62.7 | 65.4 | 66.3 | 60.1 | 52.5 | 52.5 | 47.7 | 48.7 | 69.2 |
| False friends | 2414.0 | **87.8** | 83.7 | 86.3 | 76.3 | 82.4 | 83.0 | 69.1 | 69.2 | 68.6 | 68.4 | **87.7** | 52.4 | 76.3 | 58.0 | 84.9 | 83.2 | 80.7 | 75.8 | 80.8 | 62.2 | 34.0 | 43.2 | 43.2 | 53.1 | 42.0 | 69.3 |
| Function word | 2433.0 | 82.5 | 78.0 | 73.4 | 84.1 | 79.7 | 81.4 | 83.0 | 85.7 | **86.3** | 71.8 | 65.7 | 79.3 | 69.3 | 82.7 | 64.8 | 60.3 | 65.6 | 73.2 | 56.4 | 57.0 | 56.3 | 73.7 | 73.7 | 50.3 | 49.0 | 71.3 |
| LDD & interrogatives | 1939.0 | 85.4 | 86.0 | **87.8** | 84.8 | 81.8 | 82.6 | 84.9 | 87.3 | 83.4 | **87.6** | 65.5 | 77.6 | 66.2 | **87.8** | 62.5 | 59.9 | 62.0 | 61.5 | 55.1 | 58.3 | 68.1 | 54.2 | 54.2 | 51.6 | 46.4 | 71.3 |
| MWE | 9602.0 | 82.9 | 82.9 | 82.9 | 81.2 | 80.5 | 81.4 | 82.2 | 81.6 | 77.7 | **83.9** | 77.0 | 75.6 | 74.6 | 72.3 | 75.5 | 73.1 | 72.8 | 70.9 | 67.9 | 65.1 | 69.5 | 53.5 | 53.5 | 51.7 | 51.0 | 72.8 |
| Named entity & terminology | 16284.0 | 82.8 | **84.9** | 81.6 | 80.6 | **84.9** | 84.3 | 71.6 | 72.2 | 69.5 | 71.6 | 83.0 | 71.3 | 78.8 | 70.6 | 80.3 | 78.7 | 78.3 | 72.5 | 72.9 | 67.7 | 64.1 | 47.6 | 47.6 | 53.6 | 52.1 | 72.1 |
| Negation | 346.0 | 65.3 | 59.8 | 58.7 | 49.4 | 72.3 | 67.3 | 58.7 | 57.8 | 45.4 | 44.5 | 79.5 | 49.4 | 80.3 | 41.3 | 82.9 | **83.5** | 74.3 | 72.3 | 70.5 | 72.5 | 42.5 | 49.7 | 49.7 | 49.4 | 45.7 | 60.9 |
| Non-verbal agreement | 6755.0 | **86.4** | 81.5 | 84.4 | 82.3 | 79.6 | 77.4 | 78.7 | 83.0 | 80.9 | 81.5 | 72.1 | 77.6 | 68.7 | 73.1 | 69.4 | 68.2 | 67.4 | 64.6 | 59.6 | 60.9 | 68.4 | 68.4 | 68.4 | 51.2 | 47.5 | 72.1 |
| Punctuation | 363.0 | 73.0 | 71.1 | 72.7 | 71.3 | 68.6 | **76.0** | **75.8** | 63.6 | 70.8 | 67.2 | **75.8** | 72.2 | 73.3 | 70.5 | 62.0 | 58.4 | 64.7 | 60.9 | 51.0 | 64.5 | 60.9 | 57.3 | 57.3 | 46.3 | 45.5 | 65.2 |
| Subordination | 6625.0 | 74.7 | 74.5 | 71.4 | 75.0 | 72.7 | 77.1 | **78.4** | 72.5 | 75.2 | 71.7 | 69.3 | 68.6 | 66.9 | 73.5 | 63.8 | 62.4 | 64.3 | 64.0 | 56.4 | 63.3 | 53.6 | 50.5 | 50.5 | 51.0 | 48.1 | 66.0 |
| Verb semantics | 275.0 | **88.0** | 82.2 | **88.0** | 85.5 | 86.5 | 74.2 | 79.6 | 75.3 | 80.0 | 76.0 | 53.1 | 69.8 | 55.6 | 55.3 | 60.7 | 65.1 | 53.8 | 48.7 | 68.4 | 66.5 | 72.0 | 33.5 | 33.5 | 65.5 | 67.3 | 67.4 |
| Verb tense/aspect/mood | 2994.0 | 85.0 | **86.0** | 82.6 | **86.2** | 75.5 | 79.7 | **85.8** | 82.8 | 80.7 | 79.3 | 69.7 | 72.6 | 68.7 | 70.4 | 68.1 | 66.7 | 68.8 | 63.1 | 60.1 | 55.9 | 61.6 | 47.5 | 47.5 | 50.6 | 51.3 | 69.9 |
| Verb valency | 3022.0 | **83.3** | 82.3 | 80.4 | **83.5** | 76.8 | 76.3 | 80.9 | 82.0 | 81.6 | 81.5 | 69.6 | 73.8 | 72.8 | 72.0 | 72.2 | 71.8 | 69.0 | 67.7 | 66.6 | 64.2 | 66.9 | 60.7 | 60.7 | 51.6 | 46.7 | 71.8 |
| macro avg. | 59113.0 | **82.5** | 80.6 | 80.5 | 79.6 | 79.0 | 78.9 | 77.9 | 77.9 | 75.8 | 75.6 | 72.8 | 71.1 | 71.1 | 70.7 | 70.4 | 69.0 | 68.2 | 66.2 | 65.1 | 64.5 | 62.0 | 52.6 | 52.6 | 51.6 | 49.0 | 69.8 |
| micro avg. | 59113.0 | **83.4** | 82.8 | 81.7 | 81.4 | 80.7 | 81.0 | 77.8 | 78.7 | 76.2 | 77.5 | 75.8 | 72.9 | 72.9 | 73.0 | 73.1 | 71.4 | 71.2 | 68.5 | 66.8 | 64.8 | 64.4 | 52.8 | 52.8 | 51.7 | 49.3 | 71.3 |

# B  Accuracies per phenomenon

Table 4: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-German

| ling. category | ling. phenomenon | # | XCOMET | MetricX-24 | MetricX-24-Hybrid | metametrics | MetricX-24-QE | MetricX-24-Hybrid-QE | XCOMET-QE | CometKiwi-XXL | BLEURT-20 | CometKiwi | COMET-22 | chrF++ | MEE4 | chrF | gemba | BERTScore | Yisi-1 | spBLEU | BLEU | monmonti | damonmonti | sentinel-cand-mqm | PrismRefSmall | XLsimDA | XLsimMqm | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | Lexical ambiguity | 4614 | 80 | 86 | 85 | 90 | 75 | 71 | 62 | 61 | 90 | 48 | 89 | 84 | 82 | 85 | 70 | 78 | 89 | 80 | 64 | 83 | 84 | 78 | 68 | 55 | 55 | 61 | 75 |
| Coordination \& ellipsis | Gapping | 605 | 86 | 87 | 91 | 84 | 87 | 90 | 77 | 86 | 87 | 85 | 84 | 70 | 69 | 74 | 73 | 75 | 78 | 70 | 64 | 47 | 47 | 87 | 62 | 47 | 47 | 59 | 74 |
|  | Pseudogapping | 1565 | 87 | 85 | 94 | 85 | 93 | 91 | 82 | 93 | 88 | 78 | 87 | 72 | 72 | 71 | 82 | 76 | 73 | 70 | 70 | 63 | 64 | 87 | 53 | 39 | 39 | 50 | 75 |
|  | Right node raising | 647 | 74 | 67 | 76 | 66 | 75 | 68 | 75 | 76 | 66 | 64 | 75 | 59 | 61 | 55 | 81 | 64 | 66 | 54 | 55 | 67 | 55 | 76 | 55 | 53 | 53 | 52 | 65 |
|  | Sluicing | 472 | 74 | 57 | 61 | 59 | 62 | 64 | 57 | 62 | 57 | 50 | 61 | 59 | 62 | 52 | 34 | 58 | 63 | 54 | 57 | 71 | 61 | 59 | 46 | 53 | 53 | 46 | 57 |
|  | Stripping | 545 | 59 | 60 | 66 | 66 | 54 | 64 | 59 | 59 | 69 | 50 | 66 | 53 | 56 | 47 | 40 | 62 | 71 | 48 | 56 | 64 | 71 | 59 | 36 | 57 | 57 | 44 | 59 |
|  | VP-ellipsis | 539 | 76 | 66 | 71 | 64 | 68 | 71 | 79 | 80 | 68 | 65 | 65 | 48 | 45 | 75 | 51 | 50 | 55 | 66 | 45 | 38 | 49 | 71 | 36 | 69 | 69 | 40 | 61 |
| False friends | False friends | 1389 | 74 | 78 | 80 | 78 | 52 | 74 | 71 | 85 | 73 | 72 | 86 | 77 | 80 | 74 | 82 | 69 | 73 | 74 | 64 | 64 | 65 | 77 | 74 | 47 | 78 | 59 | 72 |
| Function word | Focus particle | 333 | 62 | 61 | 65 | 58 | 61 | 47 | 61 | 55 | 57 | 69 | 71 | 77 | 70 | 74 | 31 | 69 | 62 | 74 | 59 | 35 | 40 | 31 | 47 | 56 | 56 | 53 | 57 |
|  | Question tag | 1567 | 91 | 85 | 81 | 87 | 88 | 85 | 87 | 87 | 87 | 90 | 91 | 74 | 78 | 74 | 79 | 90 | 75 | 74 | 83 | 65 | 70 | 97 | 52 | 58 | 58 | 66 | 70 |
| LDD \& interrogatives | Extraposition | 85 | 79 | 84 | 86 | 76 | 80 | 76 | 75 | 78 | 67 | 76 | 64 | 72 | 71 | 72 | 49 | 66 | 65 | 66 | 61 | 68 | 62 | 74 | 65 | 58 | 58 | 49 | 70 |
|  | Inversion | 117 | 81 | 84 | 87 | 79 | 85 | 84 | 92 | 82 | 78 | 85 | 79 | 67 | 74 | 67 | 68 | 79 | 75 | 68 | 62 | 75 | 81 | 79 | 52 | 65 | 85 | 40 | 76 |
|  | Multiple connectors | 25 | 44 | 44 | 56 | 40 | 44 | 56 | 32 | 44 | 44 | 40 | 36 | 44 | 64 | 44 | 0 | 72 | 72 | 53 | 68 | 57 | 28 | 44 | 36 | 80 | 80 | 40 | 51 |
|  | Negative inversion | 358 | 76 | 72 | 73 | 74 | 73 | 74 | 82 | 74 | 78 | 67 | 61 | 46 | 55 | 44 | 55 | 84 | 42 | 53 | 55 | 57 | 60 | 47 | 39 | 65 | 65 | 35 | 61 |
|  | Pied-piping | 49 | 98 | 100 | 94 | 94 | 100 | 83 | 98 | 88 | 94 | 90 | 92 | 71 | 80 | 71 | 69 | 85 | 67 | 88 | 88 | 71 | 76 | 63 | 59 | 71 | 71 | 51 | 82 |
|  | Polar question | 13 | 92 | 77 | 77 | 77 | 92 | 69 | 69 | 92 | 92 | 77 | 100 | 77 | 85 | 77 | 100 | 100 | 85 | 77 | 77 | 38 | 54 | 100 | 62 | 23 | 23 | 54 | 74 |
|  | Preposition stranding | 9 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 78 | 100 | 78 | 100 | 100 | 100 | 67 | 56 | 82 | 89 | 100 | 100 | 89 | 89 | 100 | 93 |
|  | Split infinitive | 88 | 91 | 98 | 100 | 94 | 100 | 100 | 100 | 94 | 100 | 86 | 89 | 56 | 62 | 51 | 60 | 67 | 67 | 50 | 50 | 75 | 89 | 100 | 59 | 67 | 67 | 41 | 76 |
|  | Topicalization | 192 | 86 | 83 | 90 | 95 | 90 | 85 | 81 | 89 | 82 | 96 | 91 | 65 | 68 | 64 | 52 | 64 | 67 | 68 | 68 | 52 | 61 | 92 | 59 | 45 | 45 | 60 | 73 |
|  | Wh-movement | 66 | 74 | 77 | 91 | 73 | 80 | 85 | 67 | 74 | 83 | 73 | 77 | 70 | 68 | 67 | 55 | 64 | 69 | 61 | 50 | 91 | 95 | 65 | 62 | 42 | 42 | 58 | 68 |
| MWE | Collocation | 506 | 90 | 90 | 88 | 88 | 79 | 78 | 86 | 81 | 84 | 91 | 84 | 80 | 82 | 77 | 76 | 78 | 83 | 74 | 75 | 77 | 67 | 80 | 59 | 62 | 62 | 52 | 78 |
|  | Compound | 257 | 94 | 92 | 95 | 95 | 94 | 98 | 86 | 88 | 91 | 88 | 84 | 80 | 72 | 72 | 80 | 78 | 91 | 78 | 75 | 77 | 74 | 96 | 67 | 67 | 67 | 46 | 78 |
|  | Idiom | 2746 | 91 | 93 | 93 | 93 | 90 | 89 | 87 | 88 | 93 | 81 | 95 | 84 | 85 | 84 | 97 | 85 | 89 | 85 | 81 | 86 | 89 | 90 | 60 | 74 | 74 | 57 | 85 |
|  | Nominal MWE | 1522 | 75 | 77 | 76 | 72 | 69 | 67 | 64 | 62 | 71 | 62 | 64 | 68 | 58 | 64 | 59 | 63 | 65 | 60 | 59 | 64 | 70 | 78 | 62 | 53 | 53 | 56 | 65 |
|  | Prepositional MWE | 353 | 92 | 88 | 80 | 89 | 76 | 72 | 86 | 82 | 82 | 84 | 91 | 69 | 74 | 73 | 69 | 75 | 79 | 78 | 68 | 67 | 51 | 79 | 52 | 59 | 59 | 50 | 74 |
|  | Verbal MWE | 432 | 84 | 88 | 88 | 83 | 88 | 81 | 78 | 88 | 80 | 72 | 65 | 66 | 68 | 64 | 84 | 72 | 69 | 65 | 65 | 71 | 66 | 70 | 58 | 86 | 86 | 57 | 75 |
| Named entity \& terminology | Date | 2010 | 69 | 84 | 75 | 89 | 69 | 66 | 56 | 63 | 82 | 64 | 78 | 76 | 73 | 76 | 64 | 77 | 76 | 71 | 80 | 65 | 66 | 75 | 58 | 64 | 64 | 56 | 71 |
|  | Domainspecific Term | 7405 | 79 | 80 | 78 | 91 | 66 | 73 | 66 | 72 | 78 | 63 | 83 | 71 | 76 | 69 | 63 | 68 | 86 | 76 | 65 | 64 | 71 | 55 | 52 | 46 | 46 | 46 | 69 |
|  | Location | 2731 | 70 | 92 | 85 | 88 | 32 | 34 | 13 | 15 | 91 | 31 | 93 | 83 | 89 | 76 | 32 | 80 | 87 | 85 | 55 | 47 | 90 | 69 | 65 | 48 | 48 | 54 | 65 |
|  | Measuring unit | 8539 | 58 | 61 | 60 | 53 | 45 | 55 | 75 | 54 | 44 | 18 | 59 | 88 | 65 | 83 | 31 | 80 | 62 | 85 | 58 | 80 | 80 | 34 | 41 | 43 | 43 | 51 | 55 |
|  | Proper name | 2206 | 76 | 74 | 76 | 72 | 75 | 73 | 75 | 73 | 69 | 44 | 70 | 63 | 64 | 64 | 32 | 65 | 72 | 66 | 59 | 63 | 71 | 71 | 49 | 63 | 63 | 48 | 67 |
| Negation | Negation | 506 | 91 | 90 | 93 | 89 | 94 | 91 | 91 | 93 | 93 | 95 | 93 | 83 | 91 | 77 | 70 | 86 | 88 | 74 | 74 | 88 | 87 | 75 | 58 | 58 | 58 | 58 | 83 |
| Non-verbal agreement | Coreference | 3340 | 85 | 80 | 78 | 82 | 86 | 79 | 80 | 85 | 82 | 65 | 89 | 77 | 75 | 78 | 97 | 68 | 67 | 76 | 79 | 47 | 55 | 79 | 57 | 67 | 67 | 56 | 72 |
|  | Genitive | 483 | 82 | 93 | 82 | 94 | 86 | 87 | 82 | 75 | 67 | 78 | 78 | 77 | 70 | 77 | 76 | 75 | 81 | 76 | 70 | 84 | 77 | 74 | 58 | 64 | 64 | 49 | 75 |
|  | Lexical Morphology/Functional shift | 2330 | 91 | 94 | 97 | 95 | 90 | 94 | 92 | 92 | 90 | 80 | 91 | 79 | 81 | 77 | 76 | 76 | 85 | 70 | 65 | 84 | 77 | 96 | 57 | 64 | 56 | 50 | 81 |
|  | Lexical Morphology/Noun formation (er) | 2067 | 72 | 77 | 70 | 76 | 66 | 64 | 68 | 63 | 77 | 65 | 68 | 63 | 60 | 65 | 55 | 61 | 64 | 61 | 54 | 64 | 65 | 65 | 45 | 74 | 74 | 44 | 65 |

Continued on next page

Table 4: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-German

| ling. category | ling. phenomenon | # | XCOMET | MetricX-24 | MetricX-24-Hybrid | metametrics | MetricX-24-QE | MetricX-24-Hybrid-QE | XCOMET-QE | CometKiwi-XXL | BLEURT-20 | CometKiwi | COMET-22 | chrF++ | MEE4 | chrF | gemba | BERTScore | YiSi-1 | spBLEU | BLEU | momonli | damonmonli | sentinel-cand-mqm | PrismRefSmall | XLsimDA | XLsimMqm | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Personal Pronoun Coreference | 4632 | 83 | 80 | 93 | 67 | 91 | 94 | 83 | 83 | 76 | 79 | 65 | 52 | 53 | 52 | 84 | 54 | 56 | 53 | 50 | 54 | 52 | 75 | 49 | 47 | 47 | 48 | 66 |
| | Possession | 555 | 89 | 88 | 88 | 86 | 89 | 88 | 88 | 90 | 79 | 87 | 82 | 70 | 76 | 74 | 83 | 72 | 73 | 68 | 67 | 81 | 76 | 66 | 49 | 49 | 49 | 50 | 75 |
| Punctuation | Substitution | 2090 | 65 | 67 | 67 | 68 | 71 | 67 | 60 | 69 | 67 | 68 | 64 | 64 | 62 | 65 | 49 | 62 | 59 | 65 | 64 | 53 | 56 | 69 | 41 | 69 | 69 | 45 | 62 |
| | Quotation marks | 2435 | 63 | 64 | 62 | 65 | 70 | 66 | 65 | 70 | 67 | 71 | 72 | 69 | 70 | 72 | 44 | 66 | 72 | 67 | 69 | 51 | 67 | 80 | 51 | 50 | 50 | 51 | 64 |
| Subordination | Adverbial clause | 583 | 92 | 90 | 92 | 89 | 93 | 93 | 65 | 94 | 83 | 84 | 87 | 69 | 72 | 68 | 81 | 66 | 73 | 67 | 75 | 62 | 66 | 88 | 57 | 59 | 54 | 51 | 77 |
| | Cleft sentence | 578 | 78 | 75 | 77 | 76 | 81 | 82 | 98 | 70 | 72 | 71 | 77 | 64 | 68 | 64 | 64 | 67 | 73 | 64 | 66 | 70 | 53 | 80 | 62 | 54 | 54 | 47 | 69 |
| | Contact clause | 788 | 98 | 91 | 94 | 96 | 97 | 97 | 89 | 97 | 96 | 99 | 95 | 74 | 79 | 78 | 88 | 80 | 90 | 72 | 73 | 77 | 78 | 97 | 62 | 54 | 54 | 52 | 82 |
| | Indirect speech | 113 | 88 | 79 | 79 | 78 | 65 | 70 | 92 | 78 | 65 | 85 | 75 | 53 | 66 | 53 | 56 | 65 | 62 | 58 | 58 | 46 | 59 | 61 | 50 | 38 | 35 | 59 | 64 |
| | Infinitive clause | 454 | 89 | 87 | 91 | 87 | 91 | 93 | 69 | 87 | 61 | 91 | 87 | 72 | 72 | 64 | 68 | 73 | 80 | 73 | 67 | 81 | 84 | 61 | 50 | 59 | 59 | 44 | 77 |
| | Object clause | 111 | 95 | 85 | 95 | 74 | 97 | 95 | 98 | 86 | 72 | 58 | 45 | 64 | 82 | 59 | 73 | 50 | 50 | 47 | 49 | 65 | 54 | 55 | 42 | 39 | 39 | 38 | 64 |
| | Pseudo-cleft sentence | 578 | 76 | 78 | 76 | 66 | 70 | 66 | 89 | 80 | 61 | 58 | 66 | 78 | 79 | 80 | 66 | 78 | 74 | 77 | 82 | 66 | 68 | 66 | 52 | 42 | 42 | 53 | 69 |
| | Relative clause | 560 | 95 | 96 | 96 | 95 | 95 | 96 | 98 | 97 | 94 | 98 | 93 | 75 | 79 | 76 | 94 | 84 | 75 | 73 | 76 | 67 | 68 | 81 | 56 | 65 | 65 | 53 | 82 |
| | Subject clause | 933 | 91 | 93 | 94 | 92 | 97 | 95 | 83 | 88 | 91 | 93 | 87 | 78 | 84 | 81 | 56 | 84 | 88 | 81 | 81 | 74 | 64 | 95 | 71 | 68 | 68 | 66 | 79 |
| Verb tense/aspect/mood | Conditional | 975 | 87 | 73 | 65 | 90 | 78 | 62 | 89 | 81 | 86 | 93 | 82 | 61 | 63 | 60 | 42 | 60 | 69 | 49 | 53 | 86 | 86 | 95 | 67 | 68 | 68 | 56 | 71 |
| | Ditransitive - conditional I progressive | 83 | 89 | 73 | 69 | 61 | 70 | 66 | 77 | 93 | 59 | 64 | 55 | 65 | 73 | 67 | 73 | 65 | 64 | 67 | 63 | 43 | 35 | 72 | 58 | 42 | 42 | 59 | 62 |
| | Ditransitive - conditional I simple | 197 | 91 | 82 | 86 | 71 | 94 | 92 | 90 | 87 | 70 | 87 | 74 | 70 | 69 | 70 | 85 | 68 | 71 | 69 | 60 | 74 | 71 | 55 | 52 | 45 | 45 | 49 | 72 |
| | Ditransitive - conditional II progressive | 130 | 91 | 92 | 92 | 87 | 81 | 95 | 65 | 92 | 85 | 72 | 78 | 78 | 82 | 69 | 48 | 75 | 76 | 61 | 72 | 71 | 71 | 61 | 55 | 51 | 51 | 58 | 75 |
| | Ditransitive - conditional II simple | 108 | 85 | 88 | 93 | 78 | 59 | 84 | 77 | 79 | 66 | 71 | 70 | 58 | 74 | 69 | 59 | 76 | 62 | 60 | 61 | 67 | 67 | 72 | 40 | 48 | 48 | 49 | 68 |
| | Ditransitive - future I progressive | 244 | 82 | 70 | 76 | 68 | 59 | 78 | 89 | 89 | 75 | 47 | 58 | 61 | 58 | 61 | 30 | 57 | 60 | 60 | 62 | 49 | 51 | 73 | 52 | 48 | 48 | 50 | 62 |
| | Ditransitive - future I simple | 217 | 78 | 70 | 77 | 61 | 63 | 70 | 91 | 76 | 61 | 47 | 45 | 54 | 54 | 53 | 28 | 59 | 54 | 51 | 58 | 38 | 38 | 46 | 41 | 34 | 34 | 51 | 55 |
| | Ditransitive - future II progressive | 210 | 94 | 91 | 82 | 89 | 88 | 86 | 89 | 87 | 67 | 90 | 66 | 68 | 81 | 76 | 92 | 73 | 75 | 75 | 75 | 31 | 65 | 58 | 46 | 68 | 68 | 46 | 75 |
| | Ditransitive - future II simple | 84 | 79 | 94 | 94 | 85 | 70 | 70 | 89 | 89 | 90 | 90 | 69 | 55 | 89 | 59 | 96 | 87 | 87 | 83 | 82 | 70 | 71 | 40 | 57 | 57 | 57 | 60 | 80 |
| | Ditransitive - past perfect progressive | 122 | 65 | 66 | 61 | 56 | 79 | 92 | 47 | 57 | 50 | 75 | 66 | 55 | 64 | 59 | 64 | 51 | 54 | 47 | 51 | 52 | 54 | 55 | 46 | 62 | 62 | 48 | 59 |
| | Ditransitive - past perfect simple | 160 | 61 | 67 | 61 | 57 | 59 | 86 | 64 | 71 | 63 | 51 | 52 | 61 | 56 | 57 | 20 | 57 | 60 | 51 | 53 | 49 | 50 | 71 | 52 | 42 | 42 | 51 | 56 |
| | Ditransitive - past progressive | 218 | 74 | 71 | 73 | 71 | 53 | 61 | 71 | 71 | 61 | 50 | 54 | 50 | 51 | 50 | 50 | 55 | 54 | 47 | 53 | 48 | 60 | 55 | 49 | 54 | 54 | 54 | 57 |
| | Ditransitive - present perfect progressive | 107 | 97 | 85 | 82 | 81 | 82 | 62 | 98 | 93 | 74 | 81 | 68 | 52 | 66 | 69 | 56 | 64 | 73 | 46 | 48 | 66 | 66 | 80 | 49 | 62 | 62 | 53 | 70 |
| | Ditransitive - present perfect simple | 185 | 90 | 71 | 77 | 68 | 60 | 82 | 99 | 96 | 66 | 40 | 52 | 52 | 54 | 52 | 19 | 61 | 57 | 52 | 56 | 43 | 50 | 62 | 35 | 41 | 41 | 46 | 59 |
| | Ditransitive - present progressive | 114 | 98 | 87 | 86 | 86 | 97 | 89 | 99 | 96 | 80 | 99 | 89 | 68 | 72 | 71 | 97 | 61 | 78 | 54 | 48 | 89 | 89 | 83 | 60 | 71 | 71 | 62 | 80 |
| | Ditransitive - simple past | 199 | 94 | 82 | 87 | 73 | 99 | 97 | 91 | 96 | 83 | 91 | 74 | 62 | 67 | 68 | 93 | 67 | 65 | 55 | 54 | 53 | 53 | 53 | 43 | 69 | 69 | 51 | 72 |
| | Ditransitive - simple present | 133 | 92 | 81 | 73 | 82 | 86 | 90 | 95 | 94 | 79 | 83 | 83 | 61 | 70 | 68 | 81 | 67 | 60 | 50 | 50 | 76 | 76 | 79 | 73 | 71 | 69 | 50 | 74 |
| | Gerund | 1119 | 98 | 98 | 98 | 98 | 98 | 98 | 97 | 95 | 95 | 97 | 95 | 73 | 82 | 78 | 87 | 79 | 82 | 73 | 74 | 80 | 69 | 60 | 66 | 22 | 22 | 52 | 80 |
| | Imperative | 259 | 90 | 75 | 81 | 83 | 90 | 77 | 88 | 90 | 80 | 89 | 86 | 72 | 78 | 76 | 88 | 69 | 74 | 68 | 62 | 79 | 69 | 63 | 57 | 50 | 50 | 59 | 75 |
| | Intransitive - conditional I progressive | 23 | 70 | 91 | 91 | 74 | 27 | 70 | 48 | 39 | 87 | 52 | 100 | 61 | 78 | 78 | 96 | 93 | 78 | 78 | 61 | 83 | 83 | 20 | 39 | 83 | 83 | 73 | 72 |
| | Intransitive - conditional I simple | 15 | 93 | 87 | 100 | 100 | 60 | 40 | 20 | 20 | 93 | 47 | 80 | 80 | 67 | 73 | 73 | 87 | 73 | 83 | 7 | 73 | 73 | 67 | 67 | 93 | 93 | 73 | 71 |
| | Intransitive - conditional II progressive | 5 | 80 | 100 | 100 | 100 | 100 | 100 | 60 | 60 | 80 | 60 | 80 | 80 | 80 | 100 | 60 | 100 | 60 | 100 | 100 | 60 | 60 | 40 | 80 | 80 | 80 | 80 | 81 |
| | Intransitive - conditional II simple | 2 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 50 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 98 |
| | Intransitive - future I progressive | 19 | 100 | 100 | 100 | 84 | 84 | 84 | 74 | 74 | 89 | 100 | 88 | 89 | 86 | 84 | 79 | 90 | 74 | 86 | 89 | 68 | 63 | 100 | 47 | 79 | 79 | 100 | 84 |
| | Intransitive - future I simple | 50 | 78 | 88 | 88 | 88 | 52 | 50 | 72 | 48 | 92 | 56 | 84 | 84 | 78 | 84 | 82 | 90 | 72 | 86 | 60 | 82 | 82 | 58 | 60 | 82 | 79 | 47 | 76 |
| | Intransitive - future II progressive | 18 | 78 | 78 | 78 | 78 | 56 | 50 | 61 | 61 | 90 | 72 | 72 | 89 | 78 | 89 | 72 | 83 | 44 | 83 | 89 | 72 | 72 | 22 | 72 | 50 | 50 | 72 | 76 |
| | Intransitive - future II simple | 15 | 93 | 93 | 93 | 93 | 93 | 93 | 87 | 93 | 80 | 93 | 67 | 93 | 93 | 93 | 80 | 93 | 67 | 100 | 100 | 80 | 80 | 20 | 100 | 13 | 13 | 80 | 81 |

524

Table 4: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-German

| ling. phenomenon | # | XCOMET | MetricX-24 | MetricX-24-Hybrid | metametrics | MetricX-24-QE | MetricX-24-Hybrid-QE | XCOMET-QE | CometKiwi-XXL | BLEURT-20 | CometKiwi | COMET-22 | chrF++ | MEE4 | chrF | gemba | BERTScore | YiSi-1 | spBLEU | BLEU | momonolt | damonolt | sentinel-cand-mqm | PrismRefSmall | XLsimDA | XLsimMqm | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Intransitive - past perfect progressive | 81 | 38 | 56 | 52 | 48 | 20 | 16 | 27 | 25 | 60 | 70 | 65 | 63 | 70 | 72 | 51 | 64 | 74 | 67 | 67 | 69 | 65 | 75 | 49 | 69 | 69 | 62 | 56 |
| Intransitive - past perfect simple | 31 | 74 | 71 | 65 | 87 | 68 | 68 | 48 | 58 | 77 | 81 | 71 | 84 | 87 | 87 | 45 | 94 | 71 | 81 | 84 | 45 | 55 | 74 | 71 | 32 | 32 | 77 | 69 |
| Intransitive - past progressive | 79 | 70 | 70 | 78 | 76 | 76 | 56 | 67 | 66 | 62 | 61 | 71 | 56 | 65 | 61 | 67 | 65 | 52 | 62 | 61 | 59 | 62 | 34 | 58 | 75 | 75 | 56 | 64 |
| Intransitive - present perfect simple | 20 | 95 | 100 | 100 | 100 | 100 | 100 | 100 | 95 | 100 | 95 | 100 | 95 | 100 | 100 | 100 | 100 | 70 | 85 | 85 | 75 | 85 | 80 | 90 | 90 | 90 | 80 | 93 |
| Intransitive - present progressive | 26 | 92 | 100 | 100 | 92 | 91 | 58 | 81 | 69 | 69 | 69 | 92 | 77 | 92 | 96 | 100 | 85 | 77 | 85 | 85 | 75 | 65 | 62 | 70 | 46 | 46 | 73 | 77 |
| Intransitive - simple past | 53 | 77 | 83 | 89 | 60 | 67 | 89 | 74 | 62 | 58 | 58 | 53 | 62 | 62 | 57 | 75 | 45 | 63 | 38 | 38 | 72 | 68 | 45 | 74 | 62 | 62 | 75 | 65 |
| Intransitive - simple present | 27 | 63 | 67 | 78 | 63 | 38 | 81 | 62 | 67 | 74 | 89 | 78 | 44 | 44 | 48 | 75 | 63 | 69 | 37 | 33 | 70 | 74 | 74 | 56 | 30 | 30 | 67 | 63 |
| Modal | 16 | 81 | 75 | 62 | 62 | 92 | 89 | 92 | 62 | 81 | 38 | 75 | 88 | 81 | 88 | 56 | 81 | 73 | 62 | 69 | 62 | 81 | 25 | 75 | 69 | 69 | 88 | 67 |
| Modal negated | 52 | 89 | 90 | 85 | 83 | 92 | 81 | 62 | 73 | 81 | 87 | 85 | 67 | 75 | 73 | 67 | 60 | 73 | 68 | 63 | 77 | 79 | 54 | 78 | 67 | 67 | 56 | 74 |
| Reflexive - conditional I progressive | 150 | 71 | 81 | 77 | 85 | 82 | 69 | 82 | 74 | 54 | 49 | 49 | 59 | 65 | 81 | 100 | 56 | 56 | 68 | 60 | 58 | 58 | 35 | 58 | 62 | 62 | 57 | 67 |
| Reflexive - conditional I simple | 141 | 78 | 78 | 69 | 77 | 67 | 71 | 62 | 68 | 59 | 38 | 45 | 61 | 61 | 60 | 82 | 56 | 50 | 65 | 62 | 71 | 62 | 35 | 65 | 59 | 59 | 47 | 62 |
| Reflexive - conditional II progressive | 204 | 92 | 88 | 88 | 89 | 72 | 91 | 89 | 91 | 69 | 78 | 61 | 72 | 74 | 74 | 95 | 65 | 66 | 67 | 70 | 63 | 58 | 36 | 72 | 52 | 52 | 55 | 72 |
| Reflexive - conditional II simple | 336 | 92 | 97 | 85 | 95 | 86 | 85 | 96 | 79 | 71 | 69 | 54 | 71 | 67 | 71 | 92 | 61 | 65 | 67 | 73 | 71 | 56 | 55 | 71 | 38 | 38 | 45 | 69 |
| Reflexive - future I progressive | 212 | 76 | 72 | 74 | 82 | 59 | 66 | 86 | 79 | 58 | 46 | 59 | 71 | 58 | 65 | 96 | 70 | 57 | 67 | 58 | 67 | 65 | 50 | 59 | 68 | 68 | 50 | 66 |
| Reflexive - future I simple | 160 | 73 | 85 | 75 | 81 | 68 | 69 | 57 | 48 | 73 | 57 | 74 | 84 | 76 | 80 | 98 | 84 | 64 | 78 | 71 | 86 | 81 | 50 | 78 | 85 | 85 | 57 | 74 |
| Reflexive - future II progressive | 158 | 73 | 80 | 70 | 76 | 66 | 68 | 68 | 64 | 70 | 82 | 71 | 81 | 73 | 75 | 68 | 62 | 78 | 68 | 65 | 67 | 67 | 59 | 61 | 83 | 83 | 51 | 70 |
| Reflexive - future II simple | 123 | 81 | 89 | 89 | 72 | 92 | 81 | 82 | 76 | 70 | 67 | 51 | 76 | 70 | 75 | 89 | 58 | 60 | 63 | 67 | 76 | 64 | 31 | 68 | 48 | 48 | 50 | 69 |
| Reflexive - past perfect progressive | 162 | 84 | 86 | 79 | 80 | 87 | 81 | 80 | 75 | 73 | 67 | 69 | 89 | 70 | 75 | 89 | 58 | 70 | 61 | 61 | 76 | 62 | 51 | 51 | 50 | 50 | 70 | 72 |
| Reflexive - past perfect simple | 169 | 75 | 77 | 73 | 80 | 74 | 68 | 79 | 85 | 66 | 76 | 57 | 92 | 54 | 56 | 67 | 53 | 57 | 51 | 50 | 66 | 52 | 52 | 51 | 58 | 58 | 49 | 63 |
| Reflexive - past progressive | 843 | 73 | 73 | 65 | 70 | 50 | 62 | 73 | 77 | 68 | 55 | 56 | 60 | 56 | 58 | 91 | 56 | 57 | 58 | 57 | 65 | 61 | 18 | 58 | 59 | 59 | 47 | 61 |
| Reflexive - present perfect progressive | 105 | 76 | 75 | 72 | 78 | 63 | 78 | 73 | 90 | 70 | 54 | 64 | 75 | 66 | 63 | 87 | 62 | 70 | 60 | 58 | 58 | 66 | 20 | 63 | 69 | 69 | 50 | 67 |
| Reflexive - present perfect simple | 127 | 60 | 72 | 68 | 74 | 81 | 68 | 59 | 61 | 64 | 61 | 59 | 76 | 66 | 80 | 73 | 58 | 65 | 63 | 56 | 51 | 49 | 37 | 63 | 53 | 53 | 51 | 61 |
| Reflexive - present progressive | 586 | 82 | 87 | 82 | 70 | 81 | 74 | 76 | 76 | 64 | 51 | 59 | 65 | 63 | 63 | 89 | 74 | 77 | 67 | 59 | 64 | 62 | 42 | 53 | 31 | 31 | 48 | 65 |
| Reflexive - simple past | 256 | 92 | 96 | 90 | 89 | 95 | 78 | 95 | 91 | 81 | 75 | 69 | 81 | 75 | 75 | 93 | 74 | 84 | 72 | 68 | 63 | 61 | 34 | 50 | 61 | 61 | 44 | 75 |
| Reflexive - simple present | 330 | 78 | 90 | 80 | 75 | 74 | 78 | 72 | 65 | 69 | 40 | 50 | 62 | 59 | 61 | 97 | 65 | 78 | 60 | 60 | 67 | 67 | 36 | 55 | 31 | 31 | 57 | 63 |
| Transitive - future II progressive | 21 | 95 | 90 | 81 | 71 | 86 | 71 | 81 | 100 | 71 | 100 | 71 | 89 | 71 | 72 | 76 | 67 | 71 | 86 | 86 | 71 | 67 | 90 | 83 | 71 | 71 | 67 | 79 |
| Transitive - conditional I progressive | 18 | 94 | 100 | 100 | 83 | 100 | 94 | 83 | 94 | 67 | 78 | 61 | 89 | 67 | 72 | 72 | 67 | 56 | 94 | 94 | 44 | 44 | 44 | 83 | 78 | 78 | 94 | 78 |
| Transitive - conditional I simple | 25 | 100 | 100 | 100 | 92 | 100 | 100 | 100 | 100 | 76 | 68 | 68 | 92 | 84 | 82 | 76 | 80 | 76 | 92 | 100 | 72 | 80 | 32 | 80 | 64 | 64 | 84 | 83 |
| Transitive - conditional II progressive | 51 | 90 | 98 | 100 | 80 | 98 | 92 | 80 | 90 | 70 | 76 | 70 | 76 | 88 | 82 | 75 | 82 | 85 | 75 | 70 | 73 | 67 | 41 | 63 | 47 | 47 | 61 | 76 |
| Transitive - conditional II simple | 20 | 100 | 100 | 100 | 100 | 88 | 100 | 100 | 100 | 70 | 80 | 70 | 85 | 85 | 90 | 100 | 80 | 85 | 69 | 70 | 65 | 60 | 10 | 75 | 50 | 50 | 85 | 80 |
| Transitive - future I progressive | 35 | 86 | 66 | 77 | 63 | 89 | 83 | 63 | 89 | 63 | 69 | 49 | 63 | 63 | 60 | 49 | 54 | 29 | 69 | 80 | 43 | 49 | 71 | 69 | 43 | 43 | 71 | 63 |
| Transitive - future I simple | 53 | 81 | 81 | 75 | 72 | 100 | 77 | 75 | 53 | 55 | 75 | 57 | 87 | 79 | 79 | 32 | 74 | 45 | 85 | 96 | 32 | 47 | 58 | 55 | 81 | 81 | 57 | 69 |
| Transitive - future II simple | 201 | 65 | 80 | 73 | 58 | 100 | 78 | 72 | 64 | 60 | 75 | 57 | 85 | 95 | 91 | 29 | 79 | 76 | 72 | 95 | 89 | 47 | 36 | 53 | 71 | 71 | 67 | 72 |
| Transitive - past perfect progressive | 18 | 83 | 67 | 61 | 67 | 100 | 44 | 72 | 22 | 44 | 78 | 72 | 83 | 78 | 83 | 28 | 78 | 61 | 72 | 72 | 44 | 56 | 44 | 83 | 78 | 78 | 47 | 72 |
| Transitive - past perfect simple | 47 | 55 | 79 | 64 | 64 | 62 | 64 | 57 | 53 | 51 | 83 | 70 | 72 | 55 | 72 | 77 | 49 | 55 | 57 | 62 | 32 | 34 | 66 | 72 | 21 | 21 | 81 | 59 |
| Transitive - past progressive | 14 | 43 | 57 | 57 | 43 | 86 | 79 | 29 | 71 | 43 | 21 | 43 | 64 | 50 | 64 | 14 | 36 | 36 | 43 | 36 | 14 | 14 | 7 | 57 | 36 | 36 | 57 | 44 |
| Transitive - present perfect progressive | 23 | 83 | 61 | 70 | 61 | 96 | 87 | 91 | 87 | 57 | 74 | 70 | 61 | 35 | 52 | 65 | 35 | 35 | 57 | 70 | 39 | 43 | 57 | 48 | 48 | 48 | 52 | 61 |
| Transitive - present perfect simple | 23 | 87 | 78 | 78 | 74 | 88 | 91 | 96 | 100 | 43 | 74 | 70 | 43 | 43 | 70 | 30 | 39 | 35 | 60 | 70 | 43 | 35 | 52 | 75 | 43 | 43 | 74 | 64 |
| Transitive - present progressive | 25 | 88 | 72 | 64 | 64 | 88 | 76 | 64 | 40 | 52 | 60 | 56 | 60 | 56 | 60 | 28 | 56 | 36 | 60 | 64 | 52 | 60 | 60 | 60 | 44 | 44 | 68 | 59 |
| Transitive - simple past | 53 | 96 | 75 | 77 | 66 | 85 | 74 | 85 | 92 | 81 | 62 | 57 | 45 | 47 | 49 | 45 | 43 | 26 | 43 | 47 | 34 | 51 | 57 | 49 | 72 | 72 | 58 | 61 |

525

Table 4: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-German

| ling. category | ling. phenomenon | # | XCOMET | MetricX-24 | MetricX-24-Hybrid | metametrics | MetricX-24-QE | MetricX-24-Hybrid-QE | XCOMET-QE | CometKiwi-XXL | BLEURT-20 | CometKiwi | COMET-22 | chrS | MEE4 | chrF | gemba | BERTScore | YiSi-1 | spBLEU | BLEU | mommonti | damonmonti | sentinel-cand-mqm | PrismRefSmall | XLsimDA | XLsimMqm | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Verb valency | Transitive - simple present | 35 | 80 | 71 | 71 | 63 | 74 | 74 | 80 | 80 | 43 | 94 | 51 | 37 | 37 | 46 | 63 | 29 | 34 | 43 | 49 | 49 | 57 | 83 | 71 | 66 | 66 | 63 | 61 |
| | Case government | 189 | 81 | 76 | 74 | 78 | 87 | 78 | 89 | 78 | 78 | 77 | 79 | 69 | 68 | 62 | 81 | 70 | 77 | 67 | 62 | 80 | 73 | 70 | 51 | 51 | 51 | 52 | 72 |
| | Catenative verb | 885 | 89 | 88 | 90 | 81 | 87 | 79 | 86 | 92 | 81 | 70 | 74 | 65 | 64 | 67 | 76 | 63 | 68 | 60 | 50 | 65 | 56 | 60 | 56 | 73 | 73 | 53 | 72 |
| | Mediopassive voice | 183 | 95 | 95 | 99 | 99 | 94 | 96 | 91 | 92 | 97 | 86 | 98 | 96 | 95 | 98 | 87 | 93 | 95 | 93 | 89 | 79 | 80 | 89 | 82 | 63 | 63 | 75 | 89 |
| | Passive voice | 176 | 77 | 84 | 81 | 84 | 83 | 82 | 65 | 62 | 67 | 83 | 78 | 76 | 47 | 76 | 61 | 67 | 74 | 69 | 44 | 75 | 57 | 85 | 49 | 47 | 58 | 55 | 64 |
| | Resultative | 1203 | 83 | 88 | 88 | 85 | 76 | 85 | 82 | 86 | 82 | 49 | 80 | 74 | 68 | 73 | 76 | 76 | 77 | 73 | 70 | 80 | 79 | 81 | 58 | 58 | 58 | 58 | 76 |
| | Semantic roles | 670 | 65 | 77 | 55 | 87 | 73 | 72 | 71 | 68 | 62 | 64 | 61 | 53 | 57 | 55 | 34 | 58 | 58 | 55 | 53 | 56 | 50 | 54 | 59 | 55 | 55 | 41 | 69 |
| Verb semantics | Verb semantics | 180 | 87 | 73 | 69 | 82 | 73 | 71 | 78 | 77 | 62 | 73 | 73 | 71 | 71 | 55 | 69 | 58 | 68 | 55 | 65 | 56 | 64 | 63 | 59 | 69 | 69 | 54 | 64 |
| macro avg. | | 78727 | 82 | 82 | 81 | 79 | 73 | 78 | 77 | 77 | 62 | 73 | 73 | 71 | 71 | 70 | 69 | 69 | 68 | 68 | 65 | 65 | 64 | 63 | 59 | 57 | 57 | 57 | 70 |
| micro avg. | | 78727 | 78 | 79 | 79 | 79 | 73 | 74 | 72 | 73 | 77 | 66 | 76 | 71 | 71 | 70 | 65 | 69 | 73 | 68 | 64 | 69 | 69 | 67 | 54 | 54 | 54 | 52 | 69 |

Table 5: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-Russian

| ling. category | ling. phenomenon | # | MetricX-24 | metametrics | XCOMET | CometKiwi-XXL | MetricX-24-Hybrid | MetricX-24-QE | COMET-22 | BLEURT-20 | MetricX-24-Hybrid-QE | XCOMET-QE | sentinel-cand-mqm | CometKiwi | YiSi-1 | BERTScore | chrS | spBLEU | chrF | BLEU | mommonti | damonmonti | gemba | XLsimDA | XLsimMqm | PrismRefSmall | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambiguity | Lexical ambiguity | 3788 | 97 | 96 | 93 | 81 | 95 | 96 | 87 | 90 | 91 | 84 | 87 | 77 | 83 | 75 | 75 | 71 | 73 | 69 | 78 | 81 | 89 | 44 | 44 | 48 | 45 | 78 |
| Coordination \& ellipsis | Gapping | 698 | 93 | 92 | 93 | 87 | 92 | 91 | 93 | 88 | 89 | 90 | 85 | 96 | 86 | 85 | 81 | 81 | 77 | 79 | 81 | 74 | 78 | 57 | 57 | 55 | 56 | 81 |
| | Pseudogapping | 381 | 71 | 67 | 55 | 64 | 70 | 67 | 62 | 63 | 70 | 54 | 54 | 57 | 60 | 59 | 57 | 54 | 57 | 55 | 55 | 53 | 53 | 39 | 39 | 40 | 40 | 57 |
| | Right node raising | 183 | 78 | 74 | 74 | 77 | 77 | 70 | 75 | 63 | 70 | 67 | 79 | 75 | 62 | 68 | 64 | 60 | 61 | 58 | 70 | 63 | 55 | 45 | 45 | 66 | 61 | 67 |
| | Sluicing | 384 | 80 | 82 | 86 | 89 | 82 | 70 | 75 | 73 | 82 | 85 | 84 | 69 | 61 | 63 | 56 | 56 | 54 | 60 | 68 | 57 | 35 | 48 | 48 | 66 | 49 | 67 |
| | Stripping | 375 | 70 | 69 | 80 | 76 | 74 | 70 | 63 | 67 | 81 | 75 | 69 | 69 | 62 | 60 | 56 | 53 | 54 | 49 | 45 | 55 | 70 | 57 | 57 | 49 | 50 | 64 |
| | VP-ellipsis | 252 | 80 | 77 | 80 | 88 | 85 | 86 | 75 | 75 | 90 | 69 | 76 | 83 | 62 | 56 | 48 | 53 | 47 | 49 | 45 | 55 | 34 | 65 | 65 | 44 | 41 | 66 |
| False friends | False friends | 2414 | 88 | 84 | 76 | 83 | 85 | 69 | 83 | 76 | 68 | 69 | 58 | 52 | 88 | 76 | 85 | 66 | 83 | 76 | 62 | 81 | 34 | 43 | 43 | 53 | 42 | 69 |
| Function word | Focus particle | 846 | 70 | 62 | 63 | 68 | 60 | 67 | 67 | 57 | 49 | 66 | 58 | 55 | 45 | 63 | 63 | 65 | 60 | 77 | 57 | 50 | 63 | 71 | 75 | 53 | 50 | 61 |
| | Question tag | 1587 | 89 | 87 | 95 | 91 | 81 | 95 | 89 | 92 | 92 | 97 | 93 | 92 | 66 | 73 | 65 | 65 | 61 | 77 | 57 | 57 | 63 | 71 | 75 | 53 | 48 | 77 |
| LDD \& interrogatives | Inversion | 333 | 82 | 88 | 84 | 73 | 90 | 89 | 79 | 83 | 71 | 78 | 85 | 67 | 71 | 68 | 68 | 66 | 68 | 67 | 61 | 60 | 67 | 47 | 47 | 56 | 52 | 71 |
| | Modifying Comparison | 90 | 68 | 71 | 74 | 87 | 69 | 78 | 52 | 100 | 71 | 74 | 98 | 56 | 44 | 41 | 33 | 29 | 29 | 28 | 56 | 67 | 56 | 73 | 73 | 37 | 40 | 61 |
| | Multiple connectors | 400 | 97 | 92 | 93 | 95 | 98 | 92 | 88 | 78 | 96 | 96 | 95 | 94 | 64 | 67 | 62 | 58 | 61 | 55 | 60 | 52 | 86 | 52 | 52 | 69 | 52 | 76 |

Continued on next page

526

Table 5: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-Russian

| ling. category | ling. phenomenon | # | MetricX-24 | metametrics | XCOMET | CometKiwi-XXL | MetricX-24-Hybrid | MetricX-24-QE | COMET-22 | BLEURT-20 | MetricX-24-Hybrid-QE | XCOMET-QE | sentinel-cand-mqm | CometKiwi | Yisi-1 | BERTScore | chrF++ | spBLEU | chrF | BLEU | mqmonli | damonmonli | gemba | XLsimDA | XLsimMqm | PrismRefSmall | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Pied-piping | 343 | 80 | 80 | 78 | 82 | 81 | 76 | 83 | 72 | 80 | 79 | 81 | 75 | 63 | 66 | 55 | 61 | 52 | 61 | 50 | 50 | 80 | 68 | 68 | 35 | 33 | 68 |
| | Preposition stranding | 393 | 90 | 90 | 89 | 90 | 92 | 93 | 88 | 86 | 90 | 90 | 92 | 84 | 75 | 75 | 70 | 72 | 67 | 69 | 64 | 61 | 66 | 58 | 58 | 48 | 46 | 76 |
| | Topicalization | 207 | 71 | 74 | 77 | 76 | 77 | 86 | 70 | 74 | 80 | 75 | 79 | 58 | 63 | 61 | 57 | 61 | 55 | 64 | 71 | 57 | 48 | 42 | 42 | 50 | 49 | 65 |
| | Wh-movement | 173 | 92 | 93 | 86 | 88 | 87 | 88 | 93 | 95 | 88 | 76 | 87 | 85 | 58 | 62 | 71 | 62 | 65 | 64 | 39 | 41 | 42 | 42 | 42 | 56 | 51 | 70 |
| MWE | Collocation | 2167 | 73 | 77 | 79 | 80 | 74 | 76 | 79 | 85 | 75 | 76 | 67 | 85 | 82 | 76 | 80 | 76 | 73 | 73 | 71 | 68 | 88 | 59 | 59 | 56 | 52 | 71 |
| | Compound | 1393 | 88 | 91 | 92 | 83 | 87 | 83 | 89 | 95 | 93 | 86 | 72 | 90 | 78 | 78 | 75 | 76 | 72 | 69 | 74 | 77 | 85 | 62 | 62 | 52 | 52 | 79 |
| | Idiom | 1784 | 100 | 98 | 95 | 93 | 100 | 99 | 95 | 95 | 99 | 91 | 90 | 99 | 88 | 79 | 78 | 73 | 74 | 73 | 57 | 64 | 57 | 67 | 67 | 51 | 49 | 83 |
| | Nominal MWE | 2166 | 75 | 72 | 68 | 67 | 77 | 68 | 71 | 73 | 74 | 56 | 52 | 53 | 68 | 68 | 75 | 70 | 74 | 67 | 57 | 57 | 57 | 43 | 43 | 48 | 51 | 64 |
| | Prepositional MWE | 1639 | 88 | 87 | 80 | 92 | 87 | 90 | 82 | 73 | 91 | 83 | 81 | 81 | 75 | 75 | 71 | 52 | 72 | 74 | 71 | 52 | 82 | 46 | 46 | 54 | 52 | 75 |
| | Verbal MWE | 453 | 68 | 63 | 72 | 79 | 60 | 64 | 69 | 68 | 63 | 70 | 65 | 60 | 69 | 69 | 67 | 61 | 64 | 64 | 75 | 66 | 59 | 29 | 29 | 41 | 48 | 62 |
| | Date | 3403 | 87 | 81 | 74 | 69 | 86 | 82 | 77 | 83 | 71 | 64 | 65 | 69 | 78 | 70 | 72 | 71 | 71 | 68 | 60 | 77 | 54 | 38 | 38 | 55 | 50 | 68 |
| | Domainspecific Term | 3471 | 90 | 97 | 89 | 72 | 89 | 62 | 95 | 95 | 69 | 64 | 84 | 70 | 88 | 82 | 86 | 82 | 83 | 74 | 71 | 72 | 70 | 55 | 55 | 53 | 50 | 76 |
| Named entity \& terminology | Measuring unit | 3510 | 63 | 72 | 85 | 57 | 59 | 52 | 77 | 73 | 54 | 64 | 53 | 58 | 63 | 81 | 81 | 82 | 79 | 81 | 54 | 56 | 64 | 45 | 45 | 57 | 56 | 64 |
| | Onomatopeia | 3401 | 86 | 86 | 85 | 84 | 86 | 84 | 86 | 87 | 84 | 80 | 73 | 78 | 85 | 78 | 81 | 79 | 80 | 62 | 78 | 80 | 74 | 57 | 57 | 47 | 48 | 77 |
| | Proper Name \& Location | 2160 | 93 | 90 | 90 | 85 | 92 | 88 | 85 | 90 | 89 | 81 | 82 | 86 | 90 | 82 | 81 | 76 | 80 | 62 | 81 | 87 | 64 | 44 | 44 | 56 | 58 | 78 |
| | Proper name | 339 | 78 | 92 | 93 | 83 | 83 | 65 | 96 | 80 | 66 | 63 | 83 | 91 | 86 | 94 | 79 | 80 | 70 | 54 | 56 | 60 | 7 | 24 | 24 | 67 | 58 | 69 |
| Negation | Negation | 346 | 65 | 60 | 49 | 59 | 59 | 58 | 67 | 72 | 45 | 45 | 41 | 49 | 79 | 74 | 83 | 74 | 84 | 72 | 73 | 71 | 42 | 50 | 50 | 49 | 46 | 61 |
| | Coreference | 526 | 86 | 82 | 83 | 81 | 84 | 83 | 74 | 69 | 83 | 81 | 80 | 77 | 63 | 61 | 57 | 58 | 57 | 66 | 57 | 60 | 58 | 58 | 58 | 44 | 35 | 68 |
| | Genitive | 2068 | 82 | 73 | 71 | 61 | 72 | 68 | 68 | 74 | 61 | 67 | 60 | 64 | 72 | 67 | 72 | 72 | 72 | 69 | 44 | 43 | 56 | 78 | 78 | 49 | 47 | 66 |
| Non-verbal agreement | Lexical Morphology/Functional shift | 1134 | 97 | 95 | 98 | 95 | 96 | 97 | 95 | 94 | 95 | 97 | 89 | 92 | 84 | 84 | 81 | 77 | 77 | 74 | 72 | 78 | 81 | 85 | 85 | 68 | 61 | 86 |
| | Lexical Morphology/Noun formation (er) | 670 | 97 | 96 | 95 | 96 | 98 | 93 | 98 | 96 | 96 | 94 | 90 | 87 | 90 | 90 | 86 | 77 | 82 | 63 | 82 | 80 | 83 | 66 | 66 | 50 | 48 | 84 |
| | Personal Pronoun Coreference | 1290 | 86 | 79 | 85 | 84 | 91 | 93 | 72 | 87 | 93 | 89 | 68 | 86 | 63 | 54 | 54 | 54 | 54 | 54 | 56 | 74 | 62 | 62 | 62 | 57 | 43 | 70 |
| | Possessive Pronouns | 521 | 80 | 79 | 77 | 78 | 80 | 77 | 70 | 75 | 85 | 81 | 64 | 68 | 69 | 64 | 64 | 61 | 64 | 63 | 50 | 50 | 54 | 45 | 45 | 51 | 45 | 66 |
| | Substitution | 546 | 77 | 75 | 74 | 76 | 78 | 80 | 76 | 74 | 81 | 78 | 82 | 76 | 61 | 65 | 67 | 62 | 67 | 63 | 73 | 56 | 66 | 48 | 48 | 46 | 44 | 67 |
| Punctuation | Quotation marks | 363 | 100 | 71 | 71 | 100 | 73 | 64 | 76 | 69 | 67 | 71 | 71 | 72 | 76 | 73 | 62 | 65 | 58 | 61 | 64 | 51 | 61 | 57 | 57 | 46 | 45 | 65 |
| | Adverbial clause | 1458 | 95 | 87 | 99 | 100 | 93 | 83 | 74 | 64 | 100 | 75 | 61 | 63 | 67 | 62 | 63 | 63 | 62 | 62 | 70 | 48 | 52 | 39 | 39 | 44 | 45 | 65 |
| | Cleft sentence | 323 | 71 | 82 | 67 | 79 | 72 | 72 | 63 | 68 | 57 | 65 | 63 | 47 | 59 | 62 | 56 | 60 | 55 | 60 | 62 | 77 | 36 | 38 | 38 | 45 | 41 | 60 |
| | Complex object | 229 | 74 | 70 | 77 | 79 | 71 | 76 | 79 | 90 | 73 | 77 | 65 | 71 | 72 | 85 | 71 | 76 | 72 | 72 | 87 | 70 | 45 | 55 | 55 | 46 | 40 | 70 |
| | Contact clause | 291 | 65 | 53 | 57 | 76 | 65 | 65 | 57 | 71 | 70 | 58 | 52 | 46 | 64 | 66 | 63 | 58 | 59 | 54 | 53 | 52 | 40 | 38 | 38 | 62 | 59 | 58 |
| Subordination | Indirect speech | 46 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 99 | 82 | 100 | 96 | 85 | 67 | 65 | 65 | 78 | 70 | 26 | 100 | 100 | 100 | 28 | 30 | 85 |
| | Infinitive clause | 305 | 95 | 87 | 99 | 95 | 93 | 75 | 89 | 85 | 88 | 99 | 95 | 91 | 71 | 69 | 67 | 59 | 64 | 62 | 70 | 67 | 60 | 33 | 33 | 48 | 47 | 74 |
| | Object clause | 276 | 71 | 82 | 93 | 100 | 72 | 57 | 79 | 68 | 57 | 76 | 64 | 97 | 61 | 70 | 61 | 68 | 59 | 64 | 59 | 51 | 87 | 87 | 87 | 35 | 29 | 70 |
| | Participle clause | 1345 | 77 | 70 | 67 | 73 | 71 | 76 | 81 | 75 | 73 | 62 | 75 | 67 | 80 | 69 | 60 | 68 | 70 | 65 | 62 | 61 | 48 | 44 | 44 | 53 | 54 | 66 |
| | Pseudo-cleft sentence | 369 | 83 | 83 | 76 | 72 | 72 | 75 | 85 | 75 | 70 | 77 | 72 | 85 | 67 | 67 | 60 | 69 | 58 | 72 | 75 | 68 | 52 | 57 | 57 | 40 | 38 | 69 |
| | Relative clause | 1088 | 62 | 76 | 77 | 68 | 57 | 61 | 74 | 68 | 60 | 70 | 60 | 50 | 71 | 76 | 68 | 69 | 68 | 65 | 73 | 63 | 51 | 67 | 67 | 46 | 43 | 65 |
| | Subject clause | 895 | 87 | 89 | 87 | 88 | 94 | 93 | 82 | 82 | 95 | 94 | 91 | 92 | 66 | 55 | 53 | 56 | 48 | 63 | 73 | 39 | 62 | 56 | 56 | 55 | 50 | 72 |
| Verb semantics | Verb semantics | 275 | 88 | 82 | 85 | 80 | 88 | 75 | 74 | 87 | 76 | 80 | 55 | 70 | 53 | 56 | 61 | 54 | 65 | 49 | 67 | 68 | 72 | 33 | 33 | 65 | 67 | 67 |
| | Conditional | 343 | 78 | 72 | 89 | 80 | 70 | 81 | 69 | 76 | 70 | 76 | 61 | 76 | 56 | 55 | 63 | 63 | 61 | 51 | 70 | 65 | 52 | 36 | 36 | 52 | 56 | 65 |
| Verb tense/aspect/mood | Ditransitive | 299 | 92 | 90 | 93 | 97 | 91 | 93 | 91 | 91 | 94 | 96 | 90 | 100 | 63 | 73 | 73 | 67 | 54 | 56 | 56 | 46 | 77 | 63 | 63 | 38 | 47 | 76 |
| | Gerund | 644 | 84 | 85 | 85 | 79 | 85 | 71 | 74 | 68 | 72 | 81 | 57 | 81 | 67 | 66 | 70 | 70 | 69 | 63 | 57 | 60 | 62 | 41 | 41 | 51 | 52 | 68 |

Table 5: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-Russian

| ling. category | ling. phenomenon | \# | MetricX-24 | metametrics | XCOMET | CometKiwi-XXL | MetricX-24-Hybrid | MetricX-24-QE | COMET-22 | BLEURT-20 | MetricX-24-Hybrid-QE | XCOMET-QE | sentinel-cand-mqm | CometKiwi | Yisi-1 | BERTScore | chrfS | spBLEU | chrF | BLEU | mommonit | daemmonit | gemba | XLsimDA | XLsimMqm | PrismRefSmall | PrismRefMedium | avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Imperative | 575 | 88 | 89 | 85 | 88 | 84 | 87 | 83 | 79 | 84 | 85 | 80 | 68 | 74 | 74 | 66 | 69 | 64 | 63 | 69 | 69 | 50 | 44 | 44 | 42 | 44 | 71 |
| | Intransitive | 103 | 94 | 91 | 87 | 98 | 94 | 90 | 94 | 81 | 94 | 88 | 93 | 93 | 81 | 68 | 68 | 65 | 65 | 54 | 58 | 60 | 56 | 42 | 42 | 37 | 41 | 73 |
| | Reflexive | 514 | 88 | 89 | 84 | 94 | 77 | 90 | 84 | 77 | 77 | 83 | 85 | 67 | 70 | 69 | 70 | 70 | 68 | 65 | 41 | 55 | 88 | 58 | 58 | 51 | 51 | 72 |
| | Transitive | 516 | 78 | 87 | 85 | 85 | 85 | 80 | 77 | 68 | 80 | 66 | 47 | 51 | 77 | 74 | 73 | 75 | 78 | 69 | 50 | 61 | 28 | 50 | 50 | 63 | 59 | 68 |
| Verb valency | Case government | 331 | 76 | 86 | 78 | 71 | 76 | 74 | 77 | 85 | 75 | 71 | 51 | 70 | 75 | 82 | 83 | 61 | 82 | 71 | 81 | 79 | 78 | 74 | 74 | 43 | 59 | 73 |
| | Catenative verb | 358 | 72 | 73 | 69 | 68 | 73 | 71 | 66 | 70 | 72 | 70 | 67 | 63 | 58 | 63 | 66 | 53 | 62 | 59 | 62 | 63 | 65 | 49 | 49 | 50 | 48 | 64 |
| | Impersonal Subject | 217 | 86 | 77 | 82 | 95 | 85 | 98 | 75 | 71 | 90 | 94 | 74 | 87 | 65 | 67 | 61 | 70 | 60 | 67 | 76 | 75 | 59 | 50 | 50 | 43 | 41 | 71 |
| | Mediopassive voice | 409 | 77 | 79 | 89 | 85 | 69 | 89 | 80 | 72 | 77 | 90 | 83 | 82 | 73 | 75 | 65 | 75 | 63 | 67 | 57 | 63 | 64 | 58 | 58 | 38 | 37 | 70 |
| | Passive voice | 228 | 94 | 89 | 87 | 84 | 92 | 98 | 88 | 84 | 90 | 84 | 69 | 82 | 79 | 83 | 75 | 75 | 74 | 84 | 73 | 83 | 66 | 74 | 74 | 52 | 44 | 79 |
| | Resultative | 660 | 91 | 86 | 91 | 87 | 91 | 88 | 78 | 80 | 88 | 87 | 82 | 66 | 72 | 75 | 79 | 70 | 76 | 66 | 65 | 67 | 73 | 62 | 62 | 64 | 55 | 76 |
| | Semantic roles | 270 | 91 | 82 | 83 | 77 | 85 | 73 | 79 | 89 | 87 | 79 | 60 | 80 | 71 | 67 | 72 | 68 | 81 | 64 | 80 | 84 | 53 | 63 | 63 | 57 | 50 | 74 |
| | Verb semantics/Verb semantics | 549 | 81 | 81 | 81 | 81 | 74 | 80 | 71 | 70 | 78 | 78 | 74 | 75 | 66 | 70 | 68 | 67 | 66 | 64 | 43 | 46 | 66 | 58 | 58 | 53 | 48 | 68 |
| macro avg. | | 59113 | 82 | 81 | 81 | 81 | 81 | 81 | 79 | 79 | 79 | 78 | 75 | 75 | 71 | 70 | 68 | 67 | 66 | 64 | 64 | 63 | 62 | 54 | 54 | 50 | 48 | 70 |
| micro avg. | | 59113 | 83 | 83 | 81 | 78 | 82 | 79 | 81 | 81 | 77 | 76 | 73 | 73 | 76 | 73 | 73 | 71 | 71 | 68 | 65 | 67 | 64 | 53 | 53 | 52 | 49 | 71 |

528

# TMU-HIT's Submission for the WMT24 Quality Estimation Shared Task: Is GPT-4 a Good Evaluator for Machine Translation?

**Ayako Sato[†], Kyotaro Nakajima[†], Hwichan Kim[†]**
**Zhousi Chen[‡], Mamoru Komachi[‡]**
[†]Tokyo Metropolitan University, [‡]Hitotsubashi University
{sato-ayako, nakajima-kyotaro, kim-hwichan}@ed.tmu.ac.jp
{zhousi.chen, mamoru.komachi}@hit-u.ac.jp

## Abstract

In machine translation quality estimation (QE), translation quality is evaluated automatically without the need for reference translations. This paper describes our contribution to the sentence-level subtask of Task 1 at the Ninth Machine Translation Conference (WMT24), which predicts quality scores for neural MT outputs without reference translations. We fine-tune GPT-4o mini, a large-scale language model (LLM), with limited data for QE. We report results for the direct assessment (DA) method for four language pairs: English-Gujarati (En-Gu), English-Hindi (En-Hi), English-Tamil (En-Ta), and English-Telugu (En-Te). Experiments under zero-shot, few-shot prompting, and fine-tuning settings revealed significantly low performance in the zero-shot, while fine-tuning achieved accuracy comparable to last year's best scores. Our system demonstrated the effectiveness of this approach in low-resource language QE, securing 1st place in both En-Gu and En-Hi, and 4th place in En-Ta and En-Te. The code used in our experiments is available at the following URL [1].

## 1 Introduction

Machine translation quality estimation evaluates translation automatically without reference translation. This practice reduces the cost of manual translation and enables efficient evaluation. The subsequent quality score flags the necessity of resorting to a more reliable translation system or revision from human post-editing. Quality estimation can be performed at various granularity levels, including word, phrase, sentence, and document.

In this paper, we describe our contribution to the QE shared task at the Ninth Machine Translation Conference (WMT24). We participate in the Task 1 of the shared task and we specifically focus on the sentence-level subtask, which involves

|          | En-Gu | En-Hi | En-Ta | En-Te |
|----------|-------|-------|-------|-------|
| baseline | 0.661 | 0.678 | 0.592 | 0.414 |
| gpt4o-mean | **0.712** | **0.735** | **0.616** | 0.457 |
| gpt4o-prob | **0.712** | 0.734 | 0.608 | **0.460** |

Table 1: Spearman's rank correlation coefficient between our predictions and human DA judgments of WMT24 test data. The best score obtained for each language pair is marked in bold.

predicting the quality score of neural MT outputs at the sentence level without access to reference translations. There are two different annotation methods for QE: Multidimensional Quality Metric (MQM) (Freitag et al., 2021) and Direct Assessment (DA) (Fomicheva et al., 2022), and we report the results of DA score prediction. Our study targets four language pairs: English-Gujarati (En-Gu), English-Hindi (En-Hi), English-Tamil (En-Ta), and English-Telugu (En-Te). The participating systems are assigned the task of predicting the quality score of each source and target sentence pair, and their performance is evaluated using Spearman's rank correlation coefficient as the primary metric, and Pearson and Kendall coefficients as supplementary metrics.

We present a system for quality estimation utilizing a large language model (LLM), inspired by the success of LLMs in regression tasks (Liu et al., 2023; Enomoto et al., 2024). Specifically, we manually designed a prompt for quality estimation and employed GPT-4o mini (OpenAI, 2024) to generate assessment scores multiple times based on this prompt. We then used either the averaged score of these generated scores or their weighted sum based on the generation probability as the final score. Evaluation experiments were conducted in both zero-shot and three-shot settings. Noticeably, we fine-tuned GPT-4o mini using the training data released at WMT23 (Kocmi et al., 2023) and as-

---

[1] https://colab.research.google.com/drive/1p8VMnAkRfuVpbvM_revV2ZaN76sSxmiE?usp=sharing

sessed its performance.

We first evaluated our systems using the development data released at WMT23. The results indicated that the Spearman's correlations in the zero- and few-shot settings ranged from 0.2 to 0.4, while those for the fine-tuned GPT-4o mini ranged from 0.4 to 0.7. Compared to a single generation, the estimated score derived from the average or weighted sum based on multiple output values was found to perform consistently better. Subsequently, we evaluated the system based on the fine-tuned GPT-4o mini using the test data from WMT24. Table 1 presents the results of our system and the baseline for Task 1 (Rei et al., 2022). The system achieved Spearman's correlation scores of 0.712, 0.735, 0.616, and 0.460 in the En-Gu, En-Hi, En-Ta, and En-Te language pairs, respectively, surpassing the baseline system's performance. We achieved the 1st place in En-Gu and En-Hi, and 4th place in En-Ta and En-Te.

## 2 Related Work

GEMBA (Kocmi and Federmann, 2023) is a translation quality metric that utilizes a large language model (LLM). It has been shown to have a high correlation with the human-rated MQM score of the WMT22 Metrics shared task. Their experiments covered three language pairs (English to German, English to Russian, and Chinese to English) of the WMT22 Metrics shared task using seven GPT variants from GPT-2 to GPT-4 models. Lu et al. (2024) investigated various prompts to improve segment-level evaluation performance. They experimented with Llama2-70B model (Touvron et al., 2023) and Mixtral-8x7b model (Jiang et al., 2024) in addition to GPT-3.5-Turbo model, and showed that the method using GPT-3.5-Turbo had the best performance. These previous studies highlight the potential of LLM evaluators as human alternatives. Our system uses the latest model, GPT-4o mini, and also estimates quality for more challenging translations for low-resource languages.

Enomoto et al. (2024) used an LLM to solve the lexical complexity prediction task. They reported a bias in the numerical values generated by an LLM in that certain values occur frequently regardless of the input. To mitigate the bias and achieve a more precise numerical output, we run the generation several times and obtain the final scores by either average or expectation weighted by generation probabilities.

## 3 Methodology

Our system uses LLMs to estimate translation quality scores. Following GEMBA (Kocmi and Federmann, 2023), to assess translation quality via prompting an LLM, the following arguments are required:

- source language name: {{source language}}
- target language name: {{target language}}
- source sentences: {{$src_1, ..., src_N$}}
- translated sentences: {{$hyp_1, ..., hyp_N$}}
- few-shot examples: {{examples}} (optional)

We define the instructions to be input into the LLM as follows:

> *Please analyze the given source and translated sentences and output a translation quality score on a continuous scale ranging from 0 to 100.*
> *Translation quality should be evaluated based on both fluency and adequacy.*
> *A score close to 0 indicates a low quality translation, while a score close to 100 indicates a high quality translation.*
> *Do not provide any explanations or text apart from the score.*
>
> {{examples}}
> {{source language}} *Sentence:* {{$src_i$}}
> {{target language}} *Sentence:* {{$hyp_i$}}
> *Score:*

The instruction template is designed to include a description of the task, the score range, and a description of the evaluation criteria. To restrict the output to numerical values only, it is important to state "Do not provide any explanations or text apart from the score." explicitly.

According to Kocmi and Federmann (2023), there are some numbers that are particularly prone to output, such as "95". To mitigate such bias in the output distribution, the final score is computed from the sampled generated results with reference to the G-Eval framework (Liu et al., 2023). $Score_{mean}$ is the simple average of the generated scores, while $Score_{prob}$ is the score weighted by the generation probabilities. Let $S = \{s_1, s_2, ..., s_n\}$ represent the set of scores generated by the prompt, and let $p(s_i)$ be the softmax output probability of each generated score.

| Method | Setting | En-Gu | | | En-Hi | | | En-Ta | | | En-Te | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $\rho$ | $r$ | $\tau$ | $\rho$ | $r$ | $\tau$ | $\rho$ | $r$ | $\tau$ | $\rho$ | $r$ | $\tau$ |
| | | | | | Single generation | | | | | | | | |
| | Zero-shot | 0.205 | 0.501 | 0.337 | 0.379 | 0.434 | 0.291 | 0.277 | 0.514 | 0.514 | 0.271 | 0.280 | 0.208 |
| | Three-shot | 0.413 | 0.491 | 0.309 | 0.390 | 0.400 | 0.282 | 0.463 | 0.420 | 0.344 | 0.294 | 0.305 | 0.226 |
| | Fine-tuned | 0.599 | 0.659 | 0.453 | 0.510 | 0.639 | 0.367 | 0.618 | 0.704 | 0.453 | 0.283 | 0.263 | 0.205 |
| | | | | | Multi generation | | | | | | | | |
| $score_{mean}$ | Zero-shot | 0.453 | 0.505 | 0.328 | 0.389 | 0.449 | 0.277 | 0.514 | 0.529 | 0.373 | 0.274 | 0.275 | 0.193 |
| | Three-shot | 0.447 | 0.512 | 0.319 | 0.422 | 0.426 | 0.294 | 0.498 | 0.428 | 0.358 | 0.290 | 0.303 | 0.205 |
| | Fine-tuned | 0.680 | **0.717** | 0.506 | 0.564 | 0.686 | 0.409 | 0.661 | **0.747** | 0.487 | 0.392 | **0.361** | 0.270 |
| $score_{prob}$ | Zero-shot | 0.451 | 0.499 | 0.323 | 0.394 | 0.447 | 0.275 | 0.519 | 0.521 | 0.368 | 0.274 | 0.276 | 0.190 |
| | Three-shot | 0.448 | 0.514 | 0.319 | 0.423 | 0.427 | 0.295 | 0.500 | 0.426 | 0.358 | 0.290 | 0.303 | 0.202 |
| | Fine-tuned | **0.683** | 0.715 | **0.508** | **0.568** | **0.690** | **0.412** | **0.663** | 0.746 | **0.489** | **0.399** | 0.360 | **0.277** |

Table 2: Spearman ($\rho$), Pearson ($r$) and Kendall ($\tau$) correlation between the proposed approaches and human DA judgments of WMT23 dev data. The best Spearman score obtained for each language pair is marked in bold. Single generation is a prediction method that uses the output value generated only once as the estimated score, and the other two are methods that calculate the average value or the expected value based on the generation probability, based on the output values by 20 times generation.

The final scores are calculated using the following formulae:

$$score_{mean} = \frac{1}{n}\sum_{i=1}^{n} s_i \qquad (1)$$

$$score_{prob} = \sum_{i=1}^{n} p(s_i) \times s_i \qquad (2)$$

In this study, the experiment is conducted with $n = 20$. Among the generated outputs, non-numeric tokens and numbers outside the specified range are excluded from $S$.

## 4 Experiments

We conduct experiments to investigate two RQs. **RQ 1**: Which methods are more effective in improving performance of low-resource language QE? **RQ 2**: How effective are sampling methods in mitigating numerical output bias?

### 4.1 Settings

**Model** We use GPT-4o mini ("gpt-4o-mini-2024-07-18") (OpenAI, 2024) for our experiments. It is priced at 15 cents per million input tokens and 60 cents per million output tokens and more than 60% cheaper than GPT-3.5 Turbo.

Our system fine-tunes GPT-4o mini. The fine-tuning process is conducted using OpenAI's API.

**Data** For the remaining sections, we only report results on the WMT23 dev dataset. The examples used for few-shot prompting are randomly obtained

from the WMT23 training dataset. WMT23 training data consists of 7,000 sentence pairs in each language and is also used for fine-tuning.

### 4.2 Results

Spearman, Pearson, and Kendall correlation coefficients between predicted and gold scores for each language pair are shown in Table 2.

#### 4.2.1 Strategies for Low-Resource Languages

For **RQ 1**, we compare the performance of the three settings: zero-shot, few-shot prompting, and fine-tuning. Few-shot in $score_{mean}$ and $score_{prob}$ improved Spearman correlation coefficients slightly by 0.033 for En-Hi and 0.016 for En-Te, while En-Gu and En-Ta scores decreased by 0.006 and 0.016 respectively. In other words, the few-shot strategy is not very effective for low-resource languages. On the other hand, in the single genaration setting, En-Gu and En-Ta improved by 0.208 and 0.186, respectively, indicating that the few-shot is more effective when the generation times are limited.

Fine-tuning improves performance in almost all evaluation metrics and is an effective measure for low-resource languages. Our systems submitted to the shared task (Table 1) are also the result of the fine-tuned models.

#### 4.2.2 Strategies for Distributional Bias

For **RQ 2**, we compare the performance of the three settings: single generation, $score_{mean}$, and $score_{prob}$. In the fine-tuned model, the difference in Spearman's rank correlation coefficient with sin-

|          | En-Gu  | En-Hi  | En-Ta  | En-Te  |
|----------|--------|--------|--------|--------|
| Manually | **0.451** | **0.394** | **0.519** | **0.274** |
| AutoCoT  | 0.444  | 0.387  | 0.514  | 0.238  |

Table 3: Spearman's rank correlation coefficient between our predictions in a zero-shot setting using two different prompt generation methods and human DA judgments of WMT23 dev data.

gle generation (average of four languages) is 0.072 for $score_{mean}$ and 0.076 for $score_{prob}$. Compared to single genaration, the other two sampling methods performed better, demonstrating the effectiveness of these methods in mitigating the effects of bias during generation. The performance of $score_{mean}$ and $score_{prob}$ is almost equal, and either method can be used.

## 5 Discussion

### 5.1 Is AutoCoT necessary for G-Eval?

In Chiang and Lee (2023), they find that the auto Chain-of-Thought (CoT) used in G-Eval does not always make G-Eval more aligned with human ratings. In this section, we examine the methods used to construct the prompt for this task.

To replicate the G-Eval framework (Liu et al., 2023) procedures, it is necessary to construct an initial prompt to generate the evaluation steps using AutoCoT. Specifically, we first manually designed a prompt that contains the definition of the QE task and the desired evaluation criteria as follows:

*You will be given a source and a translated sentence. Your task is to rate translated sentence on one metric. Please make sure you read and understand these instructions carefully. Please keep this document open while reviewing, and refer to it as needed.*

*Evaluation Criteria:*

*Translation Quality (0 - 100) - the quality of a translation based on the adequacy and fluency of the sentence.*

Then, we added a line of "*Evaluation Steps:*" to the prompt and let GPT-4 [2] generate the following evaluation steps by CoT automatically:

*Evaluation steps:*

*1. Read the source sentence and the translated sentence carefully.*

*2. Evaluate the translated sentence based on its adequacy and fluency.*

*- Adequacy: How much of the meaning expressed in the source text is also expressed in the target text? A score of 100 means all the meaning is transferred, and 0 means none of it is.*

*- Fluency: Does the translation sound like something a native speaker would say? A score of 100 means it sounds completely native, and 0 means it doesn't sound native at all.*

*3. Give the translation a score between 0 and 100, where 0 is the worst and 100 is the best.*

We compare the performance in a zero-shot setting using prompts created by AutoCoT and those created manually. As shown in Table 3, manual prompts performed better than AutoCoT for all languages. This result follows the findings of Chiang and Lee (2023), and we decided to use manually constructed prompts in our systems to get results that correlated better with human judgment.

### 5.2 Is it difficult for GPT-4 evaluators to evaluate Telugu text?

In our results, the performance on Telugu data was lower than other languages. This may be attributed to the linguistic complexity of Telugu, which features complex noun and verb conjugations, as well as its status as a low-resource language. Kishore and Shaik (2024) demonstrated that ChatGPT is less accurate in Telugu grammar and vocabulary compared to Gemini. Performance is expected to improve by using LLMs specialized for each language (e.g., Telugu GPT[3]) rather than relying on a single, generalized model.

## 6 Conclusion

Our study demonstrates the efficacy of using a LLM for sentence-level quality estimation in machine translation. By leveraging GPT-4o mini, we achieved improvements over baseline systems in predicting quality scores for various language pairs. The fine-tuned GPT-4o mini model exhibited robust performance in low-resource language QE,

---

[2] We used `gpt-4-0613` following Liu et al. (2023)

[3] https://chatgpt.com/g/g-RjoqGo7g0-telugu-gpt

with Spearman's correlation scores significantly higher than those in the zero- and few-shot settings. These findings emphasize that fine-tuning with an annotated QE dataset is crucial for enhancing performance in low-resource languages. However, in practical scenarios, creating and obtaining such datasets for low-resource languages poses significant challenges. Therefore, efforts to effectively improve performance using a small amount of data, as explored in works like (Lauscher et al., 2020; Kim and Komachi, 2023), are important directions for future research.

## Acknowledgements

## References

Cheng-Han Chiang and Hung-yi Lee. 2023. A closer look into using large language models for automatic evaluation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8928–8942, Singapore. Association for Computational Linguistics.

Taisei Enomoto, Hwichan Kim, Tosho Hirasawa, Yoshinari Nagai, Ayako Sato, Kyotaro Nakajima, and Mamoru Komachi. 2024. TMU-HIT at MLSP 2024: How well can GPT-4 tackle multilingual lexical simplification? In *Proceedings of the 19th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2024)*, pages 590–598.

Marina Fomicheva, Shuo Sun, Erick Fonseca, Chrysoula Zerva, Frédéric Blain, Vishrav Chaudhary, Francisco Guzmán, Nina Lopatina, Lucia Specia, and André F. T. Martins. 2022. MLQE-PE: A multilingual quality estimation and post-editing dataset. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4963–4974.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, Errors, and Context: A Large-Scale Study of Human Evaluation for Machine Translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts. *Preprint*, arXiv:2401.04088.

Hwichan Kim and Mamoru Komachi. 2023. Enhancing few-shot cross-lingual transfer with target language peculiar examples. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 747–767, Toronto, Canada. Association for Computational Linguistics.

Katikela Sreeharsha Kishore and Rahimanuddin Shaik. 2024. Evaluating telugu proficiency in large language models: A comparative analysis of ChatGPT and Gemini. *Preprint*, arXiv:2404.19369.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203.

Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4483–4499, Online. Association for Computational Linguistics.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-Eval: NLG evaluation using GPT-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522.

Qingyu Lu, Baopu Qiu, Liang Ding, Kanjian Zhang, Tom Kocmi, and Dacheng Tao. 2024. Error analysis prompting enables human-like translation evaluation in large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 8801–8816, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

OpenAI. 2024. GPT-4o mini: advancing cost-efficient intelligence.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine*

*Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

# HW-TSC 2024 Submission for the Quality Estimation Shared Task

**Weiqiao Shan, Ming Zhu, Yuang Li, Mengyao Piao, Xiaofeng Zhao,**
**Chang Su, Min Zhang, Hao Yang, Yanfei Jiang**
Huawei Translation Services Center, China
shanweiqiao96@gmail.com, {zhuming47, liyuang3, piaomengyao1, zhaoxiaofeng14,
suchang8, zhangmin186, yanghao30, jiangyanfei}@huawei.com

## Abstract

Quality estimation (QE) is a crucial technique for evaluating the quality of machine translations without the need for reference translations. This paper focuses on Huawei Translation Services Center's (HW-TSC's) submission to the sentence-level QE shared task, named LLMs-enhanced-CrossQE. Our system builds upon the CrossQE architecture from our submission from last year, which consists of a multilingual base model and a task-specific downstream layer. The model input is a concatenation of the source and the translated sentences. To enhance performance, we fine-tuned and ensembled multiple base models, including XLM-R, InfoXLM, RemBERT, and CometKiwi. Specifically, we employed two pseudo-data generation methods: 1) a diverse pseudo-data generation method based on the corruption-based data augmentation technique introduced last year, and 2) a pseudo-data generation method that simulates machine translation errors using large language models (LLMs). Our results demonstrate that the system achieves outstanding performance on sentence-level QE test sets.

## 1 Introduction

Quality estimation (QE) aims to automatically assess machine translation outputs without requiring reference translations (Specia et al., 2018). We report the technical details of our approach to sentence-level quality prediction and fine-grained error span detection subtasks in the WMT 2024 QE shared task. Our team, Huawei Translation Services Center (HW-TSC), participated in direct assessment (DA) score in sentence-level quality prediction and the fine-grained error span detection tasks across all language pairs. Fine-tuning pre-trained language models, which provide rich semantic information, has become the standard approach for QE tasks (Rei et al., 2020). In this paper, we present LLMs-enhanced-CrossQE, HW-TSC's system for the sentence-level QE task, which leverages multiple pre-trained language models and data

augmentation techniques. The key aspects of our system design are summarized as follows:

- **Model**: We employed our previous year's architecture, CrossQE (Tao et al., 2022), as the foundation. For every language pair, models were individually fine-tuned. Additionally, we used CometKiwi (Rei et al., 2022), a multilingual QE model, and fine-tuned it for single language pairs.

- **Data augmentation**: Based on the corruption-based data generation (CDG) method used last year (Li et al., 2023), we propose a diverse CDG (D-CDG) method. Specifically, we generate more varied corrupted translations by combining multiple error types. Additionally, we rewrite source sentences using large language models (LLMs) to create pseudo-sentences containing errors that closely resemble those produced by machine translation systems. Finally, we employ a reference-based QE model to generate pseudo scores.

- **Ensemble**: For each language pair, we ensemble eight fine-tuned models to achieve optimal performance. These checkpoints originated from four base models: XLM-R (Conneau et al., 2020), InfoXLM (Chi et al., 2021), RemBERT (Chung et al., 2020), and CometKiwi (Rei et al., 2022), and three training dataset configurations: original dataset, augmented dataset, and augmented dataset followed by the original dataset. The ensemble weight for each checkpoint was optimized with Optuna (Akiba et al., 2019). On average, eight checkpoints were used per language pair after optimization. Additionally, we experimented with a naive weight ensemble approach based on the method proposed by Yadav et al. (2024), but it did not yield significant improvements.

535

Our system ranks first in the English-Tamil direction and second in several other directions in the direct assessment quality estimation task (Zerva et al., 2024). It significantly outperforms the baseline given by the competition organizers by a large margin. Additionally, we provide detailed results of each model with and without data augmentation in Table 3. To analyze the importance of each model in the ensemble, we present the ensemble weights in Figure 2 and 1. It is worth noting that the models fine-tuned with the proposed data augmentation technique were assigned higher weights in the ensemble.

## 2 Background

### 2.1 Task Description [1]

**Sentence-level QE with direct assessment (DA) anotations**: The goal is to predict the quality score for each source-target sentence pair. The golden-truth quality scores were obtained from human translators who rated each translation from 0 to 100. The scores from three or four translators were normalized and averaged to get the final score. This year's QE shared task has four language pairs with DA quality scores: English-Hindi (en-hi), English-Tamil (en-ta), English-Telegu (en-te) and English-Gujarati (en-gu). All languages have just 7,000 training samples.

**Fine-grained error span detection**: Participants of this task need to identify the error span (start and end indices) and the error severity (major or minor).

### 2.2 Base Models

- **XLM-R** (Conneau et al., 2020): A transformer-based masked language model trained on a massive multilingual corpus with more than two terabytes of data.

- **InfoXLM** (Chi et al., 2021): A cross-lingual pre-trained model that leverages multilingual masked language modeling, translation language modeling, and cross-lingual contrast learning.

- **RemBERT** (Chung et al., 2020): A rebalanced mBERT model with factorization of the embedding layers. The input embeddings are smaller and kept for fine-tuning, while the output embeddings are larger and discarded after pre-training.

[1] https://wmt-qe-task.github.io/

- **CometKiwi** (Rei et al., 2022): A multilingual reference-free QE model that uses a regression approach and is built on top of InfoXLM. It has been trained on direct assessments from WMT17 to WMT20 and the MLQE-PE corpus.

## 3 Method

### 3.1 Model Architecture

#### 3.1.1 Task1: Sentence-level QE with direct assessment (DA)

As shown in Equation 1 and 2, the embeddings of source sentence $\mathbf{s}$ and translated sentence $\mathbf{t}$ are concatenated in both orders $[\mathbf{s}, \mathbf{t}]$ and $[\mathbf{t}, \mathbf{s}]$ to form the input of pre-trained model $f_{base}$. The output token-level embedding sequences are processed by an average pooling layer to obtain vector representations $\mathbf{h}_{s1}$ and $\mathbf{h}_{t1}$ for source and translation respectively. These feature vectors are enhanced by taking their absolute difference and element-wise multiplication, as shown in Equation 3 and 4. Finally, all feature vectors are concatenated and fed into a regression head that predicts the final score $y$ (Equation 5). This architecture enables information exchange between source and translated sentences at an early stage of the network and has proven to be significantly more effective than combining cross-lingual information after the pre-trained model.

$$\mathbf{h}_{s1}, \mathbf{h}_{t1} = f_{base}([\mathbf{s}, \mathbf{t}]) \tag{1}$$

$$\mathbf{h}_{t2}, \mathbf{h}_{s2} = f_{base}([\mathbf{t}, \mathbf{s}]) \tag{2}$$

$$\mathbf{f}_1 = [\mathbf{h}_{s1}, \mathbf{h}_{t1}, |\mathbf{h}_{s1} - \mathbf{h}_{t1}|, \mathbf{h}_{s1} \odot \mathbf{h}_{t1}] \tag{3}$$

$$\mathbf{f}_2 = [\mathbf{h}_{s2}, \mathbf{h}_{t2}, |\mathbf{h}_{s2} - \mathbf{h}_{t2}|, \mathbf{h}_{s2} \odot \mathbf{h}_{t2}] \tag{4}$$

$$y = f_{score}([\mathbf{f}_1, \mathbf{f}_2]) \tag{5}$$

#### 3.1.2 Task2: Error span detection

For this task, we speculate that the understanding ability of large models may be helpful to the task, so we use the TowerInstruct-7B-v0.2 (Alves et al., 2024) model and the GPT-4o-mini (Islam and Moushi, 2024) model to cope with this task.

### 3.2 Data Augmentation

In this year's QE shared task, we adapted two data augmentation methods. 1) Text Editing, we implemented a D-CDG method based on the CDG proposed last year (Li et al., 2023), in which we constructed more diverse translation error data by

| Method | Description |
| --- | --- |
| Deletion | A random word in the translation was deleted. |
| Insertion | A random word in the translation was selected and inserted in a random position. |
| Substitution | A random word was replaced with another word in the translation. |

Table 1: Three available text editing methods.

| Method |
| --- |
| You are a Gujarati to English machine translation system. I will give you a correct parallel data pair, rewrite the target language (English) sentence with mistakes that you may have made while doing the translation, including but not limited to incorrect words, adding extra words, Omitting crucial words, wrong numbers or dates, deleting words, exchanging the position of two words, wrong numbers, incorrect punctuation, incorrect capitalization, grammar errors. The correct parallel data is: "$SRC", "$TGT". please just output the target language with 20%, 35%, 50% mistake token of the target length. |

Table 2: A prompt example for LLMs to generate pseudo QE training data based on a sample from the Gujarati to English QE training set, $SRC and $TGT represent the source and target languages in the sample, respectively.

incorporating multiple text editing approaches. 2) LLMs-generated pseudo-data. We generated translation data with errors more similar to those produced by machine translation systems using GPT-4o-mini and constructed parallel data pairs containing translation errors through the machine translation system.

For text editing, we employed three methods proposed last year to generate translation errors: Deletion, Insertion, and Substitution. Notably, this year, we generated translation sentences with more diverse translation errors by combining these three text editing methods with a certain probability. Specifically, each time we performed a text edit, we modified the original text with equal probability by sampling a text editing method from a subset of the three available text editing methods. Additionally, we also created a version of pseudo-data by directly translating the source language into the target language and then back-translating it.

For LLM-generated pseudo-data, we constructed a prompt using the GPT-4o-mini to generate a modified source language sentence multiple times with different proportions, correlating with the number of tokens in the sentence(see Table 2). This approach yielded multiple modified source language sentences containing error tokens that closely resemble those generated by translation systems. These modified sentences were then translated into the target language using a translation system. Similar to the text editing method, we scored the pseudo-parallel translation pairs using a

reference-based QE model [2] to create pseudo QE training data. It is worth noting that we constructed the scaling factor as the ratio between the corrupted translation score and the uncorrupted translation score ($\frac{f_{QE}(s,\hat{t},t)}{f_{QE}(s,t,t)}$), following the approach from last year.

## 4 Experiments

### 4.1 Experimental setups

Our system is built on top of the COMET package [3]. We fine-tuned four pre-trained models, namely XLM-R, InfoXLM, RemBERT and CometKiWi [4], on a single Nvidia Tesla V100 GPU with a batch size of 4, gradient accumulation of 8 and mean square error loss function. We stopped the training when there was no improvement in terms of Spearman correlation on the dev set for five test runs. For each language pair, the augmented dataset from text editing method, which contains more than ten times data than the original dataset, and the augmented dataset from LLMs, which contains about three times data than the original dataset, were all pre-generated instead of generated on-the-fly to improve training efficiency. Following last year's conclusion that the pseudo-data is more effective compared with the original data, we fine-tuned four base models by pseudo-data directly. The training step took around 10 hours with the augmented

---

[2] https://huggingface.co/Unbabel/wmt22-comet-da
[3] https://github.com/Unbabel/COMET
[4] https://huggingface.co/Unbabel/wmt22-cometkiwi-da

| Method | en-hi | en-ta | en-te | en-gu | Avg. |
|---|---|---|---|---|---|
| XLM-R | 0.616 | 0.663 | 0.434 | 0.643 | 0.589 |
| + aug (D-CDG) | 0.614 (-.002) | 0.675 (+.012) | 0.449 (+.015) | 0.657 (+.014) | 0.599 (+.010) |
| + aug (LLMs) | 0.469 (-.147) | 0.617 (-.046) | 0.412 (-.022) | 0.603 (-.040) | 0.525 (-.064) |
| InfoXLM | 0.595 | 0.670 | 0.443 | 0.664 | 0.593 |
| + aug (D-CDG) | 0.608 (+.013) | 0.657 (-.013) | 0.465 (+.022) | 0.671 (+.007) | 0.600 (+.007) |
| + aug (LLMs) | 0.478 (-.117) | 0.614 (-.056) | 0.418 (-.025) | 0.629 (-.035) | 0.535 (-.058) |
| RemBERT | 0.606 | 0.671 | 0.431 | 0.688 | 0.599 |
| + aug (D-CDG) | 0.604 (-.002) | 0.672 (+.001) | 0.432 (+.001) | 0.667 (-.021) | 0.594 (-.005) |
| + aug (LLMs) | 0.458 (-.148) | 0.606 (-.065) | 0.413 (-.018) | 0.617 (-.071) | 0.524 (-.075) |
| CometKiwi | 0.590 | 0.685 | 0.451 | 0.691 | 0.604 |
| + aug (D-CDG) | 0.594 (+.004) | 0.683 (-.002) | 0.465 (+.014) | 0.696 (+.005) | 0.610 (+.006) |
| + aug (LLMs) | 0.475 (-.115) | 0.630 (-.055) | 0.420 (-.031) | 0.662 (-.029) | 0.547 (-.057) |
| Ensemble (D-CDG) | 0.652 \ **0.719** | 0.716 \ 0.675 | 0.483 \ **0.482** | 0.717 \ 0.678 | 0.642 \ 0.637 |
| Ensemble (LLMs) | - | 0.712 \ **0.683** | 0.48 \ 0.474 | 0.714 \ **0.686** | 0.635 \ 0.614 |
| Ensemble (submit) | - \ **0.719** | - \ **0.683** | - \ **0.482** | - \ **0.686** | - \ **0.643** |

Table 3: Results for sentence-level QE in terms of **Spearman** correlation. We report the performance of using D-CDG and LLM-generated pseudo-data as a data augmentation approach(aug). Except for the last three rows which show the results on the dev \ test set, other results were based on the dev set.



Figure 1: The ensemble weights for each base model.

dataset.

With four base models and two data augmentation approaches, we obtained eight checkpoints for each language pair. We ensembled these checkpoints by taking the weighted average of the predicted scores. The weights were optimized using Optuna, an automatic hyperparameter search framework. We used the Spearman correlation as the optimization objective, setting the step size to 0.05, and conducted 1000 trials on the dev set.

### 4.2 Results

#### 4.2.1 Task1

The results of sentence-level QE in terms of Spearman correlation are shown in Table 3. Without data augmentation, CometKiwi has the best average correlation of 0.604, while XLM-R, InfoXLM, and RemBERT are close behind with around 0.590.

For the two data augmentation methods, we found that the D-CDG approach led to improvements across nearly all languages and models, as shown in Table 3. Additionally, this approach outperformed the original CDG method[5]. This suggests that rewriting a sentence by combining multiple diverse text editing methods within the same sentence is more effective than using only a single text editing method. Instead, for the pseudo-data generated by LLMs, we did not observe a positive effect on the dev set in all language directions, as shown in Table 3.

Furthermore, in the model ensemble, we observed that models with the D-CDG approach played a more important role. Specifically, the In-

[5]Reported in Table 1 of last year's paper

Figure 2: The ensemble weights for different training dataset configurations. 'w/o aug' and '+ aug' mean using the original or augmented dataset respectively. '+ aug & finetune' means training on the augmented dataset and then finetuning on the original one.

| Method | en-de | | | en-hi | | | en-es | | |
|---|---|---|---|---|---|---|---|---|---|
| | F1 | recall | precison | F1 | recall | precison | F1 | recall | precison |
| Tower-Instruct-7B | 0.178 | 0.181 | 0.175 | 0.015 | 0.008 | 0.300 | 0.118 | 0.082 | 0.209 |
| GPT-4o-mini | 0.119 | 0.315 | 0.073 | 0.361 | 0.398 | 0.331 | 0.146 | 0.249 | 0.103 |
| Baseline **(test set)** | **0.192** | - | - | **0.481** | - | - | **0.161** | - | - |
| Ensemble **(test set)** | 0.178 | 0.181 | 0.175 | 0.361 | 0.398 | 0.331 | 0.141 | 0.227 | 0.102 |

Table 4: Results for error span detection in terms of F1 score.

foXLM model with D-CDG was assigned a larger weight across all languages, as shown in Figure 1. We also noticed that the assignment of higher weights to models with D-CDG in the ensemble correlated with the base model's overall importance, if a base model received substantial attention, the corresponding model with D-CDG also tended to receive more weight.

Notably, models with LLM-generated pseudo-data were not assigned higher weights in the ensemble 2. However, in the test set, models with LLM-generated pseudo-data achieved better Spearman correlation scores in two languages(en-ta and en-gu). This may be attributed to the fact that LLMs generate more diverse pseudo-data, thereby enhancing the ensemble model's generalization ability. On the other hand, there may be a large gap between the dev set and the test set, the model with text editing data is overfitted to the dev set, while the models with LLMs pseudo-data introduce some regularization ability, which makes the ensemble model achieve better results on some languages.

This may be attributed to the fact that LLMs generate more diverse pseudo-data, thereby enhancing the ensemble model's generalization ability. Additionally, the discrepancy between the development set and the test set might have caused overfitting in

models trained with text editing data. In contrast, models incorporating LLM-generated pseudo-data introduced a regularization effect, enabling the ensemble model to achieve better results in certain languages.

### 4.2.2 Task2

The results for error span detection are displayed in Table 4. In the Table, we can see that the method of using the large language model alone to detect the error segment is lower than the baseline based on cometkiwi, but it is not far from it. In addition, we can see that the method based on GPT-4o-mini is much higher than the method without LLMs in recall. That's enough to see the potential of the large language models, if human preferences can be injected for fine tuning, there is a good chance that large language models will outperform cometkiwi-based methods.

## 5 Conclusion

This paper mainly presents HW-TSC's sentence-level QE system called LLMs-enhanced-CrossQE. Using our previous year's model CrossQE as the foundation, we conducted comprehensive experiments with various pre-trained models. To further enhance the robustness of all language pairs and provide various checkpoints for model ensem-

ble, we introduced a diverse pseudo-data generation method based on the corruption-based data augmentation technique proposed last year. Our system demonstrates strong performance across all language pairs with DA annotations in the sentence-level QE task. In the future, we plan to explore the use of LLMs to generate more diverse QE pseudo-data using more effective in-context learning techniques, such as chain-of-thought (CoT) prompting, or by transferring knowledge from LLMs to QE models through direct utilization of LLM parameters. Additionally, this paper presents only a brief investigation of the error span detection task. Therefore, we plan to further explore word-level and document-level QE tasks, which can improve the interpretability of QE and hold significant promise in the era of LLMs.

# References

Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A next-generation hyperparameter optimization framework. In *Proc. ACM SIGKDD*, pages 2623–2631.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, He-Yan Huang, and Ming Zhou. 2021. Infoxlm: An information-theoretic framework for cross-lingual language model pre-training. In *Proc. NAACL*, pages 3576–3588.

Hyung Won Chung, Thibault Fevry, Henry Tsai, Melvin Johnson, and Sebastian Ruder. 2020. Rethinking embedding coupling in pre-trained language models. *arXiv preprint arXiv:2010.12821*.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Édouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proc. ACL*, pages 8440–8451.

Raisa Islam and Owana Marzia Moushi. 2024. Gpt-4o: The cutting-edge advancement in multimodal llm. *Authorea Preprints*.

Yuang Li, Chang Su, Ming Zhu, Mengyao Piao, Xinglin Lyu, Min Zhang, and Hao Yang. 2023. Hw-tsc 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 835–840.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. In *Proc. EMNLP*, pages 2685–2702.

Ricardo Rei, Marcos Treviso, Nuno M Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José GC De Souza, Taisiya Glushkova, Duarte Alves, Luísa Coheur, et al. 2022. Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task. In *Proc. WMT*, pages 634–645.

Lucia Specia, Carolina Scarton, Gustavo Henrique Paetzold, and Graeme Hirst. 2018. *Quality estimation for machine translation*, volume 11. Springer.

Shimin Tao, Su Chang, Ma Miaomiao, Hao Yang, Xiang Geng, Shujian Huang, Min Zhang, Jiaxin Guo, Minghan Wang, and Yinglu Li. 2022. Crossqe: Hw-tsc 2022 submission for the quality estimation shared task. In *Proc. WMT*, pages 646–652.

Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. 2024. Ties-merging: Resolving interference when merging models. *Advances in Neural Information Processing Systems*, 36.

Chrysoula Zerva, Frederic Blain, Jos'e G. C. de Souza, Diptesh Kanojia, Sourabh Deoghare, Nuno M. Guerreiro, Giuseppe Attanasio, Ricardo Rei, Constantin Orasan, Matteo Negri, Marco Turchi, Rajen Chatterjee, Pushpak Bhattacharyya, Markus Freitag, and Andr'e Martins. 2024. Findings of the quality estimation shared task at wmt 2024. are llms closing the gap in qe? In *Proceedings of the Ninth Conference on Machine Translation*. Association for Computational Linguistics.

# HW-TSC's Participation in the WMT 2024 QEAPE Task

**Jiawei Yu[1]\* Xiaofeng Zhao[2], Min Zhang[2], Yanqing Zhao[2], Yuang Li[2],**
**Chang Su[2], Xiaosong Qiao[2], Miaomiao Ma[2], Hao Yang[2]**
[1]School of Informatics, Xiamen University, China
[2]Huawei Translation Services Center, Beijing, China
yujiawei@stu.xmu.edu.cn
{zhaoxiaofeng14,zhangmin186,yanghao30}@huawei.com

## Abstract

The paper presents the submission by HW-TSC in the WMT 2024 Quality-informed Automatic Post Editing (QEAPE) shared task for the English-Hindi (En-Hi) and English-Tamil (En-Ta) language pair. We use LLM for En-Hi and Transformer for EN-ta respectively. For LLM, we first continue pertrain the Llama3, and then use the real APE data to SFT the pre-trained LLM. As for the transformer in En-Ta, we first pre-train a Machine Translation (MT) model by utilizing MT data collected from the web. Then, we fine-tune the model by employing real APE data. We also use the data augmentation method to enhance our model. Specifically, we incorporate candidate translations obtained from an external Machine Translation (MT) system. Given that APE systems tend to exhibit a tendency of 'over-correction', we employ a sentence-level Quality Estimation (QE) system to select the final output, deciding between the original translation and the corresponding output generated by the APE model. Our experiments demonstrate that pre-trained MT models are effective when being fine-tuned with the APE corpus of a limited size, and the performance can be further improved with external MT augmentation. our approach improves the HTER by -15.99 points and -0.47 points on En-Hi and En-Ta, respectively.

## 1 Introduction

Automatic Post-Editing (APE) is a post-processing task in a Machine Translation (MT) workflow, aiming to automatically identify and correct errors in MT outputs (Chatterjee et al., 2020a). WMT has been holding APE task competitions in different languages and fields since 2015. Different from previous years, this year's APE task is a subtask of the QE task, named Quality-informed automatic post-editing (QEAPE) (Zerva et al., 2024). It proposes to combine quality estimation and automatic

post-editing in order to correct the output of machine translation. Participants are provided with a training set comprising 7,000 instances, a development set, and a test set, with each containing 1,000 instances. Each dataset consists of triplets — the source (*src*) sentences, the corresponding machine-translation (*mt*) outputs, and the human post-edited versions (*pe*) of the translations along with sentence-level QE annotations. Additionally, participants are permitted to utilize any additional data for systems training.

Typically, training an APE model requires large amount of training data. However, obtaining *pe* is an expensive task in terms of time and money. As a result, there exists a scarcity of large-scale APE datasets.

To address this challenge, numerous data augmentation techniques have been proposed (Junczys-Dowmunt and Grundkiewicz, 2016; Negri et al., 2018; Lee et al., 2020; Wei et al., 2020; Zhang et al., 2023). Wei et al. (2020) augment the APE training data with translations generated using a different MT system. Huang et al. (2022) train an external MT to obtain more datasets consistent with APE tasks. They also use Google translation to back translate the post-edits in the training set. Deoghare and Bhattacharyya (2022) augment the APE data by generating phrase-level APE triplets using SMT phrase tables. To ensure the quality of the synthetic data, they employ the LaBSE technique (Feng et al., 2022) to filter low-quality triplets.

We first collect our pre-training MT data from NLLB (Team et al., 2022), OpenSubtitles [1], TED2020 (Reimers and Gurevych, 2020), etc. To ensure the quality of the MT data, we use the LaBSE technique (Feng et al., 2022) and filter low-quality data. In our method, we use Google translation to back translate the post-edits in the training set. Subsequently, our dataset is structured as follows: the concatenation of source sentence, back

---

[1]https://www.opensubtitles.org/en/search/subs

translation and machine translation as the input, while the post-edits serve as the reference output.

Chatterjee et al. (2020b) have proven that APE systems often make unnecessary edits to translation output. To mitigate this issue of over-correction, we employ a sentence-level QE system to determine the final output, selecting between the APE system's output and the original machine-translated (*mt*) version.

Reflecting on the historical development, 2023 is recognized as the inaugural year for large-scale models, with researchers transitioning a variety of tasks to these models, including APE. Notable studies include those that combine Neural Machine Translation (NMT) with Large Language Models (LLM) for APE (Koneru et al., 2024), and comprehensive multi-stage, multilingual large models such as Tower (Alves et al., 2024b), which integrate both MT and APE. Drawing inspiration from Tower, our evaluation utilizes the continued pre-training (CPT) and supervised fine-tuning (SFT) to explore the potential of LLM.

When being evaluated on the test set, our approach improves the HTER (Snover et al., 2006) by -15.99 points and -0.47 points on En-Hi and En-Ta, respectively.

The contributions of our work are as follows:

- We filter low-quality MT data from the collected data using LaBSE-based filtering.

- We propose an APE paradigm based on LLM, including CPT and SFT.

- We utilize Google translation to back translate the post-edits to get *src'* for data augmentation.

- We employ a sentence-level QE system to select the most appropriate output, choosing between the APE-generated output and the original translation.

## 2 Related Work

Last year's WMT23 APE shard task mainly focuses on transfer learning and data augmentation. Yu et al. (2023) use a Transformer pre-trained on the provided synthetic APE data and then fine-tuned on the real APE data. Additionally, they utilize an external MT system to generate back-translations (with Google Translate [2] run on the post-edits in

the training set). They also integrate En-Mr parallel sentences from FLORES-200 (Costa-jussà et al., 2022). R-Drop (Liang et al., 2021), which regularizes the training inconsistency induced by dropout, is used to mitigate overfitting during the training phase. Besides, they use a sentence-level QE system to select the final output between the APE-generated output and the original translation.

Moon et al. (2023) center on data filtering techniques. With a focus on removing potentially harmful material from a model training perspective, the proposed method concentrates on eliminating the two extremes of the training data distribution: the (near-) perfect MT outputs on one side, and those that require complete rewriting on the other.

Another team "kaistai" is inspired by the recent surge of (LLMs) that have been successfully applied in a variety of language generation tasks. They use an LLM with specific prompts designed to generate either (a) post-edits or (b) post-edits along with the rationales behind them.

With experience in previous competitions, we also utilize an external MT system to generate back-translations in our transformer-based system. Additionally, we adopt a sentence-level QE system to select the final output.

## 3 Dataset

### 3.1 Data source

We first collect our MT data from the web, mainly from NLLB, OpenSubtitles, TED2020, etc. Then we filter the low-quality data using LABSE. After filtering, we get 3M En-Hi and 3M En-Ta parallel MT data. We first use our filtered MT data with 3M instances to pre-train our model. Then, we use the WMT24 official En-Hi and En-Ta APE datasets for fine-tuning, which consists of a training set and a development set. The training set for both language directions contains 7,000 APE triplets.

## 4 Method

### 4.1 LABSE filter

Before using the collected MT data to pretrain our model, we filter the low-quality parallel data by using the LaBSE-based filtering (Feng et al., 2022). We do this to ensure the quality of the MT data. To do so, we first generate embeddings of the En and Hi/Ta using the LaBSE model and normalize them. Then, we compute the cosine similarity between these normalized embeddings. We select the top

70% similarity parallel sentences as our filtered MT data.

## 4.2 LLM CPT + SFT

Due to the generative nature of the APE task, we believe that LLMs are well-suited for this purpose. Based on human evaluations, we have selected the Llama3-8B-Instruct model, which possesses proficiency in Hindi, as our foundational model. Inspired by the TowerInstruct (Alves et al., 2024a), we adopted a technical approach that combines CPT and SFT. Specifically, during the CPT phase, we utilized 3 million English-Hindi parallel corpora and employed LoRA training techniques. In the SFT phase, we created a customized prompt that, along with the training set provided by the organizers, constituted our SFT training dataset. Our prompt is as follows: "You are a post-editor. You improve translations from English to Hindi using the English source and Hindi translation. Do not provide any explanation or correction." The training paradigm is structured as [prompt: src <en2hi> mt <ape> response], where the response corresponds to the labels predicted by the model.

## 4.3 Fine-tuned Transformer

We basically treat the APE task as an NMT-like problem, which takes *src* and *mt* as input and generates *pe* autoregressively. Following previous works, we use a special token *<s>* to concatenate src and mt to generate the input sentence: [*src, <s>, mt*], while the target sentence is *pe*. Initially, we pretrain the MT model using the standard Transformer (Vaswani et al., 2017) structure on 3M En-Ta MT training data. Furthermore, we fine-tune the MT model using the APE dataset with the APE training objective. To further solve the problem of data scarcity, following (Yu et al., 2023), we use the Google translation system to create the *src'* from the provided *pe* text. We simply concatenate the *src'* with the original *src* and *mt* to form the new input: [*src, <s>, src', <s>, mt*]. Then, we use it in the same way as before, aiming to have the model learn complementary information from *src* and *src'*. During inference, the same input [*src, <s>, src', <s>, mt*] is employed to generate the output, thereby enabling the utilization of the external information derived from *src'*. Since there is no *pe* during inference, we translate the given *mt* into *src'* using Google Translate.



Figure 1: This figure, adapted from (Vaswani et al., 2017) shows the architecture of our model, where mt and augmented src' are concatenated with src before being input into the encoder, and post-edits are generated with the decoder.

## 4.4 Sentence-Level Quality Estimation

We use wmt22-cometkiwi-da (Rei et al., 2022) as our sentence-level QE model, which is a COMET quality estimation model. This model can be used for reference-free MT evaluation. It receives a source sentence and the respective translation and returns a single score between 0 and 1 that reflects the quality of the translation, where 1 represents a perfect translation. We use this model to rate both the original machine translation and the output generated by our APE system. We then compare the ratings for both sequences and select the one with a higher rating as the final output.

## 5 Experiment

### 5.1 Settings

Our transformer model on En-Ta is implemented with fairseq (Ott et al., 2019). Note that the vocabulary and encoder/decoder embeddings of our model are shared between two languages and contain 30K subtokens. We use the batch size of 30,720 tokens in the pre-training stage and 8,192 tokens in

| System | En-Hi | | | | En-Ta | | | |
|---|---|---|---|---|---|---|---|---|
| | BLEU↑ | HTER↓ | ChrF↑ | COMET↑ | BLEU↑ | HTER↓ | ChrF↑ | COMET↑ |
| Baseline (Do nothing) | 39.28 | 46.36 | 59.48 | 0.81 | 70.16 | 24.71 | 81.80 | 0.91 |
| Ours | 54.50 | 30.37 | 71.06 | 0.85 | 69.64 | **24.24** | 82.36 | 0.92 |
| swetaagrawal | 58.38 | **27.08** | 73.45 | 0.86 | 70.05 | 24.54 | 82.30 | 0.92 |

Table 1: Results on the WMT24 QE-APE En-Hi and En-Ta test set. A situation with a higher BLEU score but a lower HTER indicates a better result. The official primary evaluation metric for this task is HTER.

the fine-tuning stage. We leverage FP16 (mixed precision) training technique to accelerate training process. In all stages, we apply the Adam optimizer(Kingma and Ba, 2015) with $\beta_1 = 0.9$, $\beta_2 = 0.98$ to train the model, where the inverse square root schedule algorithm and warmup strategy are adopted for the learning rate. Concretely, We use a learning rate of 5e-4 with 20k warm-up steps in the pre-training stage and a learning rate of 5e-5 with 4k warm-up steps in the fine-tuning stage. Besides, we set the dropout to 0.1 in the pre-training stage, 0.3 in the fine-tuning stage, and the value of label smoothing to 0.1 in all stages. Early stopping is adopted with patience 10 and 30 epochs during pre-training and fine-tuning, respectively. During inference, we use beam search with a beam size of 10. Finally, We employ HTER (Snover et al., 2006), BLEU (Papineni et al., 2002), ChrF (Popovic, 2015), and COMET (Rei et al., 2022) as the evaluation metrics.

Our LLM on En-Hi is implemented with Llama-Factory(Zheng et al., 2024). The base model we used is Llama3-8B-Instruction. During the CPT phase, the batch size is set to 256, the learning rate to 1e-4, and training runs for 2 epochs with a precision of bf16. The maximum sequence length is 512 and pre-training is conducted using the LoRA method with a LoRA rank of 64.In the SFT phase, the batch size remains 256, the learning rate is adjusted to 1e-5, and training extends to 8 epochs with bf16 precision. We employ the AdamW optimizer, maintain a maximum sequence length of 512, and utilize PyTorch full_shard for training.

All our transformer models are trained on a Nvidia Tesla V100 GPU with 32GB memory and our LLMs are trained on 64 D910B with 32GB memory.

## 5.2 Result

Table 1 shows the experimental results evaluated on the test set, where the baseline result is produced by directly calculating scores between the provided

MT and PE. We outperform the baseline on HTER for -15.99 and -0.47 points on the En-Hi and En-Ta language pair.

| System | En-Hi | |
|---|---|---|
| | BLEU↑ | HTER↓ |
| Baseline (Do nothing) | 30.52 | 58.44 |
| Pretrain+finetune | **49.68** | **36.01** |
| +External MT | 49.01 | 37.16 |
| +Sentence-level QE | 39.13 | 43.77 |

Table 2: Results on the WMT24 QE-APE En-Hi development set.

| System | En->Ta | |
|---|---|---|
| | BLEU↑ | HTER↓ |
| Baseline (Do nothing) | 65.31 | 29.63 |
| Pretrain+finetune | 26.33 | 57.12 |
| +External MT | 33.80 | 45.31 |
| +Sentence-level QE | **66.11** | **27.66** |

Table 3: Results on the WMT24 QE-APE En-Ta development set.

Table 2 shows the En-Hi experimental results evaluated on the dev set. The baseline denotes the test MT result. As illustrated in table 2, the HTER decreased from 58.44 to 36.01 after applying CPT+SFT, reflecting a reduction of 22.43. However, no performance improvement was observed with the addition of back-translation data. We hypothesize that this is due to the sufficiently robust performance of the CPT+SFT, which diminishes the impact of the back-translation data on further enhancement. Upon integrating QE labels, the HTER increased to 43.77 compared with CPT+SFT, an increase of 7.76. We think the QE label may not be accurate enough in En-Hi, resulting in performance loss.

Table 3 shows the En-Ta experimental results evaluated on the dev set. The first experiment is performed by fine-tuning all parameters of the pre-trained Transformer on the official training set, which increases by 27.49 in HTER compared with the baseline. Due to the lack of high-quality En-Ta MT data, the pre-training MT datasets we collected were mostly synthetic and of poor quality. This hinders the capabilities of MT models, which further results in fine-tuned APE models that also perform poorly. The experiment of adding external MT for data augmentation shows some improvement in performance. Toward the end, we utilize a sentence-level QE system to rate both the original translation and the APE output. We then select one of them with a higher rating as the final output of our APE system. With the combination of the APE model and sentence-level QE system, we see that the HTER decreases to 27.66, and the BLEU score increases to 66.11 points.

## 6 Conclusion

This paper presents our APE system submitted to the WMT 2024 QEAPE En-Hi and EN-Ta task. In our approach, we first filter low-quality MT data from the collected data using LaBSE-based filtering. Then we employ the data augmentation method to build the [*src, <s>, src', <s>, mt*] additional training datasets. Besides, We propose an APE paradigm based on LLM, including CPT and SFT. Moreover, we explore the sentence-level QE system to discard low-quality APE outputs. Evaluation of our APE system shows that our approach achieves gains on the WMT-24 APE development and test sets.

## 7 Acknowledgements

We would like to thank the anonymous reviewers. Their insightful comments helped us in improving the current version of the paper.

## References

Duarte M. Alves, José Pombal, Nuno Miguel Guerreiro, Pedro Henrique Martins, João Alves, M. Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024a. Tower: An open multilingual large language model for translation-related tasks.

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024b. Tower: An open multilingual large language model for translation-related tasks.

Rajen Chatterjee, Markus Freitag, Matteo Negri, and Marco Turchi. 2020a. Findings of the WMT 2020 shared task on automatic post-editing. In *Proc. of WMT@EMNLP*.

Rajen Chatterjee, Markus Freitag, Matteo Negri, and Marco Turchi. 2020b. Findings of the WMT 2020 shared task on automatic post-editing. In *Proc. of WMT@EMNLP*.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Y. Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *CoRR*, abs/2207.04672.

Sourabh Dattatray Deoghare and Pushpak Bhattacharyya. 2022. IIT bombay's WMT22 automatic post-editing shared task submission. In *Proc. of WMT*.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proc. of ACL*.

Xiaoying Huang, Xingrui Lou, Fan Zhang, and Tu Mei. 2022. Lul's WMT22 automatic post-editing shared task submission. In *Proc. of WMT*.

Marcin Junczys-Dowmunt and Roman Grundkiewicz. 2016. Log-linear combinations of monolingual and bilingual neural machine translation models for automatic post-editing. In *Proc. of WMT*.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *Proc. of ICLR*.

Sai Koneru, Miriam Exel, Matthias Huck, and Jan Niehues. 2024. Contextual refinement of translations: Large language models for sentence and document-level post-editing. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2711–2725, Mexico City, Mexico. Association for Computational Linguistics.

WonKee Lee, Jaehun Shin, Baikjin Jung, Jihyung Lee, and Jong-Hyeok Lee. 2020. Noising scheme for data augmentation in automatic post-editing. In *Proc. of WMT@EMNLP*.

Xiaobo Liang, Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, and Tie-Yan Liu. 2021. R-drop: Regularized dropout for neural networks. In *Proc. of NeurIPS*.

Hyeonseok Moon, Seugnjun Lee, Chanjun Park, Jaehyung Seo, Sugyeong Eo, and Heuiseok Lim. 2023. What is the resultful data?: Empirical study on the adaptability of the automatic post-editing training data. In *Proceedings of the Eighth Conference on Machine Translation, WMT 2023, Singapore, December 6-7, 2023*. Association for Computational Linguistics.

Matteo Negri, Marco Turchi, Rajen Chatterjee, and Nicola Bertoldi. 2018. ESCAPE: a large-scale synthetic corpus for automatic post-editing. In *Proc. of LREC*.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Demonstrations*, pages 48–53. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pages 311–318. ACL.

Maja Popovic. 2015. chrf: character n-gram f-score for automatic MT evaluation. In *Proc. of WMT@EMNLP*.

Ricardo Rei, Marcos V. Treviso, Nuno Miguel Guerreiro, Chrysoula Zerva, Ana C. Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte M. Alves, Luísa Coheur, Alon Lavie, and André F. T. Martins. 2022. Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task. In *Proc. of WMT*.

Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Matthew G. Snover, Bonnie J. Dorr, Richard M. Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proc. of AMTA*. Association for Machine Translation in the Americas.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proc. of NeurIPS*.

Daimeng Wei, Hengchao Shang, Zhanglin Wu, Zhengzhe Yu, Liangyou Li, Jiaxin Guo, Minghan Wang, Hao Yang, Lizhi Lei, Ying Qin, and Shiliang Sun. 2020. Hw-tsc's participation in the WMT 2020 news translation shared task. In *Proc. of WMT@EMNLP*.

Jiawei Yu, Min Zhang, Yanqing Zhao, Xiaofeng Zhao, Yuang Li, Chang Su, Yinglu Li, Miaomiao Ma, Shimin Tao, and Hao Yang. 2023. Hw-tsc's participation in the WMT 2023 automatic post editing shared task. In *Proceedings of the Eighth Conference on Machine Translation, WMT 2023, Singapore, December 6-7, 2023*, pages 926–930. Association for Computational Linguistics.

Chrysoula Zerva, Frederic Blain, Jos'e G. C. de Souza, Diptesh Kanojia, Sourabh Deoghare, Nuno M. Guerreiro, Giuseppe Attanasio, Ricardo Rei, Constantin Orasan, Matteo Negri, Marco Turchi, Rajen Chatterjee, Pushpak Bhattacharyya, Markus Freitag, and Andr'e Martins. 2024. Findings of the quality estimation shared task at wmt 2024. are llms closing the gap in qe? In *Proceedings of the Ninth Conference on Machine Translation*. Association for Computational Linguistics.

Min Zhang, Xiaofeng Zhao, Zhao Yanqing, Hao Yang, Xiaosong Qiao, Junhao Zhu, Wenbing Ma, Su Chang, Yilun Liu, Yinglu Li, Minghan Wang, Song Peng, Shimin Tao, and Yanfei Jiang. 2023. Leveraging chatgpt and multilingual knowledge graph for automatic post-editing. In *International Conference on Human-Informed Translation and Interpreting Technology (HiT-IT 2023)*. Accepted for publication.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. *CoRR*, abs/2403.13372.

# Expanding the FLORES+ Multilingual Benchmark with Translations for Aragonese, Aranese, Asturian, and Valencian

**Juan Antonio Pérez-Ortiz,**[*†] **Felipe Sánchez-Martínez,**[†] **Víctor M. Sánchez-Cartagena,**[†]
**Miquel Esplà-Gomis,**[†] **Aarón Galiano-Jiménez,**[†] **Antoni Oliver,**[‡] **Claudi Aventín-Boya,**[‡]
**Alejandro Pardos,**[◇] **Cristina Valdés,**[•] **Jusèp Loís Sans Socasau,**[□] **Juan Pablo Martínez**[△]

[†]Universitat d'Alacant {japerez,fsanchez,vm.sanchez,miquel.espla,aaron.galiano}@ua.es
[*]Valencian Graduate School and Research Network of Artificial Intelligence, ValgrAI
[‡]Universitat Oberta de Catalunya {aoliverg,caventinb}@uoc.edu
[◇]Universidad de Zaragoza apardoscalvo@gmail.com
[•]Academia de la Llingua Asturiana / Universidad de Oviedo cris@uniovi.es
[□]Institut d'Estudis Aranesi – Acadèmia Aranesa dera Lengua Occitana sanssocasau@gmail.com
[△]Academia Aragonesa de la Lengua / Universidad de Zaragoza jpmart@unizar.es

## Abstract

In this paper, we describe the process of creating the FLORES+ datasets for several Romance languages spoken in Spain, namely Aragonese, Aranese, Asturian, and Valencian. The Aragonese and Aranese datasets are entirely new additions to the FLORES+ multilingual benchmark. An initial version of the Asturian dataset was already available in FLORES+, and our work focused on a thorough revision. Similarly, FLORES+ included a Catalan dataset, which we adapted to the Valencian variety spoken in the Valencian Community. The development of the Aragonese, Aranese, and revised Asturian FLORES+ datasets was undertaken as part of a WMT24 shared task on translation into low-resource languages of Spain.

## 1 Introduction

Although notable advances have been reported in the realm of machine translation (MT) in recent years, performance for the so-called low-resource languages (Ranathunga et al., 2023) still lags behind that of languages with more extensive linguistic resources. Increasing the availability of such resources is a key factor for enabling MT to progressively and ultimately cover all languages in the world. Several relevant initiatives, such as FLORES-101 (Goyal et al., 2022), FLORES-200 (Costa-jussà et al., 2024), Seed (Maillard et al., 2023), and NTREX (Federmann et al., 2022), have recently addressed this challenge by providing multilingual datasets covering up to 200 languages. While these datasets were originally created as part of other projects, the Open Language Data Initiative[1] (OLDI) now leads a collective endeavor to ex-



Figure 1: Some of the languages of Spain: Spanish, Catalan, Galician, Basque, Asturian, Aragonese and Aranese. The language names listed are shown in the corresponding colors of their regions on the map. Note that Spanish is spoken throughout the entire country, but the map highlights the regions where other languages are co-official or regionally predominant. [Source: Wikimedia Commons, https://commons.wikimedia.org/wiki/File:Spain_languages.svg]

pand the number of supported languages. In particular, a shared task was proposed to extend OLDI's open datasets to more languages for the Ninth Conference on Machine Translation (WMT24).[2] This paper presents our efforts to extend the OLDI's FLORES+ dataset to four low-resource Romance languages spoken in Spain as part of this task, namely Aragonese, Aranese, Asturian, and Valencian. The resulting datasets can be downloaded from the repository of the PILAR corpus (Galiano-Jiménez et al., 2024).[3]

Most languages spoken in Spain belong to the Romance language family, except for Basque, which is a *language isolate* with no known genetic

---

[1]https://oldi.org

[2]https://www2.statmt.org/wmt24/open-data.html
[3]https://github.com/transducens/PILAR

547

relationship to any other language. Among the languages of Spain, Spanish is the only one with official status across the entire country. Other languages, including Catalan (primarily spoken in Catalonia, the Valencian Community, and the Balearic Islands, with a few thousand speakers in Aragon), Galician (mainly spoken in Galicia and some adjacent areas across the border), and Basque (spoken in the Basque Country and parts of Navarre), have official status in their respective autonomous regions (in certain areas for Navarre, and not in Aragon for Catalan). Additionally, Aranese is official in the Val d'Aran in accordance with the Statute of Autonomy of Catalonia. However, Asturian and Aragonese do not have official status in their respective regions, though they are recognized and protected as cultural heritage.[4] The Acadèmia Valenciana de la Llengua[5] (AVL) considers Valencian to be an alternative name for the Catalan language, as well as the term used to refer specifically to the variety of this language spoken in the Valencian Community.

Figure 1 illustrates the geographical distribution of these languages, while Table 1 provides estimates of their number of speakers in Spain. It is important to note that Spanish is spoken nationwide, while proficiency in regional languages varies among people in bilingual or diglossic areas.

We present our contribution to the FLORES+ benchmark,[6] focusing on four Romance languages spoken in Spain. Our target languages are Aragonese, Asturian, Aranese, and Valencian.

These languages were chosen due to the evident social and governmental interest in preserving them. For Asturian, there are ongoing efforts to achieve official language status in Asturias. In contrast, Aranese is already an official language in the small central Pyrenees valley where it is moslty spoken, although it has very few speakers and limited textual resources. Aragonese is recognized and protected in Aragon as a heritage language, yet political support for granting it official status remains limited. Lastly, Valencian, a variant of Catalan with notable lexical differences, was included because having MT systems generate it directly rather than adapting from Catalan significantly reduces the need for additional post-editing.

## 2 Target Languages Overview

This section provides a brief overview of the languages covered in this paper.

### 2.1 The Aragonese Language

Aragonese is a Romance language spoken in the Pyrenees valleys of Aragon, primarily in the *comarcas* of Somontano de Barbastro, Jacetania, Alto Gállego, Sobrarbe, Ribagorza and Hoya de Huesca/Plana de Uesca.

Until recently, Aragonese has had several alternative orthographic norms, none of them official. In 1987, a number of associations organized a congress where a quasi-phonetic spelling system for Aragonese was approved, called *Normas Graficas de l'Aragonés*. It was commonly used in the subsequent two decades. The *Sociedat de Lingüistica Aragonesa* (SLA, Aragonese Linguistic Society), an association established in 2004, published a set of spelling rules (SLA rules) in 2006 based on the written tradition of Medieval Aragonese. In 2010, the *Estudio de Filología Aragonesa – Academia de l'Aragonés* (EFA-ACAR), a private entity created by the II Congress on Aragonese Language, approved and published the *Propuesta Ortografica de l'Academia de l'Aragonés*, more aligned with etymology. Finally, the *Academia Aragonesa de la Lengua*[7] (AAL, Aragonese Language Academy), a public body created by the Law of the Languages of Aragon, was established in 2021, and approved in April 2023 the standard orthographic norm for Aragonese. This standard orthography has been adopted by associations (including EFA-ACAR) as well as by the main pub-

---

[4] It is also worth noting that there are other languages in Spain that could be considered, but there is less consensus on whether they are part of an enclosing language or distinct languages on their own, debates that echo the famous saying attributed to sociolinguist Max Weinreich that "a language is a dialect with an army and a navy." For example, Eonavian or Galician-Asturian is a set of Romance dialects that has been classified either as northeastern varieties of Galician, as a linguistic group of its own, or as a transitional dialect between western Asturian and Galician. Another example is the Leonese language, which is currently spoken in the northern and western parts of the historical region of León in Spain (the modern provinces of León, Zamora, and Salamanca). Leonese, however, is considered part of the Asturleonese linguistic group, along with the dialects of Asturian. Similarly, Fala (spoken in the northwestern part of Extremadura) and Cantabrian also exhibit this ambiguity being classified, depending on the perspective, as distinct languages, hybrid dialects, or intermediate linguistic varieties.

[5] https://www.avl.gva.es

[6] FLORES+ is OLDI's evaluation benchmark for multilingual machine translation, building upon the FLORES-200 benchmark (Costa-jussà et al., 2024), which spans over 200 languages and comprises 2 009 sentences split into development (dev) and devtest sets. Texts come from English Wikimedia (more exactly, Wikinews, Wikijunior and Wikivoyage).

[7] https://www.academiaaragonesadelalengua.org

| Language | Speakers | Bibliographic reference |
|---|---|---|
| Spanish | 47 000 000 | (Instituto Nacional de Estadística, 2022, pages 6–7) |
| Catalan (incl. Valencian) | 9 000 000 | (Generalitat de Catalunya, 2007, Table 2) |
| Valencian | 3 500 000 | (Generalitat Valenciana, 2021, page 6) |
| Galician | 2 100 000 | (Observatorio da Lingua Galega, 2007, page 11) |
| Basque | 1 200 000 | (Depto. de Cultura y Política Lingüística, 2023, page 2) |
| Asturian | 250 000 | (Llera Ramo, 2018, Figure 6) |
| Aragonese | 25 000 | (Reyes et al., 2017, Table 5) |
| Aranese | 4 500 | (Generalitat de Catalunya, 2019, page 4) |

Table 1: Approximate number of speakers *in Spain* for some of the languages spoken in the country. When the reference includes different figures for the number of speakers depending on their ability to speak, understand, read, or write the language, we provide the data for the number of people who can speak it, including those with even a basic level of proficiency.

lishers in Aragonese. In August 2024, the Government of Aragon established this standard spelling for administrative and educational uses.[8]

A recent study reveals that Aragonese remains alive across the entire area where it is traditionally spoken (Eito et al., 2024). As of today, it is estimated that around 8 000 people use it on a daily basis, according to the most optimistic estimates, within a broader group of approximately 25 000 people who claim to have knowledge of the language (Reyes et al., 2017). Its use is much more prevalent in family and neighbourly interactions, especially among older individuals and in rural communities. However, the use of the language declines significantly in more formal contexts and outside the immediate environment of the speakers. Consequently, intergenerational transmission is severely at risk; some initial measures, such as teaching in schools or the language limited presence in the media, may not be sufficient to ensure the language survival.

### 2.2 The Aranese Language

Aranese is a standardized form of the Pyrenean Gascon variety of the Occitan language spoken in the Val d'Aran in northwestern Catalonia. Aranese is one of the three official languages in Catalonia, alongside Catalan and Spanish.[9] In the case of

Aranese, its official status is limited to the territory of the Val d'Aran.

The Aranese language is regulated by the *Institut d'Estudis Aranesi–Acadèmia Aranesa dera Lengua Occitana*.[10] The main goal of this institution is to establish and update the linguistic norms of the Aranese variety of Occitan and ensure that the standardization process is consistent throughout its linguistic area.

Regarding education, the Occitan language is included, along with Catalan, Spanish, English, and optionally French, in the school system of the valley at all compulsory levels. However, the presence of Aranese in higher education, such as in the baccalaureate program, is virtually non-existent. According to the Linguistic Census of Aranese 2001,[11] only 34.2% of the population in the Val d'Aran have Aranese as their mother tongue, and only 25.8% use Aranese exclusively at home. The estimated number of speakers is approximately 4 500 (Generalitat de Catalunya, 2019).

### 2.3 The Asturian Language

Asturian is one of the Iberian Romance languages, closely related to Galician-Portuguese and Castilian Spanish, and a language historically influencing the Silver Way[12] to the Portuguese border (in fact, Asturian-based Mirandese is official in Miranda do Douro, in Portugal).

---

[8] https://www.boa.aragon.es/cgi-bin/EBOA/BRSCGI?CMD=VEROBJ&MLKOB=1347070761212

[9] The Val d'Aran was relatively well connected to Occitania in the north, but remained isolated to the south until the construction of the Vielha Tunnel, first in 1948 and then with a new one in 2007. The changes experienced by the valley due to its shift toward the south following the opening of the Vielha Tunnel, and especially the development of tourism, have led to a massive influx of migrants initially from the rest of Spain and later from other countries. This migration has resulted in drastic changes in the linguistic practices of the

Aranese people. Occitan has increasingly been confronted with Spanish and Catalan, losing the dominance it once had in past centuries.

[10] https://www.institutestudisaranesi.cat

[11] https://llengua.gencat.cat/web/.content/documents/altres/arxius/aran_cens.pdf

[12] The Silver Way is a route of the *Camino de Santiago* that runs from southern Spain to Santiago de Compostela.

The *Academia de la Llingua Asturiana*[13] (Academy of the Asturian Language) is an official institution of the Government of the Principality of Asturias that promotes and regulates the Asturian language. Its main objectives include researching and standardizing the Asturian language, developing language usage norms and dictionaries, promoting its use and education, compiling its lexicon, fostering research related to Asturian, awarding literary prizes, and protecting the language rights of Asturian users. As a result, the Asturian language has clear and widely accepted spelling and orthographic rules.

Although it does not have official language status, Asturian is protected under the Statute of Autonomy of Asturias. In many schools, children can take Asturian-language classes, and in some schools, it is offered as an elective language.

In a report for the BBC, the President of the Academy of the Asturian Language pointed out that "at present about 250 000 people are able to understand, speak, read and write in Asturian, that is, a 25% of the population of the region" (Hernández, 2022; Llera Ramo, 2018). As the latest 3rd Sociolinguistic Study on the Asturian Language (Llera Ramo, 2018) points out, 87% of the population identify with both Asturian and Spanish identities. Of these, 64% have a balanced sense of both identities, while 18% feels predominantly Asturian over Spanish. Literacy rates are high: 90% of the population report understanding Asturian, though 29% describe their knowledge as passive, while 38% are able to read and 25% to write.

## 2.4 The Catalan and Valencian Languages

Catalan is a Romance language with official status in three autonomous communities in Spain: Catalonia, the Balearic Islands and the Valencian Community. It is also the official language of Andorra, a small state in the Pyrenees, and it has semi-official status in the Italian *comune* of Alghero, in the island of Sardinia. Catalan is also spoken in the south of France and in some parts of the Spanish regions of Aragon and Murcia.

Catalan is usually named Valencian in the Valencian Community, but both terms academically refer to the same language. Catalan is divided into two major dialect groups: Eastern and Western (see Figure 2). The main difference is in the pronunciation of unstressed 'a' and 'e'. In Eastern dialects, these

Figure 2: Catalan's two main dialects (Eastern/Western) shown divided by the dashed line. [Source: modified from the original in Wikimedia Commons, `https://commons.wikimedia.org/wiki/File:Catalan_dialects-en.png`]

sounds have merged into /ə/, while in Western dialects, they remain distinct as /a/ and /e/. There are also some other differences in pronunciation, verb forms, and vocabulary. Western Catalan includes the Northwestern Catalan and Valencian dialects. The Eastern group includes four dialects: Central Catalan, Balearic, Rossellonese, and Algherese.

The *Institut d'Estudis Catalans*[14] in Catalonia and the *Acadèmia Valenciana de la Llengua*[15] are the two institutions that regulate the Catalan and Valencian varieties, respectively. The relationship between these two institutions has not always been easy, but in recent years, they have achieved some degree of coordination.

In 2018[16], 94.4% of the population aged 15 or older in Catalonia could understand Catalan, 81.2% could speak it, 85.5% could read it, and 65.3% could write it. In the Valencian Community, 79.4% of the population aged 15 or older could understand Valencian, 54.9% could speak it, 60.9% could read it, and 44.4% could write it (Generalitat Valenciana, 2021). Valencian speakers make up approximately 3.5 million of the total 9 million Catalan speakers in Spain (Generalitat Valenciana, 2021; Generalitat de Catalunya, 2007).

Although standard Catalan and Valencian are highly similar and largely mutually intelligible, the differences and the effort required to adapt Cata-

---

[13] `https://alladixital.org`

[14] `https://www.iec.cat`

[15] `https://www.avl.gva.es`

[16] `https://www.idescat.cat/indicadors/?id=basics&n=10367&lang=en`

lan texts into Valencian are significant to justify the creation of a dedicated FLORES+ dataset for Valencian. This will support the development and evaluation of MT systems specifically designed for this variant.

## 3 Development of the FLORES+ Datasets

This section describes the workflow and resources used to produce the FLORES+ datasets for Aragonese, Aranese, Asturian and Valencian.

### 3.1 The Aragonese FLORES+ Dataset

A first draft of the FLORES+ dev and devtest sentences for Aragonese was initially obtained from Spanish using the Spanish–Aragonese[17] rule-based machine translation system Apertium (Forcada et al., 2011). This translation was subsequently post-edited by specialists proficient in Aragonese. Finally, the post-edited translation was reviewed by a member of the Academia Aragonesa de la Lengua. In this last step, the reviewer had also in view the English version of the FLORES+ dataset; an important number of translations were modified for better agreement with the English original data.

The translators and reviewers relied on the following resources to inform their decisions:

- The orthography of Aragonese[18] published by the Academia Aragonesa de la Lengua in 2023.

- The grammar of Aragonese published by the EFA-ACAR.[19]

- The dictionary of Aragonese,[20] published by the EFA-ACAR.

- The Aragonese dictionary *Tresoro d'a Luenga Aragonesa*,[21] a lexicographical research project developed by the *Instituto de Estudios Altoaragoneses*, in collaboration with the Government of Aragon.

- The Spanish–Aragonese and Aragonese–Spanish dictionary *Aragonario*,[22] developed under the project Linguatec and published by the Government of Aragon.

**Justification for the use of MT.** While it might be argued that using machine translation (MT) in a dataset like FLORES+, conceived as a benchmark for MT systems, is not the ideal procedure, our decision is justified by three main reasons: first, the two-step workflow of MT followed by post-editing is common practice for this language, with many existing texts produced in this manner; second, the scarcity of qualified linguists and translators for Aragonese made it difficult to complete the task within the required timeframe without the support of an MT system; third, rule-based systems like Apertium typically exhibit strong performance for closely related Romance languages, resulting in minimal translationese that does not significantly impact the overall quality.

### 3.2 The Aranese FLORES+ Dataset

Similarly to Aragonese (see above), the data for Aranese was initially obtained by translating the sentences in the Catalan FLORES+ dev and devtest datasets using the Apertium rule-based MT system (Forcada et al., 2011) for Catalan–Aranese.[23] The revision process was then performed in two steps. Firstly, a professional reviewer with wide experience in translation and revision with proficiency in Aranese was presented with the French, Catalan and Occitan versions of the FLORES+, along with the machine translated version into Aranese. Finally, the post-edited translation was reviewed by different individuals, who are native speakers, from the Institut d'Estudis Aranesi (IEA). The use of MT is motivated by the same reasons as those for Aragonese (see the end of the previous section).

During the post-editing and the subsequent review process, the guidelines provided by the IEA were strictly followed to ensure that the translation into Aranese aligned with their recommendations. Specifically, the translators and reviewers based their decisions on the following resources published by the IEA:

---

[17]https://github.com/apertium/apertium-spa-arg, release 0.5.0, the latest release available at the time of translation.
[18]https://academiaaragonesadelalengua. org/sites/default/files/ficheros-pdf/ ortografia-de-laragones_web_an.pdf
[19]http://www.academiadelaragones.org/biblio/ Edacar10_GBAprovisional.pdf
[20]http://www.academiadelaragones.org/biblio/ Edacar13.pdf
[21]http://diccionario.sipca.es/fabla/faces/ index.xhtml

---

[22]https://aragonario.aragon.es/
[23]https://github.com/apertium/apertium-oci-cat, release 1.0.8.

- Dictionary for Aranese;[24] erratum 2021;[25] extensions 2021,[26] 2022[27] and 2023.[28]

- Grammar for Aranese.[29]

- Orthographic vocabulary of Aranese.[30]

- Winter sports vocabulary of Aranese.[31]

- Computer and informatics vocabulary of Aranese.[32]

- The Aranese verbs.[33]

### 3.3 The Asturian FLORES+ Dataset

The dev and devtest datasets we are contributing for Asturian are a corrected version of the original FLORES+ dataset, which were initially translated from English by Meta using professional translators as part of their *no language left behind* initiative (Costa-jussà et al., 2024). We had this translation into Asturian reviewed by native speakers, some of whom are members of the *Academia de la Llingua Asturiana*, philologists and a renowned writer, translator and activist for the Asturian language. The revision process was carried out twice by different people. In the first round, the reviewers were presented with the Spanish text and the existing version of the Asturian FLORES+, and in the second round, with the Spanish FLORES+ sentences and the first revised version.

During the review process, special attention was paid to adhering to the guidelines provided by the Academia to ensure that the translation aligns with

their recommendations by relying on the following resources:

- Diccionariu de la Llingua Asturiana.[34]

- Gramática de la Llingua Asturiana.[35]

- Normes ortográfiques.[36]

### 3.4 The Valencian FLORES+ Dataset

The dataset for Valencian was created by adapting the existing Catalan version of the FLORES+ devtest set.[37] The work was carried out by a single native speaker of the Valencian variant, who holds a university degree and has experience translating into Valencian and reviewing texts in this language.

To expedite the process, LanguageTool[38] was employed to apply an initial set of changes to the original text. These changes were then manually reviewed, and additional necessary modifications were made, considering a list of lexical items that differ across the various Catalan dialects (see below). In cases of uncertainty, the dictionary of the Acadèmia Valenciana de la Llengua was consulted. The following resources were used to guide the modifications to the original Catalan text:

- List of lexical items that differ across the different dialects of the Catalan language.[39]

- Diccionari normatiu valencià.[40]

- Linguistic criteria for the institutional use of the Valencian dialect in the Valencian universities.[41]

### 3.5 Modifications Across Versions of the Datasets

As previously mentioned, the datasets for Aragonese, Aranese, and Asturian have undergone several iterations during their development. Table 2 shows a comparison of the translation error rate (TER) between these versions, labeled as v1,

---

[24] http://www.institutestudisaranesi.cat/wp-content/uploads/2021/04/DICCIONARI-DER-ARAN%C3%89S.pdf

[25] http://www.institutestudisaranesi.cat/wp-content/uploads/2020/12/ERRATA-WEB.pdf

[26] http://www.institutestudisaranesi.cat/wp-content/uploads/2020/12/500-1-g%C3%A8r-2021-WEB.pdf

[27] http://www.institutestudisaranesi.cat/wp-content/uploads/2022/01/WEB-AMPLIACION-01-01-2022.pdf

[28] http://www.institutestudisaranesi.cat/wp-content/uploads/2023/03/Ampliacion-2023.pdf

[29] http://www.institutestudisaranesi.cat/wp-content/uploads/2021/04/gramatica-aranes.pdf

[30] http://www.institutestudisaranesi.cat/wp-content/uploads/2019/11/33520-vocabulari-ortografic.pdf

[31] http://www.institutestudisaranesi.cat/wp-content/uploads/2019/11/33288-vocabulari-esports-iuern.pdf

[32] http://www.institutestudisaranesi.cat/wp-content/uploads/2018/11/TECNOLOGIA.pdf

[33] http://www.institutestudisaranesi.cat/wp-content/uploads/2019/11/33287-els-ve%CC%80rbs.pdf

[34] https://diccionariu.alladixital.org/

[35] https://alladixital.org/wp-content/uploads/2022/08/Gramatica-de-la-Llingua-Asturiana.pdf

[36] https://alladixital.org/wp-content/uploads/2024/01/Normes-Ortografiques-8a-edicion-FINAL-3.pdf

[37] It is important to note that Valencian is the only target language for which only one of the two FLORES+ datasets (the devtest set) has been translated, instead of both.

[38] https://languagetool.org

[39] https://ca.wikipedia.org/wiki/Llista_diat%C3%B2pica_del_l%C3%A8xic_catal%C3%A0

[40] https://www.avl.gva.es/lexicval

[41] https://sl.ua.es/en/assessorament/documentos/criteris-linguistics.pdf

| Language | v1 → v2 | v2 → v3 | v1 → v3 |
|---|---|---|---|
| Aragonese dev | 4.8 | 22.2 | 24.3 |
| Aragonese devtest | 4.0 | 26.8 | 28.4 |
| Aranese dev | 1.8 | 54.2 | 54.7 |
| Aranese devtest | 0.1 | 70.3 | 70.2 |
| Asturian dev | 6.0 | 4.4 | 10.0 |
| Asturian devtest | 5.5 | 0.1 | 5.6 |

Table 2: TER scores comparing different versions (v1, v2, and v3) of translations for Aragonese, Aranese, and Asturian. v1 represents the original translations, v2 the post-edited or improved versions (depending on the language; see the main text), and v3 the standardized versions produced under the supervision of language academies.

v2, and v3.[42] The TER metric (Snover et al., 2006) quantifies the number of edits required to transform one sentence into another, referred to as the *reference*. It is calculated by dividing the total number of edits (insertions, deletions, substitutions, and shifts) by the total number of words in the reference. TER is a widely accepted metric in machine translation evaluation, providing insight into how close a system's output is to a human-generated reference. Lower TER scores correspond to higher translation quality. In this study, we use TER to compare different versions of the datasets for each language, offering an estimate of the extent of changes introduced by the reviewers during the revision process. In our case, higher TER scores correspond to a larger number of edits.

Version v1 refers to the original translations. As already mentioned, for Aragonese and Aranese, the translations in v1 were initially produced by Apertium, whereas v2 represents the post-edited translations of v1. For Asturian, the v1 corresponds to the version already included in the FLORES+ dataset, and v2 incorporates some normative and stylistic corrections to improve the quality of the translations. Finally, v3 refers in all cases to a version that was further adapted to conform to the linguistic standards promoted by the respective language academies. The revisions carried out under the direct supervision of these academies were designed to ensure the highest quality and compliance with current linguistic norms.

The data reveals that in most cases, the second round of quality checks leading to v3 introduced significant differences compared to v2. This is par-

ticularly notable for Aragonese and Aranese, where the adaptation to standard forms resulted in substantial changes, as reflected in the TER scores. Asturian is an exception, especially in the devtest set, where the differences between v2 and v3 were minimal, suggesting that the initial revisions in v2 already addressed most of the necessary improvements. This underscores the importance of multiple quality review stages, as the adaptation to standardized forms can introduce notable refinements to the translations.

## 4 Validation of the FLORES+ datasets

The dev sets of FLORES+ for Asturian, Aragonese, and Aranese were distributed and utilized in an WMT24 shared task[43] which focused on translating from Spanish into low-resource languages of Spain (Sánchez-Martínez et al., 2024). These resources proved valuable for the participants as validation sets for training neural MT systems. The corresponding devtest subsets were then used by the organizers of the shared task to evaluate the submitted systems. To ensure fairness in the evaluation, the devtest sets were withheld until the end of the submission period. Several submitted systems outperformed the respective Apertium-based baselines, achieving higher performance scores based on the BLEU, chrF++ and TER metrics.

The Valencian FLORES+ dataset was utilized in the development of a Spanish–Valencian neural MT system. One experiment involved training models with varying proportions of training data in the Valencian and Central Catalan dialects. As expected, results indicated that the higher the proportion of Valencian data in the training set, the higher the BLEU score achieved on the Valencian FLORES+ dataset and the greater the presence of Valencian dialectal forms in the output.

As a tentative indicator of the baseline performance expected on the newly created devtest sets, Table 3 presents standard automatic evaluation metrics for Apertium-based systems (Forcada et al., 2011) translating the Spanish side of the devtest. The results suggest that the use of MT (specifically the Apertium system) during the initial phase of corpus creation likely contributed to the significantly higher evaluation scores for the Aragonese and Aranese Apertium-based systems compared to Asturian. This points to a potential bias in the dataset

---

[42]SacreBLEU (Post, 2018) TER signature: nrefs:1 | case:lc | tok:tercom | norm:no | punct:yes | asian:no | version:2.0.0

[43]https://www2.statmt.org/wmt24/romance-task.html

| Apertium system | BLEU | chrF++ | TER |
|---|---|---|---|
| Spanish–Aragonese | 61.1 | 79.3 | 27.2 |
| Spanish–Aranese | 28.8 | 49.4 | 72.3 |
| Spanish–Asturian | 17.0 | 50.8 | 80.4 |

Table 3: Automatic evaluation scores of the Apertium-based systems on the devtest FLORES+ data. The Apertium data correspond to the following releases: `spa-arg`, 0.6.0; `oc-es`, 1.0.8; `spa-ast`, 1.1.1.

favoring rule-based MT approaches over non-rule-based systems, and even human translations, when evaluating the devtest translations from Spanish into these languages. Nonetheless, results from the WMT24 shared task (Sánchez-Martínez et al., 2024) show that neural systems can still outperform Apertium, even for Aragonese and Aranese. This suggests that the use of MT in the early stages of creating the Aragonese and Aranese FLORES+ datasets did not significantly compromise the utility of the data.

## 5 Concluding remarks

In this paper, we have detailed the development of the FLORES+ datasets for four low-resource Romance languages spoken in Spain: Aragonese, Aranese, Asturian, and the Valencian variant of Catalan. Each dataset was meticulously curated through a two-step manual review process involving native speakers and professionals from the respective language academies. The creation of these datasets is particularly significant given the ongoing efforts to revitalize and promote these languages. Our work also highlights several key challenges in resource development for low-resource languages, such as the scarcity of expert translators, or the constraints of limited time and funding for dataset production.

## Acknowledgements

## References

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Jeff Wang, and NLLB Team. 2024. Scaling neural machine translation to 200 languages. *Nature*, 630(8018):841–846.

Depto. de Cultura y Política Lingüística. 2023. *VII Encuesta Sociolingüística 2021. Comunidad Autónoma de Euskadi. Informe Resumen*. Departamento de Cultura y Política Lingüística, Gobierno Vasco, Viceconsejería de Política Lingüística, Donostia-San Sebastián.

Antonio Eito, José Ángel Iranzo, and Chaime Marcuello. 2024. Hablando aragonés: análisis de la encuesta de uso de la lengua aragonesa en la zona de uso predominante y su evolución. Unpublished.

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online. Association for Computational Linguistics.

Mikel L. Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M. Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine Translation*, 25(2):127–144.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Pan-iberian language archival resource.

Generalitat de Catalunya. 2007. *El catalá, llengua d'Europa*. Generalitat de Catalunya, Departament de la Vicepresidència, Secretaria de Política Lingüística, Barcelona.

Generalitat de Catalunya. 2019. *Els usos lingüístics de la població de l'Aran: Principals resultats de l'Enquesta d'usos lingüístics de la població. 2018*. Generalitat de Catalunya, Barcelona.

Generalitat Valenciana. 2021. *Conocimiento y uso social del valenciano: síntesis de resultados*. Generalitat Valenciana, Conselleria d'Educació, Cultura i Esport, València.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Alicia Hernández. 2022. Dónde se habla el bable y por qué dicen que es un idioma en peligro. https://www.bbc.com/mundo/noticias-internacional-59547573.

Instituto Nacional de Estadística. 2022. Encuesta de características esenciales de la población y viviendas. Año 2021. Datos definitivos. Nota de prensa.

Francisco J. Llera Ramo. 2018. *III Estudio Sociolingüístico de Asturias 2017: Avance de resultados*. Academia de la Llingua Asturiana, Uvieu. Estaya Sociollingüística, colección 7.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzman. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

Observatorio da Lingua Galega. 2007. *Situación da lingua galega na sociedade: Observación no ámbito da cidadanía: resumo executivo*. Observatorio da Lingua Galega, Xunta de Galicia, Santiago de Compostela.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Surangika Ranathunga, En-Shiun Annie Lee, Marjana Prifti Skenduli, Ravi Shekhar, Mehreen Alam, and Rishemjit Kaur. 2023. Neural machine translation for low-resource languages: A survey. *ACM Comput. Surv.*, 55(11).

Anchel Reyes, Chabier Gimeno, Miguel Montañés, Natxo Sorolla, Pep Espluga, and Juan Pablo Martínez. 2017. *L'aragonés y lo catalán en l'actualidat: analisi d'o censo de población y viviendas de 2011*. Seminario Aragonés de Sociolingüística, Asociación Aragonesa de Sociolochía, Universidad de Zaragoza. Primera parte, febrero 2017.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Felipe Sánchez-Martínez, Juan Antonio Pérez-Ortiz, Aarón Galiano-Jiménez, and Antoni Oliver. 2024. Findings of the WMT 2024 Shared Task on Translating into Low-Resource Languages of Spain: Blending rule-based and neural systems. In *Proceedings of the Ninth Conference on Machine Translation (WMT24)*, Miami, Florida, USA. Association for Computational Linguistics.

# The Bangla/Bengali Seed Dataset Submission to the WMT24 Open Language Data Initiative Shared Task

**Firoz Ahmed\***
University of Florida
firozahmed@ufl.edu

**Nitin Venkateswaran\***
University of Florida
venkateswaran.n@ufl.edu

**Sarah Moeller**
University of Florida
smoeller@ufl.edu

## Abstract

We contribute a seed dataset for the Bangla/Bengali language as part of the WMT24 Open Language Data Initiative shared task. We validate the quality of the dataset against a mined and automatically aligned dataset (NLLBv1) and two other existing datasets of crowdsourced manual translations. The validation is performed by investigating the performance of state-of-the-art translation models fine-tuned on the different datasets after controlling for training set size. Machine translation models fine-tuned on our dataset outperform models tuned on the other datasets in both translation directions (English-Bangla and Bangla-English). These results confirm the quality of our dataset. We hope our dataset will support machine translation for the Bangla/Bengali community and related low-resource languages.

## 1 Introduction

The Indian sub-continent is an area of rich linguistic diversity (Saxena and Borin, 2006; Hock and Bashir, 2016), and it is not uncommon for a language in this region to have both millions of speakers and insufficient resources for NLP development. Bangla/Bengali [ben] is one such language, ranked the 7th most spoken language in the world in the latest Ethnologue list of 200 most spoken languages (Eberhard and Fennig, 2024), and classified in the taxonomy of Joshi et al. (2020) as a Rising Star, "let down by insufficient efforts in labeled data collection" despite a "strong web presence and thriving online community". This classification contrasts squarely with that of languages such as Standard German, a "winner" in the Joshi et al. taxonomy because it has heavy investments in resources and technology, despite a ranking of 12 in the Ethnologue 200, below Bangla/Bengali.

The relative lack of data resources poses a challenge for neural machine translation (MT) efforts in Bangla/Bengali. While creating large-scale datasets of parallel sentences would be the next step towards improving Bangla/Bengali MT, efforts to create these types of resources have only recently been made (Hasan et al., 2020; Siripragada et al., 2020; Ramesh et al., 2022). Such efforts often must use automated methods to crawl and align the texts between language pairs, with manual checks and reviews being prohibitively expensive. There has been little work comparing larger datasets with smaller, professionally translated and manually curated datasets to investigate how the differences between these two types of dataset could impact the quality of machine translation.

This paper describes the results of one such manual effort, creating translation pairs between English and Bangla/Bengali for a smaller dataset and verifying the quality of those translations. Maillard et al. (2023) shows the sizeable impact of these smaller datasets on MT quality via bilingual and multilingual translation experiments, with the high quality manually translated datasets outperforming even a back-translation data augmentation approach with larger train set sizes. Continuing this line of reasoning, we hypothesize that models trained using a smaller but professionally translated dataset of Bangla/Bengali would perform better than models trained on larger, automatically mined and aligned parallel texts with little to no human intervention or review, once training sizes are controlled for. To this end, we created a smaller dataset of manual translations to test our hypotheses, and explored different training set sizes from larger datasets to check their equivalencies against our smaller dataset.

Our main contributions are as follows:

1. We contribute to the open datasets of the Open Language Data Initiative (OLDI) and produce a seed dataset for Bangla/Bengali by translat-

---

\* These authors contributed equally

ing the English seed dataset.

2. We carry out fine-tuning translation experiments to show that models tuned on our smaller, manually translated dataset outperform, or are on par with, models tuned on samples of comparable sizes from another dataset that has been automatically mined and aligned (NLLBv1).

3. We compare our dataset with other manually translated datasets for Bangla/Bengali available via OPUS (Tiedemann, 2012), and show that our dataset outperforms both corresponding and larger samples from these datasets (1.5x, 2x larger than our dataset) for a majority of pre-trained models in our experiments.

## 2 Related Work

Machine translation (MT) efforts in Bangla/Bengali currently rely on creative methods such as data augmentation and multilingual transfer to approach state-of-the-art MT. For example, Mondal et al. (2024), a recent work, uses back-translation to augment training data for English-Bengali transformer-based MT. Laskar et al. (2022a) augment data for English-Bengali MT using an SMT-based phrase-pair injection approach (Sen et al., 2021), and transliterate English texts into Bengali script as a transfer mechanism to share subword-level information between source and target sentences. Jasim et al. (2020) use a partial back-translation method by translating only selected phrases to the source language, achieving competitive results for Bengali MT on the WAT2018 (Nakazawa et al., 2018) test set. Laskar et al. (2022b) investigate knowledge transfer among Indic languages for neural MT, including Bengali, by transliterating all Indic languages into English script to share subword information during training. Bala Das et al. (2023) build a transformer-based multilingual neural MT system for 15 Indic language pairs, including Bengali, and English with shared encoder-decoders and transliteration schemes for related languages. Gala et al. (2023) build a multilingual NMT system for 22 Indic languages including Bengali.

Efforts to create large-scale datasets of parallel sentences for Bangla/Bengali have only recently been made. For example, Hasan et al. (2020) create a dataset of 2.75 million sentence pairs for machine translation, using an automated sentence segmenta-tion toolkit and an ensemble of aligners for bitext alignment. Siripragada et al. (2020) collect parallel corpora across 10 Indian languages, including Bengali, by crawling two Government of India websites and applying document and sentence level alignment methods, producing 126.7K parallel texts for Bengali-English. Ramesh et al. (2022) create a dataset of parallel texts for 11 Indian languages, including 8.6 million parallel texts for Bengali-English, by crawling news and education/MOOC websites such as Coursera and passing the data through automated pipelines. Schwenk et al. (2021) mine billions of parallel texts from the web for multiple languages, which include approximately 10 million sentence pairs for Bengali-English aligned using LASER embeddings (Artetxe and Schwenk, 2019). As part of the No Language Left Behind project, NLLB Team et al. (2022) mine 62 million sentence pairs for Bengali-English using LASER3 embeddings (Heffernan et al., 2022); the NLLBv1 dataset created from the project is the largest known dataset of parallel texts for this language pair to date.

## 3 Meet The Data

Here we describe the language, the data collection process, and the format of the dataset.

### 3.1 Language description

Bangla/Bengali (ISO-639-3:ben, glottocode:beng1280, ISO-15924:Beng), an Indo-Aryan language, is the official and national language of Bangladesh and an official language of the state of West Bengal and other states in India. The language is commonly referred to as Bengali within the Indian states, and as Bangla in the nation of Bangladesh. The standardized dialects spoken in these two regions differ mainly in the morpho-phonological space. For example, Bangla has separate objective and genitive case markings for nouns and pronouns while Bengali has syncretized forms for these. The mid-back-rounded vowel (/ɔ/) is more common word-finally in Bangla than in Bengali. Despite "numerous small differences", both dialects have been called "indisputably the same language" (David, 2015). We refer to the language as Bangla for the rest of this paper, since the translations were produced in this dialect.

The script system of Bangla is similar to that of other South Asian languages in being an abugida

| অন্ধকার | প্রায় | হলো | | হলো। |
|---------|-------|-----|---|------|
| ɔndhɔkar | prae | ho-l-o | | holo |
| dark | almost | become-PST-3.NHON | | REDUP |

'It was just about **to get dark**. (Lit. The dark almost **happened happened**.)'

| আমার | গাড়ি | পরিষ্কার | করতে | হবে। |
|------|------|---------|------|------|
| amar | gaṛi | poriṣkar | kor-te | hɔ-b-e |
| 1SG.GEN | car | clean | do-IPFP | become-FUT-3.NHON |

'I need **to clean** my car.'

Figure 1: Glossed examples in Bangla script using reduplication and conjunct verbs; examples from (David, 2015)

system organized by syllables with two forms for each vowel viz., the independent and diacritic forms, and with a system of conjunct characters for complex consonant segments. The script (Fig. 1) is represented in Unicode with range 0980-09FF[1], which we use for our translations.

Bangla has certain features which make translation between Bangla and English a challenging task. These include rich morphological systems of inflection, derivation, and reduplication, a rich case system, a system of light verbs and conjunct verbs, and a system of noun classifiers. All these features are less prevalent in English. While Bangla is an SOV type language, scrambling of constituents within and across clauses for the purpose of altering information structure is common. This can pose a challenge for neural translation systems (Belinkov and Bisk, 2017). An in-depth description of the features of Bangla/Bengali can be found in David (2015).

## 3.2 Data Collection and Translation

The Bangla sentences in our dataset were manually translated from the English sentences in the Seed dataset v2.0 (Maillard et al., 2023) maintained by the Open Language Data Initiative. Details about the sourcing and composition of the dataset are described in Maillard et al. (2023). One native speaker of Bangla, an author of this paper fluent in English with graduate-level linguistic training and experience in professional translation from English to Bangla, translated all 6,193 sentences in the dataset. The Avro keyboard for Windows[2] with Unicode support was used to generate the Bangla translations. The translation guidelines[3] supplied

by the Open Language Data Initiative were followed during the translation process.

## 3.3 Data Format

The Bangla translations are stored as a text file with a single line per translation, containing sentences in the same order as in the English seed dataset[4]. We follow the dataset formatting guidelines provided by the Open Language Data Initiative[5].

## 4 Experimental Validation

We compare our Bangla translations of the Seed dataset with the following three datasets. All datasets were downloaded from OPUS.

**NLLBv1** (NLLB Team et al., 2022). This is the largest available collection of automatically aligned Bangla-English sentence pairs with a wide range of text domains.

**Joshua-IPC** (Post et al., 2012). This is a dataset of parallel sentences for six Indic languages including Bangla. It was crowd-sourced by the authors via Amazon Mechanical Turk for translation experiments using the Joshua statistical MT system (Weese et al., 2011). Sentences for the Indic languages were extracted from the top 100 viewed Wikipedia pages for the language, and four English translations sourced for each sentence.

**TED2020** (Reimers and Gurevych, 2020). This is a dataset of crawled and aligned subtitles of TED Talks for the month of July 2020 across multiple languages, with subtitling carried out by a global community of volunteer translators[6]. We downloaded all 10,519 sentence pairs for Bangla-English, with translations in the English to Bangla direction.

## 4.1 Controlling for training set sizes

To facilitate comparisons with our translations, we control for training set sizes by sampling 1K, 3K and 6K sentence pairs from all datasets, similar to the approach used in Maillard et al. (2023). We select these training sizes to test whether models trained on smaller samples of our translations outperform models trained on samples of corresponding sizes from the other datasets. In addition, we sample 9K and 12K sentence pairs from the NLLBv1 and Joshua-IPC datasets, and 9K and the full 10,519 sentence pairs from TED2020. We

---

[1] https://www.unicode.org/charts/PDF/U0980.pdf
[2] https://www.omicronlab.com/avro-keyboard.html
[3] Translation guidelines:https://oldi.org/guidelines

[4] github:openlanguagedata/seed/blob/main/seed/eng_Latn
[5] Formatting guidelines: https://oldi.org/guidelines
[6] https://www.ted.com/participate/translate

compare results from these larger sizes with results trained on 6K sentence pairs from our translations.

We sample using five different seeds for each training size where possible, averaging results across all seeds instead of relying on results from a single sample per training size.

### 4.2 Translation models

We fine-tune translation models on existing pre-trained multilingual models and one pre-trained monolingual model in both directions (Bangla-English and English-Bangla). Only the sampled sentence pairs described in section 4.1 are used to fine-tune the models and no additional data is used. All models are fine-tuned using the HuggingFace `transformers` library (Wolf et al., 2020) and use a linear learning rate schedule with an initial rate of 1e-6, with warmup. The following pre-trained models are used.

**NLLB-200**. The state-of-the-art NLLB-200 model (NLLB Team et al., 2022) is pre-trained on 200 languages, including Bangla and English. The `nllb-200-1.3B` dense model with 1.3 billion parameters is used for our fine-tuning.

**mBART50** (Tang et al., 2020). This is a multilingual seq2seq model primarily intended for the task of machine translation through multilingual fine-tuning. This model is pre-trained on 50 languages, including Bangla and English. We use the `mbart-large-50` model.

**mT5** (Xue et al., 2021). We experiment with the multilingual variant of the text-to-text transformer pretrained on a Common Crawl based dataset containing 101 languages, including Bangla and English. The `mt5-large` model is used.

**BanglaT5** (Bhattacharjee et al., 2023). To investigate the impact of a pre-trained monolingual model on translation quality, we fine-tune the BanglaT5 model pretrained on the Bangla2B+ corpus (Bhattacharjee et al., 2022). We select this model based on its open-source availability and relatively large pre-training corpus size of 27.5GB. We note that future work could include experiments with other open-source pre-trained monolingual models as and when they become available.

### 4.3 Evaluation Metrics

We evaluate all models with the Bangla/Bengali and English datasets from the FLORES+ evaluation benchmark for multilingual machine translation (NLLB Team et al., 2022), maintained by the Open Language Data Initiative. We use the development set for model tuning and early stopping, and the test set to report translation metrics.

We report the chrF++ scores (Popović, 2017) calculated using the `sacrebleu` toolkit (Post, 2018), since the chrF-based score is known to correlate well with human rankings especially for morphologically rich languages like Bangla, outperforming BLEU (Popović, 2015). BLEU has also been shown to be less useful for morphologically complex languages, with language-specific customizations showing better correlations with human rankings (Chauhan et al., 2021; Bouamor et al., 2014).

## 5 Experiment Results

Comparing the results on our dataset against the others, we can confirm that our manual translations yielded high quality parallel sentences between English and Bangla. Tables 1 and 2 in the Appendix show the fine-tuned chrF++ scores in the English-Bangla and Bangla-English directions. Here we discuss the results in Figures 2, 3, 4 and 5, displayed below.

Displaying the results in the English to Bangla direction, Figure 2 shows that models fine-tuned on our translations outperform models tuned on samples of corresponding sizes from the other datasets. This demonstrates the high quality of our translations. In the case of NLLB-200 our translations are on par with the NLLBv1 samples. The results, interestingly, also hold for the 1K and 3K sample sizes showing that smaller samples of our translations are also effective. Given that it is likely the NLLBv1 dataset was used to pre-train the NLLB-200 models, it is not surprising that the NLLB-200 models fine-tuned on the NLLBv1 samples in our experiments show good performance despite the automatic sentence alignment process.

In Figure 3, it can be seen that in the Bangla to English direction, models tuned on our translations outperform other models across all corresponding sample sizes, except for NLLB-200 where models tuned on the TED2020 dataset show the best performance. The average performance gap for the 6K sample size between our translations and the other datasets is 5.34 points for the mBART50 model and 4.52 points for the mT5 model. These wide margins show that the quality of the dataset used for fine-tuning can make a sizeable difference even for multilingual models with large pre-training corpora in English.

Considering the test of our translations against

Figure 2: Averaged fine-tuned English to Bangla chrF++ scores on the FLORES+ test set for the 1K, 3K and 6K training set sizes. Models tuned on the Bangla seed dataset (red) outperform, or are on par with, models tuned on the other datasets across pre-trained model types and training sizes. Scores are averaged across five random samples per training set size and dataset.



Figure 3: Averaged fine-tuned Bangla to English chrF++ scores on the FLORES+ test set for the 1K, 3K and 6K training set sizes. Models tuned on the Bangla seed dataset (red) outperform models tuned on the other datasets across pre-trained model types and training set sizes, except for the NLLB-200 models. Scores are averaged across five random samples per training set size and dataset.

Figure 4: Averaged fine-tuned English to Bangla chrF++ scores on the FLORES+ test set for the 6K, 9K and 12K training set sizes; the complete TED2020 dataset is used in the 12K case. Models tuned on the Bangla seed dataset (red) outperform models tuned on other datasets of larger training sizes (9K, 12K) across pre-trained model types, except for the NLLB-200 models tuned on the NLLBv1 data. Scores are averaged across five random samples per training set size and dataset.



Figure 5: Averaged fine-tuned Bangla to English chrF++ scores on the FLORES+ test set for the 6K, 9K and 12K training set sizes; the complete TED2020 dataset is used in the 12K case. Models tuned on the Bangla seed dataset (red) outperform models tuned on other datasets of larger training sizes (9K, 12K) across pre-trained model types, except for the NLLB-200 models. Scores are averaged across five random samples per training set size and dataset.

samples of larger sizes from the other datasets, i.e 9K and 12K sentence pairs. The results in Figures 4 and 5 show that our translations outperform these larger training samples across all pre-trained model types and datasets except for the NLLB-200 models. The NLLB-200 models tuned on larger samples of the NLLBv1 dataset in the English to Bangla direction scored better than our translations. This is as expected, given the possible overlap between the NLLBv1 datasets used to pre-train and fine-tune the models in our experiments. The fact that models tuned on our translations scored better than models tuned on larger samples from the other datasets is another demonstration of the higher quality of our dataset.

# 6 Conclusion

We have created a high quality dataset of Bangla-English seed translations to contribute to the Open Language Data Initiative, paving the way for more translations between Bangla and other languages, including low-resource ones, that are supported by the initiative. We have demonstrated the high quality of our translated dataset by comparing it with a larger dataset that was mined and automatically aligned, as well as with two datasets of crowdsourced and reviewed translations. The models tuned on our dataset outperform models tuned on the other datasets after controlling for training set size. We hope that our dataset will support ongoing research in machine translation for the Bangla/Bengali community and other low-resource languages.

## Acknowledgements

We would like to thank the organizers of the Open Language Data Initiative Shared Task for their support throughout the process.

## Ethics Statement

`Licensing` We release the dataset under a CC-BY-SA-4.0 license.

# References

Mikel Artetxe and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the Association for Computational Linguistics*, 7:597–610.

Sudhansu Bala Das, Atharv Biradar, Tapas Kumar Mishra, and Bidyut Kr. Patra. 2023. Improving multilingual neural machine translation system for indic languages. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 22(6).

Yonatan Belinkov and Yonatan Bisk. 2017. Synthetic and natural noise both break neural machine translation. *ArXiv*, abs/1711.02173.

Abhik Bhattacharjee, Tahmid Hasan, Wasi Ahmad, Kazi Samin Mubasshir, Md Saiful Islam, Anindya Iqbal, M. Sohel Rahman, and Rifat Shahriyar. 2022. BanglaBERT: Language model pretraining and benchmarks for low-resource language understanding evaluation in Bangla. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1318–1327, Seattle, United States. Association for Computational Linguistics.

Abhik Bhattacharjee, Tahmid Hasan, Wasi Uddin Ahmad, and Rifat Shahriyar. 2023. BanglaNLG and BanglaT5: Benchmarks and resources for evaluating low-resource natural language generation in Bangla. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 726–735, Dubrovnik, Croatia. Association for Computational Linguistics.

Houda Bouamor, Hanan Alshikhabobakr, Behrang Mohit, and Kemal Oflazer. 2014. A human judgement corpus and a metric for Arabic MT evaluation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 207–213, Doha, Qatar. Association for Computational Linguistics.

Shweta Singh Chauhan, Philemon Daniel, Archita Mishra, and Abhay Kumar. 2021. Adableu: A modified bleu score for morphologically rich languages. *IETE Journal of Research*, 69:5112 – 5123.

Anne Boyle David. 2015. *Descriptive Grammar of Bangla*. De Gruyter Mouton, Berlin, München, Boston.

Gary F. Simons Eberhard, David M. and Charles D. Fennig, editors. 2024. *Ethnologue: Languages of the World*, twenty-seventh edition. SIL International, Dallas, TX, USA.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Tahmid Hasan, Abhik Bhattacharjee, Kazi Samin, Masum Hasan, Madhusudan Basak, M. Sohel Rahman, and Rifat Shahriyar. 2020. Not low-resource anymore: Aligner ensembling, batch filtering, and new datasets for Bengali-English machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,

pages 2612–2623, Online. Association for Computational Linguistics.

Kevin Heffernan, Onur Çelebi, and Holger Schwenk. 2022. Bitext mining using distilled sentence representations for low-resource languages. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2101–2112, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Hans Henrich Hock and Elena Bashir, editors. 2016. *The Languages and Linguistics of South Asia*. De Gruyter Mouton, Berlin, Boston.

Binu Jasim, Vinay Namboodiri, and C V Jawahar. 2020. PhraseOut: A code mixed data augmentation method for MultilingualNeural machine tranlsation. In *Proceedings of the 17th International Conference on Natural Language Processing (ICON)*, pages 470–474, Indian Institute of Technology Patna, Patna, India. NLP Association of India (NLPAI).

Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.

Sahinur Rahman Laskar, Pankaj Dadure, Riyanka Manna, Partha Pakray, and Sivaji Bandyopadhyay. 2022a. English to Bengali multimodal neural machine translation using transliteration-based phrase pairs augmentation. In *Proceedings of the 9th Workshop on Asian Translation*, pages 111–116, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

Sahinur Rahman Laskar, Riyanka Manna, Partha Pakray, and Sivaji Bandyopadhyay. 2022b. Investigation of multilingual neural machine translation for indian languages. In *Workshop on Asian Translation*.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzmán. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

Subrota Kumar Mondal, Chengwei Wang, Yijun Chen, Yuning Cheng, Yanbo Huang, Hong-Ning Dai, and H. M. Dipu Kabir. 2024. Enhancement of english-bengali machine translation leveraging back-translation. *Applied Sciences*, 14(15).

Toshiaki Nakazawa, Katsuhito Sudoh, Shohei Higashiyama, Chenchen Ding, Raj Dabre, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, and Sadao Kurohashi. 2018. Overview of the 5th workshop on Asian translation. In *Proceedings of the 32nd*

*Pacific Asia Conference on Language, Information and Computation: 5th Workshop on Asian Translation: 5th Workshop on Asian Translation*, Hong Kong. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Matt Post, Chris Callison-Burch, and Miles Osborne. 2012. Constructing parallel corpora for six Indian languages via crowdsourcing. In *Proceedings of the Seventh Workshop on Statistical Machine Translation*, pages 401–409, Montréal, Canada. Association for Computational Linguistics.

Gowtham Ramesh, Sumanth Doddapaneni, Aravinth Bheemaraj, Mayank Jobanputra, Raghavan AK, Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Mahalakshmi J, Divyanshu Kakwani, Navneet Kumar, Aswin Pradeep, Srihari Nagaraj, Kumar Deepak, Vivek Raghavan, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh Shantadevi Khapra. 2022. Samanantar: The largest publicly available parallel corpora collection for 11 Indic languages. *Transactions of the Association for Computational Linguistics*, 10:145–162.

Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4512–4525, Online. Association for Computational Linguistics.

Anju Saxena and Lars Borin, editors. 2006. *Lesser-Known Languages of South Asia*. De Gruyter Mouton, Berlin, New York.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021. CCMatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500, Online. Association for Computational Linguistics.

Sukanta Sen, Mohammed Hasanuzzaman, Asif Ekbal, Pushpak Bhattacharyya, and Andy Way. 2021. Neural machine translation of low-resource languages using smt phrase pair injection. *Natural Language Engineering*, 27(3):271–292.

Shashank Siripragada, Jerin Philip, Vinay P. Namboodiri, and C V Jawahar. 2020. A multilingual parallel corpora collection effort for Indian languages. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3743–3751, Marseille, France. European Language Resources Association.

Y. Tang, C. Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. Multilingual translation with extensible multilingual pretraining and finetuning. *ArXiv*, abs/2008.00401.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Jonathan Weese, Juri Ganitkevitch, Chris Callison-Burch, Matt Post, and Adam Lopez. 2011. Joshua 3.0: Syntax-based machine translation with the thrax grammar extractor. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 478–484, Edinburgh, Scotland. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

## A    Appendix

### A.1    Fine-tuned English to Bangla chrF++ scores

Table 1 shows the fine-tuned English to Bangla chrF++ scores across all model types, datasets, and training set sizes.

### A.2    Fine-tuned Bangla to English chrF++ scores

Table 2 shows the fine-tuned Bangla to English chrF++ scores across all model types, datasets, and training set sizes.

### A.3    `sacrebleu` version string

The `sacrebleu` version string is provided below for reproducibility:
`nrefs:1|case:mixed|eff:yes|nc:6|nw:2|space:no|version:2.4.2`

| mBART50 | | | | | |
|---|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 16.6 | 20.12 | 22.04 | 23.6 | 24.72 |
| Joshua-IPC | 14.22 | 16.58 | 17.48 | 18.42 | 18.44 |
| TED2020 | 17.42 | 20.8 | 22.48 | 23.82 | 24.04 |
| Bangla seed | **23.12** | **25.14** | **27.82** | – | – |

| mT5 | | | | | |
|---|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 12.12 | 13.62 | 14.12 | 13.66 | 14 |
| Joshua-IPC | 11.2 | 12.46 | 12.34 | 12.32 | 12.28 |
| TED2020 | 15.74 | 15.22 | 15.1 | 15.28 | 14.9 |
| Bangla seed | **18.2** | **17.88** | **18.04** | – | – |

| NLLB-200 | | | | | |
|---|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 34.5 | 36.84 | **38.04** | **39.3** | **39.6** |
| Joshua-IPC | 28.32 | 28.76 | 28.92 | 29.3 | 29 |
| TED2020 | 31.6 | 33.7 | 34.26 | 34.76 | 34.64 |
| Bangla seed | **35.12** | **36.9** | 37.14 | – | – |

| BanglaT5 | | | | | |
|---|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 26.76 | 27.36 | 27.36 | 27.52 | 27.32 |
| Joshua-IPC | 24.98 | 25 | 25.06 | 25.14 | 25.08 |
| TED2020 | 29.64 | 29.7 | 29.76 | 29.94 | 30.2 |
| Bangla seed | **34.4** | **33.92** | **33.64** | – | – |

Table 1: Average fine-tuned English to Bangla chrF++ scores on the FLORES+ test set. Scores are averaged across five random samples per training set size and dataset

| mBART50 | | | | |
|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 24.54 | 28.62 | 29.94 | 31.14 | 32.2 |
| Joshua-IPC | 24.96 | 27.68 | 29.18 | 30.08 | 29.76 |
| TED2020 | 25.46 | 27.5 | 29.12 | 29.58 | 29.54 |
| Bangla seed | **30.62** | **33.64** | **34.76** | – | – |

| mT5 | | | | |
|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 21.42 | 25.68 | 26.6 | 26.66 | 26.84 |
| Joshua-IPC | 24.06 | 26.46 | 26.74 | 26.74 | 26.78 |
| TED2020 | 23.88 | 25.84 | 26.46 | 26.76 | 26.48 |
| Bangla seed | **30.46** | **31.66** | **31.12** | – | – |

| NLLB-200 | | | | |
|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 48.16 | 48.22 | 48.3 | 48.08 | 48.14 |
| Joshua-IPC | 44.36 | 44.42 | 44.54 | 44.48 | 44.14 |
| TED2020 | **49.88** | **49.74** | **49.94** | **49.88** | **49.7** |
| Bangla seed | 48 | 47.78 | 47.84 | – | – |

| BanglaT5 | | | | |
|---|---|---|---|---|
| | 1K | 3K | 6K | 9K | 12K |
| NLLBv1 | 23.2 | 22.5 | 23.08 | 23.4 | 23.2 |
| Joshua-IPC | 21.74 | 21.82 | 21.66 | 21.98 | 21.68 |
| TED2020 | 24.38 | 24.28 | 24.66 | 24.44 | 24.3 |
| Bangla seed | **28.48** | **28.1** | **27.92** | – | – |

Table 2: Average fine-tuned Bangla to English chrF++ scores on the FLORES+ test set. Scores are averaged across five random samples per training set size and dataset

# A high-quality Seed dataset for Italian machine translation

**Edoardo Ferrante**
Council for Ligurian Linguistic Heritage
`info@conseggio-ligure.org`

## Abstract

This paper describes the submission of a high-quality translation of the OLDI Seed dataset into Italian for the WMT 2024 Open Language Data Initiative shared task.

The base of this submission is a previous version of an Italian OLDI Seed dataset released by Haberland et al. (2024) via machine translation and partial post-editing. This data was subsequently reviewed in its entirety by two native speakers of Italian, who carried out extensive post-editing with particular attention to the idiomatic translation of named entities.

## 1 Language overview

This paper presents an Italian version of the OLDI Seed dataset (Maillard et al., 2023; NLLB Team et al., 2024).

Italian is a Romance language, recognised as an official language of the Italian Republic, the Republic of San Marino and the Canton of Ticino in Switzerland (Maiden, 2014). Modern Italian fundamentally represents cultured Florentine, as first attested by 14th century authors (Dante, Petrarch and Boccaccio) and later scholars (Coletti, 2022). Although it is a variety of Tuscan, standard Italian is purged of the more typical features of Tuscan, at a phonetic level represented above all by the so-called *gorgia* (i.e. the fricative pronunciation of certain occlusive consonants in intervocalic position) (Marotta, 2008).

The presence of a curated Italian version in the Seed dataset is of great importance for the regional languages of the Italian peninsula, six of which are already represented in the same dataset.[1] The creation of an Italian version enables the training of machine translation models for these languages to and from Italian, a direction which is more culturally relevant than English-centric MT, as the vast majority of speakers of such languages (or prospect learners) are also native Italian speakers (Haberland et al., 2024; Ramponi, 2024).

## 2 Data creation

The original source of the data was an initial Italian version of the Seed dataset released by Haberland et al. (2024). The authors created it by machine translating the original English version with an OpusMT bilingual English-Italian model (Tiedemann and Thottingal, 2020), combined with partial post-editing. Through personal correspondence with the authors we learned that their post-editing, which only affected a small percentage of the overall data, involved two steps:

1. A check of the length ratios of Italian and English sentences, followed by manual checking and post-editing of sentence pairs with outlier length ratios.

2. A spellchecker run using Hunspell (Ooms et al., 2017), followed by manual checking and post-editing of sentence pairs where spelling mistakes were found.

The submission described in this paper constitutes a further refinement of the dataset of Haberland et al. (2024), in order to bring it to a level that could be seen as comparable to that of translations produced by highly proficient bilingual individuals.

This project involved the participation of two annotators, henceforth A1 and A2, both native speakers of Italian with a university level of education. The refinement process followed these steps:

1. A manual, sequential review of the entire dataset by A1, followed by post-editing where necessary.

2. Following Haberland et al. (2024), a targeted review of sentence pairs with outlier length ratios, followed by post-editing where necessary.

---

[1] These are Friulian, Ligurian (Genoese), Lombard, Sicilian, Sardinian, Venetian (Maillard et al., 2023; NLLB Team et al., 2024).

3. A targeted review of sentences involving specific subsets of the corpus which were found to have a high incidence of mistranslated strings: date and time expressions, sentences about mathematics and sentences about the history of cinema.

4. A final targeted review of sentence pairs which were found to be of low-quality using a series of Quality Estimations approaches using LLMs, as described in Zhao et al. (2024).

Apart from the first item above, which was carried out by annotator A1 alone, the workload for all subsequent tasks was split equally between both annotators.

## 3 Experimental validation

In order to experimentally validate the quality of this Seed dataset, we replicate the baseline experiments of Haberland et al. (2024), by training an Italian-Ligurian machine translation model on a combination of the 6,193 paired Italian-Ligurian sentences from the Seed data and the same 1,520 paired Italian-Ligurian sentences from the Tatoeba project[2] used by the authors. The translation model is trained using Fairseq (Ott et al., 2019), with the exact same architecture and overall setup of Haberland et al. (2024).

| Model | FLORES |
|---|---|
| NLLB-3.3B | 13.9 |
| Haberland et al. (2024) | 14.5 |
| **Ours** | **15.0** |

Table 1: Italian-Ligurian translation performance measured with BLEU on FLORES `devtest`.

In Table 1 we compare the BLEU scores[3] obtained by three models on the FLORES (NLLB Team et al., 2024) `devtest` data. The first model, provided only for context, is the massively multilingual 3.3B version of NLLB (NLLB Team et al., 2024), which was trained on much larger amounts of data but without any direct Italian-Ligurian supervision. The second is the baseline model of Haberland et al. (2024). The final row reports the performance of the best of three training runs of our model, which is a re-training of Haberland et al.'s,

the only difference being the use of the improved Seed data.

As can be observed in the results, our model achieves a performance of 15 BLEU points on the FLORES devtest set, 1.1 points higher compared to NLLB-3.3B (NLLB Team et al., 2024) and half a point higher compared to the baseline model of Haberland et al. (2024). The relatively small degree of improvement compared to the latter baseline can be attributed to the fact that, in general, machine translation for a high-resource language pair such as English-Italian is of high quality, so that manual post-editing (especially in a formal domain such as Wikipedia text) leads to only minor changes.

This result, although numerically marginal, confirms that our post-editing of the seed data for improved idiomaticity does not hurt the downstream performance of models trained on it but does, in fact, slightly improve it.

## 4 Data samples

We provide a selection of samples whose translations proved to be particularly hard for the OpusMT bilingual English-Italian model.

## Acknowledgments

We thank Haberland et al. (2024) for the helpful discussions about post-editing and for helping us replicate their experiments.

## References

Vittorio Coletti. 2022. *Storia dell'italiano letterario*. Einaudi.

Christopher R. Haberland, Jean Maillard, and Stefano Lusito. 2024. Italian-Ligurian machine translation in its cultural context. In *Proceedings of the 3rd Annual Meeting of the Special Interest Group on Under-resourced Languages @ LREC-COLING 2024*, pages 168–176, Torino, Italia. ELRA and ICCL.

M. Maiden. 2014. *A Linguistic History of Italian*. Longman Linguistics Library. Taylor & Francis.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzmán. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

Giovanna Marotta. 2008. Lenition in tuscan italian (gorgia toscana).

---

[2] https://tatoeba.org/
[3] SacreBLEU (Post, 2018) signature `nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.4.0`.

| English (original) | Italian |
|---|---|
| He made a series of two-reel comedies, including One Week (1920), The Playhouse (1921), Cops (1922), and The Electric House (1922). | Realizzò una serie di commedie a due bobine, tra cui Una settimana (1920), Il teatro (1921), Poliziotti (1922) e La casa elettrica (1922). |
| The development of a regulatory framework concerning genetic engineering began in 1975, at Asilomar, California. | Lo sviluppo di un quadro normativo sull'ingegneria genetica è iniziato nel 1975, ad Asilomar, in California. |
| But the next major advance in the theory was made by Georg Cantor; in 1895 he published a book about his new set theory, introducing, among other things, transfinite numbers and formulating the continuum hypothesis. | Ma il successivo importante passo avanti nella teoria fu compiuto da Georg Cantor, che nel 1895 pubblicò un libro sulla sua nuova teoria degli insiemi, introducendo, tra l'altro, i numeri transfiniti e formulando l'ipotesi del continuo. |
| Aside from Steamboat Bill, Jr. (1928), Keaton's most enduring feature-length films include Our Hospitality (1923), The Navigator (1924), Sherlock Jr. (1924), Seven Chances (1925), The Cameraman (1928), and The General (1926). | Oltre a Io... e il ciclone (1928), tra i lungometraggi più duraturi di Keaton vi sono La legge dell'ospitalità (1923), Il navigatore (1924), Sherlock Jr. (1924), Le sette probabilità (1925), Il cameraman (1928) e Come vinsi la guerra (1926). |
| These chains of extensions make the natural numbers canonically embedded (identified) in the other number systems. | Queste catene di estensioni rendono i numeri naturali canonicamente immersi (identificati) negli altri sistemi numerici. |

Table 2: Dataset samples.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2024. Scaling neural machine translation to 200 languages. *Nature*, 630(8018):841–846.

Jeroen Ooms et al. 2017. Hunspell: High-performance stemmer, tokenizer, and spell checker.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–

191, Brussels, Belgium. Association for Computational Linguistics.

Alan Ramponi. 2024. Language Varieties of Italy: Technology Challenges and Opportunities. *Transactions of the Association for Computational Linguistics*, 12:19–38.

Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-MT — Building open translation services for the World. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal.

Haofei Zhao, Yilun Liu, Shimin Tao, Weibin Meng, Yimeng Chen, Xiang Geng, Chang Su, Min Zhang, and Hao Yang. 2024. From handcrafted features to llms: A brief survey for machine translation quality estimation.

# Correcting FLORES Evaluation Dataset for Four African Languages

**Idris Abdulmumin**[1*+], **Sthembiso Mkhwanazi**[2], **Mahlatse S. Mbooi**[2],
**Shamsuddeen Hassan Muhammad**[3*+], **Ibrahim Said Ahmad**[4*+], **Neo Putini**[5],
**Miehleketo Mathebula**[1], **Matimba Shingange**[1], **Tajuddeen Gwadabe**[*+], **Vukosi Marivate**[1,6]

[1]Data Science for Social Impact, University of Pretoria, [2]Council for Scientific and Industrial Research, South Africa,

[3]Imperial College, London, [4]Northeastern University, [5]University of KwaZulu-Natal, [6]Lelapa AI, [*]HausaNLP, [+]MasakhaneNLP

correspondence: idris.abdulmumin@up.ac.za

## Abstract

This paper describes the corrections made to the FLORES evaluation (dev and devtest) dataset for four African languages, namely Hausa, Northern Sotho (Sepedi), Xitsonga, and isiZulu. The original dataset, though groundbreaking in its coverage of low-resource languages, exhibited various inconsistencies and inaccuracies in the reviewed languages that could potentially hinder the integrity of the evaluation of downstream tasks in natural language processing (NLP), especially machine translation. Through a meticulous review process by native speakers, several corrections were identified and implemented, improving the overall quality and reliability of the dataset. For each language, we provide a concise summary of the errors encountered and corrected and also present some statistical analysis that measures the difference between the existing and corrected datasets. We believe that our corrections improve the linguistic accuracy and reliability of the data and, thereby, contribute to a more effective evaluation of NLP tasks involving the four African languages. Finally, we recommend that future translation efforts, particularly in low-resource languages, prioritize the active involvement of native speakers at every stage of the process to ensure linguistic accuracy and cultural relevance.

## 1 Introduction

Low-resource languages, especially from Africa, are greatly under-represented in the Natural Language Processing (NLP) landscape, and this is primarily due to the absence of sufficient resources for both training and evaluation (Adelani et al., 2022; Kreutzer et al., 2022). Various efforts have been made to create such resources and these include initiatives from organizations such as Lacuna[1] that fund new and qualitative open datasets, and communities such as Masakhane, HausaNLP,

the University of Pretoria's Data Science for Social Impact (DSFSI) Research Group, and other individual initiatives (Abdulmumin et al., 2022; Parida et al., 2023). For machine translation evaluation, the FLORES dataset (Goyal et al., 2021; NLLB Team et al., 2022) is widely accepted as a benchmark for evaluation, especially because it was the first of its kind for many languages and enables many-to-many evaluation, making it easier to evaluate say a Hausa to Sepedi translation system without pivoting through a high resource language, e.g., English. Recently, the MAFAND dataset (Adelani et al., 2022) was created, but it only allows bilingual evaluation and is limited to the news domain.

While all these resources are being developed, it is imperative to review them for validation to ensure that they meet the expected standard of accuracy and representation. A revealing work by Kreutzer et al. (2022), albeit on mostly web-crawled datasets, found that many of the datasets that are being relied upon for low-resource languages are littered with significant errors such as misalignments, incorrect translations, and other issues. The significance of evaluation datasets make them even more deserving of such reviews especially by literate native speakers that know how these languages are written and spoken. This paper, therefore, presents a comprehensive review and correction of the public FLORES evaluation datasets for four African languages: Hausa, Northern Sotho, Xitsonga and isiZulu. We also provide the corrected datasets for future evaluation tasks[2].

## 2 The FLORES Evaluation Dataset

The FLORES evaluation dataset consists of the first FLORES-101 (Goyal et al., 2021) and the subsequent more expanded FLORES-200 (NLLB Team et al., 2022) that included more languages.

---

[1]https://lacunafund.org/

[2]https://github.com/dsfsi/flores-fix-4-africa

**FLORES-101:** This was the original evaluation data and was created by translating the English dataset collected from Wikipedia, consisting of several topics and domains, into 101 mostly low-resource languages. The dataset was the first available evaluation benchmark for several low-resource languages and it enabled the evaluation of many-to-many translation systems without pivoting through another high-resource language such as English. Several quality control mechanisms were put in place to ensure that the final dataset was of acceptable quality. To determine if translations are good enough for inclusion in FLORES-101, a 20% sample of the dataset were reviewed by language-specific evaluators who assess the quality using a Translation Quality Score (TQS) on a 0 to 100 scale, with a score of 90% deemed acceptable. Errors such as grammar, punctuation, spelling, and mistranslation were examined, and each was assigned a severity level of minor, major, or critical. Three of the four languages in this paper were included in this dataset–Hausa (`hau`), Northern Sotho (`nso`) and Zulu (`zul`).

**FLORES-200:** This dataset expanded FLORES-101 to over 200 languages, including our fourth target language–Xitsonga (`tso`). In this data, a more comprehensive process was developed to ensure the quality of the translations. Specifically, professional translators and reviewers aligned on language standards before the translators translated the sentences. Afterwards, automated checks were first conducted and then followed by manual checks by independent reviewers. Translations that were found lacking quality were sent back for post-editing. Similarly to FLORES-101, translations scoring above 90% TQS were included in the FLORES-200.

## 2.1 Problems Identified in FLORES

Prior to this work, we have not found any published work that carefully reviews and attempt to correct mistakes in the FLORES evaluation dataset. However, some issues have been raised on the FLORES' public GitHub repositories.[3] Some of these issues include near-identical translations in several dialects of Arabic: Mesopotamian (`acm_arb`), Ta'izzi-Adeni (`acq_arb`), Najdi (`ars_arb`), and Moroccan (`ary_arb`) Arabic dialects were found to

be too similar to Standard Arabic (`arb`),[4,5] unspecifying the "orthography" and "variety" used in Lombard (`lmo_latn`) and Sardinian (`srd_latn`),[6,7] unmatched quotation marks,[8] and using Mandarin Chinese in Traditional Chinese Script (`zho_Hant`) for Cantonese (`yue_Hant`) translations.

## 3 Focus Languages and Evaluation

### 3.1 Languages Covered

In this work, the public[9] FLORES dev and devtest splits of Hausa, Northern Sotho (Sepedi), Xitsonga and isiZulu were reviewed and corrected by native speakers of the languages. A description of each language is presented in Appendix A.

### 3.2 Correction Guidelines

For reviewing and subsequently correcting the identified errors in the datasets, the participants were given the following guidelines.

**Reviewing:** At this stage, the participants identified sentences in both data splits that require reviewing.

- **Read the original text:** carefully read the original English text to understand the intended meaning and context.
- **Compare with translated text:** compare each sentence or phrase in the original English text with its corresponding translation. Pay attention to both the overall meaning and the nuances of the language.
- **Check for accuracy:** look for errors, inaccuracies, or deviations from the original meaning in the translation. This includes mistranslations, omissions, additions, and grammatical mistakes.
- **Evaluate clarity and cohesion:** assess whether the translated text is clear and coherent in the target language. Ensure that it flows naturally and is easy for a target language-speaking audience to understand.

---

[3]https://github.com/openlanguagedata/flores

[4]https://github.com/openlanguagedata/flores/issues/8

[5]https://github.com/facebookresearch/flores/issues/64

[6]https://github.com/openlanguagedata/flores/issues/5

[7]https://github.com/openlanguagedata/flores/issues/6

[8]https://github.com/facebookresearch/flores/issues/36

[9]https://github.com/openlanguagedata/flores

**Correcting the translations:** To correct the translations, we followed the guidelines provided in the shared task description.[10] The participants were trained on and encouraged to follow these guidelines when correcting the identified incorrect translations.

## 3.3 The Annotators

The correction task was conducted by volunteer annotators that focused on their native languages. These annotators were a mix of university students and researchers holding first, second and third degrees in computing and linguistics.

## 3.4 Evaluating the Corrections

To determine the amounts of corrections and subsequent differences between the original and corrected data, we used the following metrics. The computations were conducted only on the instances that were corrected. We used the original dataset as the supposed predictions and for the reference translations, we used the corrected data. We used the Natural Language Toolkit (NLTK) (Bird and Loper, 2004) for all tokenization.

**Token Difference:** This is the difference between the number of all tokens in the original and corrected datasets.

**Token Divergence:** This was used to measure the difference or dissimilarity between two sets of tokens. Given $T_o$ and $T_c$ as the set of tokens in the original and corrected datasets respectively, the following formula was used:

$$\text{divergence} = \frac{|T_o - T_c| + |T_c - T_o|}{|T_o \cup T_c|} \quad (1)$$

The formula calculates the proportion of tokens that are different between the two sets relative to the total number of unique tokens across both texts. Higher divergence score indicates that the two texts are quite different, suggesting significant changes or corrections were made.

**Translation Edit Rate:** (Snover et al., 2006) is a metric used in machine translation and other natural language processing tasks to measure the number of edits required to change a system-generated

translation into a reference translation, and is computed using the following formula.

$$\text{TER} = \frac{\text{\# of edits}}{\text{\# of words in ref. translation}} \quad (2)$$

The fewer the edits, the better the translation quality and a higher TER score indicates lower quality in the predicted translations.

**BLEU:** (Papineni et al., 2002) is an n-gram based metric that indicates the quality of generated machine translations. The BLEU is computed as follows:

$$\text{BLEU} = BP \times \exp\left(\sum_{n=1}^{N} w_n \log p_n\right) \quad (3)$$

where $BP$ is the Brevity Penalty and is used to penalize instances where shorter translations are generated when the reference is comparably longer; $p_n$ is the precision between the candidate translation and a set of ground truths; and $w_n$ is the n-gram weights.

**COMET:** (Rei et al., 2020) is a metric that leverages pre-trained neural models and cross-lingual word embeddings to evaluate the quality of machine translation systems. We used the pre-trained models provided by Wang et al. (2024).

## 4 Error Analysis

Tables 1 and 2 present how similar, or different, the original sentences were to the corrections. Some of the errors found are analyzed below per language.

**Hausa (hau)**  A significant part of the translations were suspected to have been automatically generated, as many of them appeared incoherent or unclear. To investigate this, we conducted a comparison with translations from the Hausa FLORES dataset and new translations generated by Google Translate. The comparison revealed that, although there were limited exact matches[11], several incorrect lexical choices in the dataset's translations aligned with those produced by Google Translate, supporting the suspicion that the translations may have been automatically generated. It is important to note that other translation tools may exist for Hausa that we did not evaluate. Furthermore, several sentence-level translations from Google Translate were found to be more qualitative and coherent

---

[10] https://oldi.org/guidelines#translation-guidelines

[11] Google Translate may have evolved since the creation of the dataset.

| lang | dev (997 sentences) | | | | | devtest (1,012 sentences) | | | | |
|------|------------|-----------|-----------|----------|--------|------------|-----------|-----------|----------|--------|
|      | # corr. (%) | # tokens$_o$ | # tokens$_c$ | $\Delta$ tokens | % div. | # corr. (%) | # tokens$_o$ | # tokens$_c$ | $\Delta$ tokens | % div. |
| hau  | 632 (63.4) | 17,948 | 18,073 | 125 | 24.7 | 70 (6.9)  | 2,006 | 1,978 | 28 | 49.2 |
| nso  | 67 (6.7)   | 2,226  | 2,271  | 45  | 28.9 | 62 (6.1)  | 2,082 | 2,105 | 23 | 28.0 |
| tso  | -          | -      | -      | -   | -    | 83 (8.2)  | 2,919 | 2,947 | 28 | 27.4 |
| zul  | 190 (19.1) | 3,605  | 3,588  | 17  | 23.7 | 226 (22.3)| 4,414 | 4,396 | 18 | 31.8 |

Table 1: Data statistics; # corr. (%) → number of sentences requiring at least one correction (percentage of original data); # tokens$_o$ → original token count; # tokens$_c$ → corrected token count; $\Delta$ tokens → token count difference; % div. → percentage of token divergence.

| lang. | dev | | | | devtest | | | |
|-------|-----|-----|------|-------|------|-----|------|-------|
|       | TER | | BLEU | COMET | TER | | BLEU | COMET |
|       | Score | # Edits | | | Score | # Edits | | |
| hau | 19.2 | 3,107 | 72.0 | 54.1 | 40.4 | 711 | 56.6 | 42.1 |
| nso | 22.4 | 472   | 68.5 | 55.2 | 21.2 | 409 | 71.8 | 55.9 |
| tso | -    | -     | -    | -    | 20.9 | 547 | 73.9 | 58.4 |
| zul | 17.2 | 524   | 76.3 | 53.0 | 23.6 | 879 | 70.6 | 53.0 |

Table 2: Similarities between the original and corrected FLORES evaluation data on the four African languages – original as predictions; corrected as reference translations.

than those in the current dataset. For an illustration, we examine sentences from the dev and devtest sets, see Table 3.

In several instance, named entities were translated instead of reusing them as they are due to the lack of their equivalents in Hausa. This is illustrated in the first example provided in Table 3. Planned Parenthood appears as an organization that was not supposed to be translated (and may only be explained as *hukuma mai kula da tsarin iyali*). The words in the organization name were translated as *Iyayen Tsararru*, with their literal word translations (*iyaye* → parents, *tsararru* → planned) instead of the name of the organization as a named-entity. In the second example, the phrase "standard business attire" was translated as *Kaya masu kala ɗaya su ne cikakkun tufafin mu'amala* instead of *kayan sawa na aiki da aka saba dasu*. The first translation is at best an incorrect explanation of the English phrase. And these are just two examples of the many we found in the dataset.

In addition to these severe mistakes, the dataset was littered with a lot of inconsistencies especially in the use of the standardized Hausa alphabets. Special characters are often ommitted and instead replaced with their normalized equivalents, e.g., ɓ → b, ɗ → d, etc. In some few places, the special ƴ is written as 'y which is acceptable.

**Northern Sotho (nso)** Several key challenges and areas for improvement were identified and cor-

rected, focusing on vocabulary consistency, syntax, spelling, and the accurate conveyance of technical terms. Most of the text was accurately translated and, for the text with problems, only small changes were required to make it more accurate. Some of the words like "*safatanaga* and *disafatanaga*" have generally maintained lexical consistency although they were wrongly translated. These have been corrected to "*sefatanaga* or *difatanaga* (plural)".

Although sometimes Sepedi uses borrowed words for many technical and scientific terms, things such as pavement do have a translation which could be "*tsela ya maoto or tselanathoko*". These could have been used instead of borrowing the pavement term to say *pabamente*. The use of a borrowed term could have been from the available corpus or from learned behaviour for borrowing unknown English terms. Another example is the word college which was translated to *colleje*, but Sepedi has a standard borrowed translation: "*kholetšhe*".

Addressing spelling errors and ensuring proper spacing between words are vital for readability and comprehension. For instance, the word "*tswarelo*" was corrected to "*tshwarelo*" to reflect the proper spelling. Similarly, "*patlaladitše*" was adjusted to "*phatlaladitše*", and "*bontša*" to "*bontšha*". Additionally, "*mephutso*" should be spelt as "*meputso*", and "*delo*" should be corrected to "*selo*". Spacing was required when using "*begona*" so that it is "*be gona*" and similar adjustments were made. These adjustments are crucial to maintain lexical consis-

| SN | English | Wrong Translation in FLORES | Corrected Translation |
|---|---|---|---|
| 1. | Komen's policy disqualified Planned Parenthood due to a pending investigation on how Planned Parenthood spends and reports its money that is being conducted by Representative Cliff Stearns. | Manufar Komen ta hana Iyayen Tsararru sanadiyyar binciken kashe kuɗi kan yadda Tsararren Iyaye yake ciyarwa kuma ta ba da rahoton kuɗaɗɗinta wanda Wakilin Cliff Stearns ke gudanarwa. | Manufar Komen ta dakatar da chanchantar Planned Parenthood sanadiyyar binciken da akeyi akan yanda Planned Parenthood take kashewa da kuma bayar da ba'asin kuɗin ta wanda Wakili Cliff Stearns yake gudanarwa. |
| 2. | Suits are standard business attire, and coworkers call each other by their family names or by job titles. | Kaya masu kala ɗaya su ne cikakkun tufafin mu'amala, kuma abokan aiki kan kira junansu da sunan iyalinsu ko da muƙaman aiki. | Kwat sune kayan sawa na aiki da aka saba dasu kuma abokan aiki suna kiran juna ne da sunan gidansu ko kuma matsayin da mutum yake kai. |

Table 3: Some Hausa Examples of incorrect and inconsistent translations in FLORES dev and devtest.

tency and to ensure that translations are accurate and easily understood.

Some terms were left out, like "scientific" as "*tša bo ramahlale*" when scientific tools were talked about, and this greatly affected the meaning of the sentence. Additionally, in another instance, a sentence describing the use of Caesarean section to give birth to Nadia was misleading. Incorrectly, it implied that Nadia was both the baby being born and the individual undergoing the operation. This was corrected to have the intended meaning.

**Xitsonga (`tso`)** Some of the problems identified in the Xitsonga translations included problems to do with vocabulary accuracy and the use of borrowed words. Among the errors that were identified is the translation of "Type 1 diabetes" to "*vuvabyi bya chukela bya Type 1*". The correct phrase should therefore be "*vuvabyi bya chukela bya muxaka wo sungula*", which captures the type of diabetes and avoid misunderstanding. Similar trends raise the importance of using proper terms that might fit local context as opposed to directly translating English words.

Another problem was that translations were mostly uniform, without contextual variations. Even here, the words "*xiyenge xa tlilinikhali na sayense*" (clinical and scientific division) were used wrongly. The word actually is "*xiyenge xa vutshila ni ntokoto bya sayense*" (clinical and scientific division), but this clearly passes on the intended meaning. Moreover, the use of pluralization of terms was arbitrary. While the singular form of the term "worker" is "*mutirhi*", the plural form should be "*vatirhi*", and the singular form of "methods" is "*maendlelo*", which should be in plural throughout instead of appearing in single forms.

Spelling problems and the usage of borrowed terms can have a substantial influence on the correctness of Xitsonga translations. One of the most illustrative examples of such incongruity of terms is that the English word "channel" has been translated as "*chanele*". Instead, the work should have used the original term "*nongonoko*" in order to ensure a perfect linguistic and connotative translation. To avoid generation of wrong impressions, the phrase borrowed from IsiZulu as used to mean "President" had to be replaced by the word "*murhangeri wa tiko*" from Xitsonga. Deficient spelling, as in the case of writing "*dokodela*" instead of "Dr", and examples of slang such as using "*mwana wa*" instead of the formal "*muongori*" indicate how borrowing and spelling mistakes reduced the quality of the translations. Fluency and correct spelling as well as using the native language correctly are a necessity to maintain the translated material's effectiveness.

**isiZulu (`zul`)** Similar to the errors identified in the other languages above, isiZulu translations displayed several common challenges. These included inconsistencies in vocabulary, syntax errors, and issues with the accurate expression of technical and scientific terms. The agglutinative nature of isiZulu and its conjunctive writing style further worsen these issues, leading to specific translation errors related to morphology and orthography.

A closer examination of these challenges reveals issues such as in the translation of "Around 11:29, the protest moved up Whitehall, ..." which was initially rendered as "*Ngawo-11:29 ababhikishi baya Odongeni Olumhlophe, ...*". This translation contains two key issues. First, "*Ngawo-11:29*" should have been "*Ngabo-11:29*" to correctly match the grammatical structure for time

expressions in isiZulu. Second, the literal transliteration of "Whitehall" as "*Odongeni Olumhlophe*" failed to integrate properly into the sentence. The correct approach would involve incorporating the place name with the locative prefix "e-" to produce "e-Whitehall.". This prefix addition is required in conjunctive languages when using borrowed words or terms, but MT systems often fail to capture these variations. Additionally, another common issue was the unnecessary borrowing of words from English, despite the availability of standardized isiZulu terms. This was particularly evident with month names, scientific terms, and country names, where inconsistencies were frequent—one translation might use "January," another "*uJanuwari*," and yet another "*uMasingana*" Another example of this can be seen with the country name "Spain," which was inconsistently translated as both "Spain" and "*Speyini*" in different sections. Similar inconsistencies occurred with attempts to translate organizational names or acronyms, leading to partial translations that disrupted the linguistic flow.

To address the inconsistencies, standardized isiZulu terms were consistently applied throughout the translations. For instance, month names such as "*uMasingana*" replaced the inconsistent use of "January" and "*uJanuwari*" In dealing with organizational names and acronyms and countries' names, the approach was to fully translate these entities or retain their original form consistently, avoiding partial translations that could disrupt the flow.

In addition to the inconsistencies with terminology, other errors were also identified and addressed. These included issues with verb conjugation, where incorrect tenses or forms were initially used, and the improper handling of borrowed words that did not align with isiZulu's morphosyntactic rules. Minor spelling errors and incorrect use of prefixes or suffixes were also corrected to ensure that the translations were both grammatically accurate and easily understood.

## 5    Conclusion

In this work, we highlight the importance of qualitative evaluation datasets for low-resource languages and present our findings from a comprehensive review of the FLORES dataset for four African languages: Hausa, Northern Sotho, Xitsonga, and isiZulu. The original translations were marred by vocabulary inconsistencies, syntax errors, and in-

accurate technical terms. After making necessary corrections, we measured the amount of edits and resulting difference between the improved datasets and the original using metrics like BLEU, TER, and COMET, which showed that significant improvements were made. The results presented highlight the need for ongoing refinement and human oversight in developing accurate translation datasets for underrepresented languages. For future work, we intend to expand the corrections to more African languages.

## References

Idris Abdulmumin, Satya Ranjan Dash, Musa Abdullahi Dawud, Shantipriya Parida, Shamsuddeen Muhammad, Ibrahim Sa'id Ahmad, Subhadarshi Panda, Ondřej Bojar, Bashir Shehu Galadanci, and Bello Shehu Bello. 2022. Hausa visual genome: A dataset for multi-modal English to Hausa machine translation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6471–6479, Marseille, France. European Language Resources Association.

David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruiter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajuddeen Gwadabe, Freshia Sackey, Bonaventure F. P. Dossou, Chris Emezue, Colin Leong, Michael Beukman, Shamsuddeen Muhammad, Guyo Jarso, Oreen Yousuf, Andre Niyongabo Rubungo, Gilles Hacheme, Eric Peter Wairagala, Muhammad Umair Nasir, Benjamin Ajibade, Tunde Ajayi, Yvonne Gitau, Jade Abbott, Mohamed Ahmed, Millicent Ochieng, Anuoluwapo Aremu, Perez Ogayo, Jonathan Mukiibi, Fatoumata Ouoba Kabore, Godson Kalipe, Derguene Mbaye, Allahsera Auguste Tapo, Victoire Memdjokam Koagne, Edwin Munkoh-Buabeng, Valencia Wagner, Idris Abdulmumin, Ayodele Awokoya, Happy Buzaaba, Blessing Sibanda, Andiswa Bukula, and Sam Manthalu. 2022. A few thousand translations go a long way! leveraging pre-trained models for African news translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3053–3070, Seattle, United States. Association for Computational Linguistics.

Steven Bird and Edward Loper. 2004. NLTK: The natural language toolkit. In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pages 214–217, Barcelona, Spain. Association for Computational Linguistics.

David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2022. *Ethnologue: Languages of the World*, 25 edition. SIL International, Dallas, Texas.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Kr-

ishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2021. The flores-101 evaluation benchmark for low-resource and multilingual machine translation.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.

Ximbani Eric Mabaso. 2018. Xitsonga in south africa. *The Social and Political History of Southern Africa's Languages*, pages 311–330.

Stuart Mesham, Luc Hayward, Jared Shapiro, and Jan Buys. 2021. Low-resource language modelling of south african languages. *arXiv preprint arXiv:2104.00772*.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Derek Nurse and Gérard Philippson. 2006. *The bantu languages*, volume 4. Routledge.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Shantipriya Parida, Idris Abdulmumin, Shamsuddeen Hassan Muhammad, Aneesh Bose, Guneet Singh Kohli, Ibrahim Said Ahmad, Ketan Kotwal, Sayan Deb Sarkar, Ondřej Bojar, and Habeebah Kakudi. 2023. HaVQA: A dataset for visual question answering and multimodal research in Hausa language. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10162–10183, Toronto, Canada. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

StatsSA. 2022. Statistics South Africa.

Elsabé Taljard and Sonja E Bosch. 2006. A comparison of approaches to word class tagging: Disjunctively vs. conjunctively written bantu languages. *Nordic journal of African studies*, 15(4).

Jiayi Wang, David Adelani, Sweta Agrawal, Marek Masiak, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Xuanli He, Sofia Bourhim, Andiswa Bukula, Muhidin Mohamed, Temitayo Olatoye, Tosin Adewumi, Hamam Mokayed, Christine Mwase, Wangui Kimotho, Foutse Yuehgoh, Anuoluwapo Aremu, Jessica Ojo, Shamsuddeen Muhammad, Salomey Osei, Abdul-Hakeem Omotayo, Chiamaka Chukwuneke, Perez Ogayo, Oumaima Hourrane, Salma El Anigri, Lolwethu Ndolela, Thabiso Mangwana, Shafie Mohamed, Hassan Ayinde, Oluwabusayo Awoyomi, Lama Alkhaled, Sana Al-azzawi, Naome Etori, Millicent Ochieng, Clemencia Siro, Njoroge Kiragu, Eric Muchiri, Wangari Kimotho, Toadoum Sari Sakayo, Lyse Naomi Wamba, Daud Abolade, Simbiat Ajao, Iyanuoluwa Shode, Ricky Macharm, Ruqayya Iro, Saheed Abdullahi, Stephen Moore, Bernard Opoku, Zainab Akinjobi, Abeeb Afolabi, Nnaemeka Obiefuna, Onyekachi Ogbu, Sam Ochieng', Verrah Otiende, Chinedu Mbonu, Yao Lu, and Pontus Stenetorp. 2024. AfriMTE and AfriCOMET: Enhancing COMET to embrace under-resourced African languages. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5997–6023, Mexico City, Mexico. Association for Computational Linguistics.

## A Description of the Target Languages

**Hausa (`hau`):** Hausa is a widely spoken language across West Africa, particularly in Nigeria, Niger, Cameroon, and Ghana. It is spoken by approximately 77 million people worldwide, primarily in West Africa (Eberhard et al., 2022). Hausa ranks as the second most spoken language in Africa and 27th globally. The language belongs to the Chadic branch of the Afroasiatic language family, and it has a rich history of written communication. It was first written in Arabic script known as Ajami, reflecting the language's connection to Arabic, with many Hausa words borrowed from Arabic due to historical contact and influence. Today, the Boko script (also known as Roman script), which uses Latin characters, is the most common writing system for Hausa. This script excludes the letters p, q, v, and x, and includes additional consonants (ɓ, ɗ, ƙ, ƴ, kw, ƙw, gw, ky, ƙy, gy, sh, ts) and vowels (long a, i, o, u, e, and two diphthongs ai and au). Hausa follows a Subject-Verb-Object (SVO) sentence structure.

**Northern Sotho (`nso`):** Northern Sotho, also known as Sepedi or Sesotho sa Leboa, is one of the official languages of South Africa and is spoken primarily by the Bapedi people in Limpopo Province. It is a Bantu language that belongs to the Sotho-Tswana group and shares linguistic similarities with Sesotho (Southern Sotho) and Setswana. Sepedi is known for its rich oral tradition that includes folklore, proverbs, and praise poetry that have played a significant role in the preservation of cultural heritage (Nurse and Philippson, 2006). Sepedi is written using the Latin alphabet, with the standard 26 letters and a few additional characters such as the "š" which are adapted to its unique sounds. The language primarily follows a Subject-Verb-Object word order in sentence structure.

**Xitsonga (`tso`):** Xitsonga, or Tsonga, is a Bantu language that is mainly spoken in South Africa and more especially in the Limpopo province and parts of Mpumalanga province. The language is estimated to be spoken by about 2.3 million people in South Africa. Xitsonga belongs to the Niger-Congo language family, specifically the Tshwa-Ronga subgroup, and is characterized by the extensive use of prefixes and suffixes to convey meaning (Mabaso, 2018). This linguistic feature can impact the accuracy of translations, especially when dealing with

| Language | Sentence |
|---|---|
| English | I know them |
| Hausa | Na san su |
| Northern Sotho | Ndza va tiva |
| Xitsonga | Ke a ba tseba |
| isiZulu | Ngiyabazi |

Table 4: The grammatical structure of different languages.

technical and scientific concepts. It also feature a complex system of writing and syntax, which are prerequisites to clear and concise language usage. Xitsonga is currently used in education and media section in South Africa, thus is regarded as relevant in cultural linguistic practices. That is why, the language being mentioned as a part of the country's multiple languages system emphasizes its relevance and application in different phases of the people's activity.

**isiZulu (`zul`):** Zulu or isiZulu (in Zulu) is one of the 12 official languages in South Africa, and it is considered to be the most widely spoken indigenous language in the country. It constitutes approximately a quarter of the population, with around 15.1 million speakers out of the population of 62 million people (StatsSA, 2022). IsiZulu is part of the Nguni language family, which is made up of a group of closely related Bantu languages belonging to a larger Niger-Congo language family, and they are widely spoken across Southern Africa (Mesham et al., 2021). These languages are particularly notable for their complex morphology, characterized by agglutinative morphology and conjunctive orthography. Agglutinative morphology means that words are typically formed by combining multiple small meaning-carrying units, known as morpheme. Conjunctive orthography means that the morphemes are glued together to form a word, rather than writing them with spaces in between, as seen in disjunctive orthography, commonly associated with the Sotho group, as well as Tshivenda and Xitsonga in South Africa indigenous languages (Taljard and Bosch, 2006). To illustrate this distinction, consider the example in Table 4 which examines the different grammatical structures of the phrase *I know them*.

Table 4 shows that while the phrase's meaning is consistent across languages, the writing systems vary: in disjunctive orthography, morphemes are

separated by spaces, while in conjunctive orthography, as in isiZulu, they are joined into a single word. For example, in the phrase *I know them*, each morpheme serves a specific grammatical function–'I' as the subject, 'know' as the verb, and 'them' as the object. In disjunctive orthography, these morphemes are written separately, making each unit distinct. In conjunctive orthography, they are combined into one continuous word, but the meaning remains intact. These orthographic variations pose challenges for machine translation systems, which must accurately process morphemes in different writing systems to produce accurate translations.

# Expanding FLORES+ Benchmark for more Low-Resource Settings: Portuguese-Emakhuwa Machine Translation Evaluation

**Felermino D. M. A. Ali[1,2,3,5], Henrique Lopes Cardoso[1,2], Rui Sousa-Silva[3,4]**

[1]Laboratório de Inteligência Artificial e Ciência de Computadores (LIACC / LASI)

[2]Faculdade de Engenharia da Universidade do Porto, Rua Dr. Roberto Frias, 4200-465 Porto, Portugal

[3]Centro de Linguística da Universidade do Porto (CLUP)

[4]Faculdade de Letras da Universidade do Porto, Via Panorâmica, 4150-564 Porto, Portugal

[5]Faculdade de Engenharia da Universidade Lúrio, Pemba 3203, Mozambique

{up202100778, hlc}@fe.up.pt, rssilva@letras.up.pt

## Abstract

As part of the Open Language Data Initiative shared tasks, we have expanded the FLORES+ evaluation set to include Emakhuwa, a low-resource language widely spoken in Mozambique. We translated the *dev* and *devtest* sets from Portuguese into Emakhuwa, and we detail the translation process and quality assurance measures used. Our methodology involved various quality checks, including post-editing and adequacy assessments. The resulting datasets consist of multiple reference sentences for each source. We present baseline results from training a Neural Machine Translation system and fine-tuning existing multilingual translation models. Our findings suggest that spelling inconsistencies remain a challenge in Emakhuwa. Additionally, the baseline models underperformed on this evaluation set, underscoring the necessity for further research to enhance machine translation quality for Emakhuwa. The data is publicly available at https://huggingface.co/datasets/LIACC/Emakhuwa-FLORES

## 1 Introduction

Evaluation datasets and benchmarks are essential for advancing Natural Language Processing (NLP) models. They provide the necessary tools for assessing model performance and guiding further enhancements. However, the scarcity of evaluation datasets and benchmarks for low-resource languages has significantly hindered the progress of NLP technologies in these languages. Recognizing this challenge, the FLORES+ evaluation set has emerged as a critical tool for the Machine Translation (MT) community, especially in low-resource languages. It promotes a more inclusive approach to language technology development across diverse linguistic landscapes. This work focuses on expanding the FLORES+ (NLLB Team et al., 2022) evaluation set to include Emakhuwa, a low-resource language spoken in Mozambique

by approximately 9 million people. Our dataset consists of the *dev* and *devtest* sets managed by the Open Language Data Initiative[1] (OLDI), which contain 997 sentences and 1012 sentences, respectively. Throughout our data collection process, we implemented robust quality assurance mechanisms, including thorough post-editing. The resulting dataset features multiple reference translations derived from these post-editing efforts.

## 2 Related Works

The Flores v1.0 MT evaluation set was introduced by Guzmán et al. (2019). This initial version focused on two language pairs: Nepali–English and Sinhala–English, with the data divided into *dev*, *test*, and *devtest* splits. After its release, the dataset was gradually expanded to include more languages. A significant expansion came with the work of Goyal et al. (2021), who introduced Flores-101, extending the evaluation set to support 101 languages. Further expansion was done with the release of Flores-200 by the NLLB team (NLLB Team et al., 2022) in 2022, which increased the language coverage to 204 languages. Additional contributions include Doumbouya et al. (2023), who added the Nko language, as well as AI4Bharat et al. (2023), who incorporated Bodo, Dogri, Meitei, Sindhi, and Goan Konkani into the dataset. These contributions have significantly broadened the opportunities for low-resource languages in MT, allowing researchers to track the progress of MT systems on these expanded evaluations. However, the coverage remains limited, especially considering that there are over 7,000 languages worldwide. One such language that remains underserved is Emakhuwa, which still lacks datasets for MT.

---

[1]https://oldi.org/

579

## 3 Emakhuwa

Emakhuwa, alternatively referred to as Makua, Macua, or Makhuwa, belongs to the Bantu language family and is predominantly spoken in the northern and central regions of Mozambique, specifically in the Nampula, Niassa, Cabo Delgado, and Zambezia provinces. There are eight variants of Emakhuwa, with Emakhuwa-Central (ISO 639-3 code *vmw*) being the standard variety (Ngunga and Faquir, 2014).

Emakhuwa follows the Subject-Verb-Object (SVO) structure, use a Latin scripts (ISO 15924 *Latn*), and is gender-neutral. Furthermore, similarly to other languages in the Bantu family, it is linguistically rich, with complex morphology featuring agglutinative and tonal attributes.

### 3.1 Challenges in Emakhuwa

Emakhuwa digital resources are scarce, and the spelling standards are still under development. While a fully standardized system is not yet in place, the existing guidelines (Ngunga and Faquir, 2014) offer a critical framework for contemporary written communication in Emakhuwa. One problem stressed in official standardization (Ngunga and Faquir, 2014) is the lack of guidance on tonal marking. Consequently, existing materials exhibit inconsistent spelling, particularly when marking tone, which is essential in Emakhuwa for disambiguation. To give an example, let us consider two words carrying distinct meanings: *omala* and *omaala* / *omàla*; *omala* means "to finish", while *omaala* / *omàla* means "to silence" or "to hush." In this case, the tonal marker *aa* / *à* clarifies the intended meaning.

Spelling variations are largely evident in existing Emkahuwa text corpora, where some use diacritics (e.g., *à*, *è*, *ì*, *ò*, *ù*) and consonantal sounds (e.g, *kh*, *nn*) for tonal marking, others use vowel lengthening (e.g., *aa*, *ee*, *ii*, *oo*, *uu*), and some even use a combination of methods. Emakhuwa's agglutinative nature with complex morphology further amplifies spelling discrepancies. Since tonal variations often occur at the morpheme level, different combinations of morphemes result in varied spellings of the same word.

These spelling inconsistencies create significant obstacles for language technology processes. They lead to data sparsity, as some spelling variants appear less frequently, which impairs the model's ability to learn the language's nuances effectively.

This sparsity inflates the vocabulary size and can result in reduced performance of language technologies.

An additional challenge in Emakhuwa that contributes to inconsistencies is the adaptation of loanwords. Emakhuwa text corpora frequently contain Portuguese loanwords with inconsistent adaptations due to the absence of standardized guidelines for integrating borrowed terms (Ali et al., 2024). These loanwords are adapted in one of three ways: phonetically to match Portuguese pronunciation, in alignment with Emakhuwa phonotactics, or retained unchanged from Portuguese.

## 4 Methodology

We chose to translate the *devtest* and *dev* sets from Portuguese (*pt*) into Emakhuwa (*vmw*) because our translators were only proficient in these two languages. We focus specifically on the central variant of Emakhuwa, as it is the standard and established language variant.

The translators were selected based on their proficiency in these languages and their proven experience in Portuguese-Emakhuwa translation. In total, we collaborated with five experts: two were assigned the tasks of translation and revision, while the remaining collaborators were responsible for evaluating the translations (refer to Table 6 in the appendix for more details).

In general, we implemented the workflow as a peer review process, divided into three main steps: Data Preparation, Translation, and Validation. Below is a detailed description of each step (refer to Figure 1).

### 4.1 Data Preparation

We compile the sentences in *devtest* and *dev* sets as segments and then load them to the Matecat[2] CAT (Computer-Assisted Translation) tool. Before assigning translation tasks, we prepare a guideline and glossary. The guidelines were adapted from the Open Language Data Initiative guidelines[3], written in Portuguese and suggesting that the translated text should adhere to the latest orthography standards of the central variant of Emakhuwa. On the other hand, the glossary was built by digitizing existing bilingual dictionaries and the glossary of Political, Sports, and Social Concepts from Radio of Mozambique (Moçambique E.P., 2016). We con-

---

[2]https://www.matecat.com/
[3]https://oldi.org/guidelines

Figure 1: Workflow

ducted a small workshop to familiarize the team with the guidelines and gather feedback to improve them. The translation team found the glossary helpful, as it prevented using loanwords for existing Emakhuwa terms and ensured consistency in translations.

## 4.2 Translation

Translation tasks were subdivided between two translators: one worked on the *devtest* segments, and the other on the *dev* segments. Once all segments were translated, they were submitted to our spell checker system for an automatic check to identify potential misspellings (refer to Figure 10 in the appendix). We then provided feedback to the translators, asking them to review and refine their work if necessary.

## 4.3 Validation

The validation corresponds to two steps: revision and Judgments.

### 4.3.1 Revision

Following the translation step, we swapped the translated works between the two translators, asking them to post-edit each other's translations on the Matecat platform. Table 1 provides the Quality Report generated by Matecat, which includes various metrics used to evaluate the translation based on the revisions made. The report indicates that the reviewer working on *devtest* made more suggestions. A closer examination of the error categories on *devtest* (refer to Figure 9 in the appendix) reveals that most of the issues identified in the translation fell under the category of "Language Quality", meaning grammar, punctuation, and spelling errors. On the other hand, the reviewer of the *dev* set

identified mostly errors related to "terminology and language consistency", suggesting that the translator was not consistently using the proper terms and maintaining uniformity throughout the text.

| | dev | devtest |
|---|---|---|
| Post-Editing Effort | 99% | 95% |
| Time to edit | 02m38s | 05m42s |
| Quality score | 23.31 | 54.22 |
| Avg. Edit Distance | $0.23 \pm 1.77$ | $7.09 \pm 11.94$ |

Table 1: Matecat's quality report post revision.

### 4.3.2 Judgments

Once all segments have been revised, we perform a second translation quality assessment using a Direct Assessment (DA) pipeline similar to the one described by Guzmán et al. 2019. Judgments were collected using our annotation tool (see Figure 7 in the appendix), and involve the following aspects.

**Direct Assessment**   Three different raters evaluate the translation adequacy (i.e., the perceived translation quality) on a scale from 0 to 100. A score of 0 means that "no meaning was preserved in the translation". Scores from 1 to 34 - "the translation preserves some of the source meaning but loses significant parts", 35 and 67 - "the translation retains most of the source meaning", 68 to 99 - "the translation is consistent with the source text", and a score of 100 means "the translation is perfect". These quality intervals are inspired by the study of Wang et al. 2024.

**Control**   To ensure raters' attentiveness and improve consistency during the evaluation, we included control instances with incorrect translation pairs. These incorrect pairs were generated using

the Madland-400-3bt[4] model (Kudugunta et al., 2024), a multilingual MT system that supports the Emetto variant of Emakhuwa (ISO 639-3 *mgh*). While this model typically performs poorly when translating from Portuguese to Emakhuwa, it produces similar words that can mislead inattentive annotators. Based on these control translations, we provided feedback to the evaluators as they progressed in their tasks. We used emojis to give the feedback in our annotation tool: a 🙂 appeared if less than 25% of control translations were incorrectly rated (i.e. scores above 34 points), 😠 if 25%-50% are incorrectly rated, 😡 if 50%-75% are incorrectly rated, and 😭 if more than 75% of control translations are rated too highly.

**Post Editing** During the validation phase, we asked evaluators to post-edit translations with lower scores to enhance fluency and better align them with the source sentence's meaning. However, this task was made optional to prevent evaluators from inflating scores to avoid additional post-editing work.

**Standard orthography** To assess the perceived usage of standard orthography, raters also judged whether the translated text used standard orthography on a scale from 1 (not using standard orthography) to 5 (entirely written in standard orthography).

Finally, we calculate the average score for each segment. We then returned segments scoring below 70 to the translator for reworking. Figure 2 shows the histogram of the average translation scores.



Figure 2: Averaged Translation Quality Score Histogram on both *dev* and *devtest* sets. Translations with an average score below 70 (indicated by the red line) were returned to the translator for rework.

### 4.4 Analysis

Figures 3 and 4 show the raw scores per annotator for Direct Assessments. Given the mean scores, in both the *test* and *devtest* sets, Annotator 1 and Annotator 2 gave higher quality scores, while Annotator 3 was more critical but still within the spectrum of acceptable translations. This suggests a generally positive perception of the translations produced. Figure 5 displays the Direct Assessment scores on the control set. Annotator 1 and Annotator 3 have median scores below the threshold of 34 points, suggesting that, as expected, they have generally assessed the control translations as low quality. Annotator 2, however, has a median score above the threshold, suggesting a trend to a more positive assessment compared to the other two annotators and was less attentive among the annotators.



Figure 3: Direct Assessment adequacy scores per annotator on *dev* set



Figure 4: Direct Assessment adequacy scores per annotator on *devtest* set

Table 2 provides the reliability results for adequacy and standard orthography usage assessments. The inter-class correlation for adequacy is 0.67 for

Figure 5: Direct Assessment adequacy scores per annotator on *control* set



Figure 6: Assessment of standard orthography usage on the control set.

*dev* and 0.66 for *devtest*, suggesting moderate reliability. However, the inter-class correlations for standard orthography usage are lower, with values of 0.35 for *dev* and 0.27 for *devtest*, indicating considerable disagreement among annotators. This discrepancy highlights the ongoing lack of clarity regarding Emakhuwa spelling standards, as further illustrated in Figure 6, which depicts the varying assessments of standard orthography.

| | Adequacy | | Orthography | |
| | dev | devtest | dev | devtest |
|---|---|---|---|---|
| ICC | 0.67 | 0.66 | 0.35 | 0.27 |
| CI | [0.63, 0.71] | [0.62, 0.7] | [0.27, 0.42] | [0.18, 0.35] |

Table 2: Intraclass Correlation Coefficient (ICC) and Confidence Interval (CI) Results for Adequacy and Orthography usage annotation.

### 4.5 Dataset Collected

Table 3 presents the statistics for the *devtest* and *dev* sets resulting from the completion of the translation tasks. The *dev* set comprises 997 sentence pairs,

| | dev | devtest |
|---|---|---|
| #ref. | 997 | 1,012 |
| #ref. words | 18,673 | 21,011 |
| #post-edited refs | 1,848 | 1,889 |

Table 3: Statitics for the resulting dataset sets

while the *devtest* set contains 1,012 sentence pairs. A sample of the dataset is displayed in Table 7 in the appendix.

## 5 Experiments

This section describes the experiment involving training neural MT models using the training sets described below. Then, we performed a comprehensive benchmark evaluation using the evaluation sets introduced in this study.

### 5.1 Training Data

To train the models, we used the data outlined below:

- Ali et al. (2021) dataset: This subset comprises parallel data in Portuguese and Emakhuwa from different sources, including online texts from the Jehovah's Witness, the African Story Book websites, and Optical Character Recognition (OCR) extracted texts. The corpus contains diverse writing styles, spelling styles, and genres.

- Parallel News: This subset consists of news articles translated from Portuguese into Emakhuwa.

The dataset includes around 63k training parallel sentences and 964 validation parallel sentences, spanning a range of topics (see Table 4), where a significant portion of the data comes from the religious domain, mainly consisting of translations of biblical texts.

| | Sentences | | Tokens | |
| Source | Train | Dev | *pt* | *vmw* |
|---|---|---|---|---|
| Ali et al. (2021) | 46,454 | 399 | 1,104,279 | 951,520 |
| News | 17,403 | 565 | 596,066 | 541,598 |
| **Total** | 63,857 | 964 | 1,700,345 | 1,493,118 |

Table 4: Training and Validation data statistics

| | | dev | | | | devtest | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Single Ref. | | Multi Ref. | | Single Ref. | | Multi Ref. | |
| | | BLEU | chrF | BLEU | chrF | BLEU | chrF | BLEU | chrF |
| **transformer-base** | | | | | | | | | |
| **Baseline** | pt→vmw | 3.7 | 30.67 | 3.95(+0.25) | 31.32(+0.65) | 3.27 | 29.23 | 3.57(+0.3) | 29.84(+0.61) |
| | vmw→pt | 4.36 | 25.48 | - | - | 2.93 | 23.96 | - | - |
| **Multilingual Language Models** | | | | | | | | | |
| **afri-byT5** | pt→vmw | 10.32 | 41.88 | 10.81(+0.49) | 42.64(+0.76) | 7.03 | 35.87 | 7.73(+0.7) | 36.72(+0.85) |
| | vmw→pt | **22.45** | **47.31** | - | - | 13.74 | **37.78** | - | - |
| **afri-mT5** | pt→vmw | 5.66 | 35.37 | 5.96(+0.3) | 36.01(+0.64) | 4.7 | 32.7 | 5.06(+0.36) | 33.25(+0.55) |
| | vmw→pt | 12.12 | 38.18 | - | - | 7.39 | 32.92 | - | - |
| **byT5** | pt→vmw | **10.66** | **42.37** | **11.2**(+0.54) | **43.16**(+0.79) | **7.49** | 36.33 | **8.13**(+0.64) | **37.15**(+0.82) |
| | vmw→pt | 22.24 | 47.01 | - | - | **14.1** | 37.75 | - | - |
| **mT0** | pt→vmw | 5.52 | 30.33 | 5.76(+0.24) | 30.9(+0.57) | 4.69 | 27.89 | 5.02(+0.33) | 28.36(+0.47) |
| | vmw→pt | 17.46 | 38.92 | - | - | 10.63 | 32.69 | - | - |
| **mT5** | pt→vmw | 6.76 | 34.09 | 7.18(+0.42) | 34.8(+0.71) | 5.67 | 31.67 | 6.06(+0.39) | 32.23(+0.56) |
| | vmw→pt | 15.42 | 37.58 | - | - | 9.65 | 32.22 | - | - |
| **Many-to-Many Multilingual Translation Language Models** | | | | | | | | | |
| **M2M100** | pt→vmw | 8.25 | 39.22 | 8.79(+0.54) | 40.14(+0.92) | 6.92 | 36.33 | 7.57(+0.65) | 37.19(+0.86) |
| | vmw→pt | 21.08 | 45.31 | - | - | 13.67 | 37.46 | - | - |
| **NLLB** | pt → vmw | 8.19 | 41.44 | 8.74(+0.54) | 42.32(+0.88) | 5.88 | 36.13 | 6.34(+0.46) | 37.01(+0.88) |
| | vmw→pt | 17.41 | 42.88 | - | - | 10.35 | 35.05 | - | - |

Table 5: BLEU and chrF scores for various models on *dev* and *devtest* splits, for single and multiple references

## 5.2 Setup

We trained MT models in both directions, *pt-vmw* (Portuguese to Emakhuwa) and *vmw-pt* (Emakhuwa to Portuguese), using two approaches: training a vanilla transformer model and fine-tuning existing multilingual language models.

**Training** We adopt the transformer architecture (Vaswani et al., 2017), implemented through the OpenNMT toolkit (Klein et al., 2017). The model consists of an encoder and decoder comprising 6 layers, 8 heads, and 512 hidden units in the feed-forward network. We used an embedding size of 512 dimensions for both source and target words and a batch size of 32. We applied layer normalization and added dropout with a 0.1 probability to the embedding and transformer layers. Additionally, the Adam optimizer (Kingma and Ba, 2014) was used, and a learning rate of 0.0002. The checkpoints were saved every 1000 updates. We preprocess the input, applying the Byte Pair Encoding subword segmentation.

**Fine-tuning Multilingual Models** Multilingual language models are one of the most prominent approaches to low-resource languages nowadays since it enables knowledge transfer among related languages, making cross-lingual transfer and zero-shot learning possible.

In our experiments, we fine-tuned various multilingual language models that are well-established in the literature, namely: mT5 (Xue et al., 2021), byT5 (Xue et al., 2022), and the multilingual translation models M2M-100 (Fan et al., 2021) and NLLB (NLLBTeam et al., 2024). Specifically, we use mT5-base (580M parameters), byT5-base (580M parameters), M2M-100 (418M parameters), and NLLB-200's distilled variant (600M parameters). Additionally, we also fine-tuned the African-centric language models, namely, AfribyT5 (580M parameters) and AfrimT5 (580M parameters) by Adelani et al., 2022.

## 5.3 Evaluation

To assess the systems' performance, we used the SacreBLEU toolkit (Post, 2018) to compute the BLEU (Papineni et al., 2002) and ChrF scores (Popović, 2015).

## 6 Results and Discussion

Results are presented in Table 5. Our baseline results, derived from a vanilla transformer-base model, set a foundational performance benchmark. On the *devtest* set, the baseline model achieved a BLEU score of 2.93 and a ChrF score of 23.96 for the *vmw → pt* translation direction. These modest scores underscore the limitations of the vanilla

transformer-base model in handling the complexities of translation tasks involving low-resource languages like Emakhuwa.

However, introducing multilingual language models enhanced translation performance, particularly in the *vmw → pt* direction. Among these, models based on byT5 demonstrated superior performance. For instance, the fine-tuned byT5 model achieved a BLEU score of 14.1 and a ChrF score of 37.75 on the *devtest* set, which marks a substantial improvement over the baseline. This highlights the advantage of leveraging tokenization-free approaches, which are better suited for handling the morphological richness and orthographic variations characteristic of Emakhuwa.

Across Table 5, our results show that while BLEU scores remained relatively low in the *pt → vmw* translation direction, ChrF were consistently higher. This discrepancy between BLEU and ChrF scores suggests that BLEU may be disproportionately penalizing spelling variations and minor orthographic differences, which are more prevalent in Emakhuwa translations. ChrF, on the other hand, being more sensitive to character-level $n$-grams, captures better the quality of translations. Nevertheless, further studies need to be done to assess the correlation of these automatic metrics with human evaluations.

**Using multiple references**   Notably, using multiple references improved scores for both BLEU and ChrF across all models. Specifically, BLEU scores increased by +0.24 to +0.54 on the *dev* set and by +0.3 on the *devtest* set.

## 7   Conclusion

In conclusion, this study expanded the FLORES+ evaluation set to include Emakhuwa, a low-resource language spoken in Mozambique. By translating the *dev* and *devtest* sets from Portuguese to Emakhuwa. We discussed key challenges such as spelling inconsistencies and loanword adaptations, which are prevalent due to Emakhuwa's underdeveloped spelling standards. Our rigorous methodology, involving translation, post-editing, and validation, ensured high-quality datasets used to benchmark neural MT models. The results indicate that incorporating multiple reference translations can enhance translation quality, particularly in languages with underdeveloped orthographies such as Emakhuwa. The dataset is publicly available, providing a valuable resource for future research in low-resource language MT.

## References

David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruiter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajuddeen Gwadabe, Freshia Sackey, Bonaventure F. P. Dossou, Chris Emezue, Colin Leong, Michael Beukman, Shamsuddeen Muhammad, Guyo Jarso, Oreen Yousuf, Andre Niyongabo Rubungo, Gilles Hacheme, Eric Peter Wairagala, Muhammad Umair Nasir, Benjamin Ajibade, Tunde Ajayi, Yvonne Gitau, Jade Abbott, Mohamed Ahmed, Millicent Ochieng, Anuoluwapo Aremu, Perez Ogayo, Jonathan Mukiibi, Fatoumata Ouoba Kabore, Godson Kalipe, Derguene Mbaye, Allahsera Auguste Tapo, Victoire Memdjokam Koagne, Edwin Munkoh-Buabeng, Valencia Wagner, Idris Abdulmumin, Ayodele Awokoya, Happy Buzaaba, Blessing Sibanda, Andiswa Bukula, and Sam Manthalu. 2022. A few thousand translations go a long way! leveraging pre-trained models for African news translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3053–3070, Seattle, United States. Association for Computational Linguistics.

AI4Bharat, Jay Gala, Pranjal A. Chitale, Raghavan AK, Sumanth Doddapaneni, Varun Gumma, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages.

Felermino D. M. A. Ali, Andrew Caines, and Jaimito L. A. Malavi. 2021. Towards a parallel corpus

of portuguese and the bantu language emakhuwa of mozambique.

Felermino Dario Mario Ali, Henrique Lopes Cardoso, and Rui Sousa-Silva. 2024. Detecting loanwords in emakhuwa: An extremely low-resource Bantu language exhibiting significant borrowing from Portuguese. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4750–4759, Torino, Italia. ELRA and ICCL.

Moussa Doumbouya, Baba Mamadi Diané, Solo Farabado Cissé, Djibrila Diané, Abdoulaye Sow, Séré Moussa Doumbouya, Daouda Bangoura, Fodé Moriba Bayo, Ibrahima Sory Conde, Kalo Mory Diané, Chris Piech, and Christopher Manning. 2023. Machine translation for nko: Tools, corpora, and baseline results. In *Proceedings of the Eighth Conference on Machine Translation*, pages 312–343, Singapore. Association for Computational Linguistics.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Michael Auli, and Armand Joulin. 2021. Beyond english-centric multilingual machine translation. *Journal of Machine Learning Research*, 22(107):1–48.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzman, and Angela Fan. 2021. The flores-101 evaluation benchmark for low-resource and multilingual machine translation.

Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc'Aurelio Ranzato. 2019. The FLORES evaluation datasets for low-resource machine translation: Nepali–English and Sinhala–English. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6098–6111, Hong Kong, China. Association for Computational Linguistics.

Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980.

Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. OpenNMT: Open-source toolkit for neural machine translation. In *Proceedings of ACL 2017, System Demonstrations*, pages 67–72, Vancouver, Canada. Association for Computational Linguistics.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi Stella,

Ankur Bapna, and Orhan Firat. 2024. Madlad-400: a multilingual and document-level large audited dataset. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.

R. de Moçambique E.P. 2016. Glossários de conceitos políticos, desportivos e sociais (português-línguas moçambicanas). Retrieved from `http://197.249.65.29/moodle/file.php/1/Glosario_RMe.pdf`.

Armindo Ngunga and Osvaldo Faquir. 2014. *Padronização da Ortografia de Línguas Moçambicanas: Relatório do VI Seminário*. Centro de Estudos das Línguas Moçambicanas.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

NLLBTeam, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2024. Scaling neural machine translation to 200 languages. *Nature*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on*

*Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Jiayi Wang, David Adelani, Sweta Agrawal, Marek Masiak, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Xuanli He, Sofia Bourhim, Andiswa Bukula, Muhidin Mohamed, Temitayo Olatoye, Tosin Adewumi, Hamam Mokayed, Christine Mwase, Wangui Kimotho, Foutse Yuehgoh, Anuoluwapo Aremu, Jessica Ojo, Shamsuddeen Muhammad, Salomey Osei, Abdul-Hakeem Omotayo, Chiamaka Chukwuneke, Perez Ogayo, Oumaima Hourrane, Salma El Anigri, Lolwethu Ndolela, Thabiso Mangwana, Shafie Mohamed, Hassan Ayinde, Oluwabusayo Awoyomi, Lama Alkhaled, Sana Al-azzawi, Naome Etori, Millicent Ochieng, Clemencia Siro, Njoroge Kiragu, Eric Muchiri, Wangari Kimotho, Toadoum Sari Sakayo, Lyse Naomi Wamba, Daud Abolade, Simbiat Ajao, Iyanuoluwa Shode, Ricky Macharm, Ruqayya Iro, Saheed Abdullahi, Stephen Moore, Bernard Opoku, Zainab Akinjobi, Abeeb Afolabi, Nnaemeka Obiefuna, Onyekachi Ogbu, Sam Ochieng', Verrah Otiende, Chinedu Mbonu, Yao Lu, and Pontus Stenetorp. 2024. AfriMTE and AfriCOMET: Enhancing COMET to embrace underresourced African languages. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5997–6023, Mexico City, Mexico. Association for Computational Linguistics.

Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. ByT5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

| Name | Role | Tasks | Expertise | Alias |
|------|------|-------|-----------|-------|
| Araibo Suhamihe | Translator | Translate *devtest*, Revise *dev* | Professional experience | Translator1 |
| Salustiano Eurico Ramos | Translator | Translate *dev*, Revise *devtest* | Professional experience | Translator2 |
| Gito Anastácio Anastácio | Evaluator | Evaluate and post-edit *devtest* / *dev* | Professional experience | Annotator1 |
| Júlio José Paulo | Evaluator | Evaluate and post-edit *devtest* / *dev* | Professional experience | Annotator2 |
| Vasco André António | Evaluator | Evaluate and post-edit *devtest* / *dev* | Professional experience | Annotator3 |

Table 6: Translation Team



Figure 7: Annotation Tool User Interface.

| | | |
|---|---|---|
| *pt → vmw* | | |
| **Source** | *pt* | **A camada é mais fina debaixo dos mares e mais espessa abaixo das montanhas**. |
| Translation | *en* | *It is thinner under the maria and thicker under the highlands.* |
| **References (*vmw*)** | A | Mpattapatthaaya tiwoyeva vathi wa mphareya ni yowoneya vathi wa miyaako. |
| | B | Mpattapatthaaya ti'yottettheeya othi wa iphareya ni yookhoomala vathi wa miyaako. |
| | **baseline** | Khalai atthu yahikhotta vathi-va, khukelela vasulu vaya. |
| | **afri-byT5** | Okhala wira okathi wa okathi ole ti wootepexa ottuli wa iphareya ni otepexa ottuli wa miyako. |
| | **afri-mT5** | Nthowa nenlo ninkhala ntoko nsuwa ntoko nsuwa ni ninkhala ntoko nsuwa ni ninkhala ntoko nsuwa ni ninkhala ntoko nsuwa. |
| **Systems (*vmw*)** | **byT5** | Ekamada eyo yootepa omalela vathi va iphareya ni yootepa omalela vathi va miyaako. |
| | **mT0** | Okhala wira ematta eyo enniphwanyaneya ottuli wa maasi, nto ematta eyo enniphwanyaneya ottuli wa maasi. |
| | **mT5** | Ekatana eyo ti yootepa otthuneya ovikana maasi ni yootepa otthuneya ovikana maasi. |
| | **M2M100** | Ekaaxa ele ti yootepa otthuneya vathi va ephareya ni yootepa otthuneya vathi va mwaako. |
| | **NLLB** | Mukattelo ti woorekama vathi vathi wa ophareya ni wootepa maasi vathi wa miyaako. |
| *vmw → pt* | | |
| **Source** | *vmw* | **Mpattapatthaaya tiwoyeva vathi wa mphareya ni yowoneya vathi wa miyaako.** |
| **References** | *pt* | A camada é mais fina debaixo dos mares e mais espessa abaixo das montanhas. |
| | **baseline** | A sua <unk> ainda é a propriedade que existe no <unk> sistema de coisas <unk>. |
| | **afri-byT5** | A sua aliança é pequena sobre o mar e visível das montanhas. |
| | **afri-mT5** | A sua vantagem é pequena sobre o mar e pequena sobre os oceanos. |
| | **byT5** | O amigo é pequeno sobre o mar e visível sobre os montes. |
| **Systems (*pt*)** | **mT0** | O companheiro é pequeno na água e pequeno na água. |
| | **mT5** | O seu amigo é pequeno na água e pequeno na água. |
| | **M2M100** | A arca é pequena debaixo do mar e visível debaixo das montanhas. |
| | **NLLB** | A bacia é barata no fundo do mar e muito clara no fundo das margens. |

| | | |
|---|---|---|
| *pt → vmw* | | |
| **Source** | *pt* | **Todos os cidadãos da cidade do Vaticano são católicos romanos.** |
| Translation | *en* | *All citizens of Vatican City are Roman Catholic.* |
| **References (*vmw*)** | A | Atthu ootheene opooma wo Vatikaanu anatiini a ekirixitawu ya katolika. |
| | B | Atthu ootheene opooma ya oVatikaanu anatiini a ekirixitawu katolika. |
| | **baseline** | Anammuttettheni otheene a epooma ya Vatoolika aari aRoma. |
| | **afri-byT5** | Atthu otheene a epooma ya oVaticano ti makatooliku a oRoma. |
| | **afri-mT5** | Otheene a epooma ya Vatikaano ti maKatoliko a oRoma. Otheene atthu otheene a epooma ya Vatikaano ti maKatoliko romano. |
| **Systems (*vmw*)** | **byT5** | Atthu otheene a epooma ya oVatikano ti makatooliku a oRoma. |
| | **mT0** | Atthu otheene a epooma ya oVaticano ti maKatolika a oRoma. |
| | **mT5** | Atthu otheene a epooma ya oVaticano ari maKristau a oRoma. |
| | **M2M100** | Atthu otheene a epooma ya oVaticano ari maKatoolika a oRoma. |
| | **NLLB** | Atthu ootheene anikhala epooma ya Vatikaano ti makatooliku a orooma. |
| *vmw → pt* | | |
| **Source** | *vmw* | **Atthu ootheene opooma wo Vatikaanu anatiini a ekirixitawu ya katolika.** |
| **References** | *pt* | Todos os cidadãos da cidade do Vaticano são católicos romanos. |
| | **baseline** | Todos na cidade do Vaticano apela a terra de <unk>. |
| | **afri-byT5** | Toda a população na cidade do Vaticano realiza a religião católica. |
| | **afri-mT5** | Todos os cidadãos em Vaticano são religiosos da igreja católica. |
| | **byT5** | Toda a população na cidade do Vaticano é religiosa da cristã católica. |
| **Systems (*pt*)** | **mT0** | Todos os cidadãos da cidade de Vaticane são cristãos da igreja católica. |
| | **mT5** | Todos os cidadãos na cidade de Vaticano são cristãos católicos. |
| | **M2M100** | Todos os cidadãos do Vaticano são cristãos católicos. |
| | **NLLB** | Todos na cidade do Vaticano são religiosos católicos. |

Table 7: Example of source-reference sentences pairs from *devtest* and outputs from translating source text using models discussed in Section 5.2

Figure 8: Matecat User Interface



Figure 9: Matecat Quality Report after revision, categorized by the following translation issue typologies: 1) Style (readability, consistent style, and tone); 2) Tag issues (mismatches, whitespaces); 3) Translation errors (mistranslation, additions or omissions); 4) Terminology and translation consistency; 5) Language quality (grammar, punctuation, spelling). The error point count corresponds to the number of segments found with any of the issues described above.

| 2592086107 | Entre os acusados estão dois vice-presidentes da Fifa, o uruguaio Eugenio Figueredo e Jeffrey Webb, das Ilhas Caimão e que é também presidente da Concacaf (Confederação de Futebol da América do Norte, Central e Caraíbas), assim como o paraguaio Nicolás Leoz, ex-presidente da Confederação da América do Sul (Conmebol), o antigo presidente da Confederação Brasileira de Futebol José María Marín, o membro do comité da Fifa para os Jogos Olímpicos Rio2016, o costarriquenho Eduardo Li, membro do comité executivo da Fifa, e Jack Warner, de Trinidad e Tobago, antigo vice-presidente do organismo e ex-presidente da Concacaf. |
|---|---|
| Por favor reveja Maiusculos no texto traduzido " <br><br> Reveja! Os números no texto traduzidos não batem com o do texto inicial" | Eriyari ya alokohiwa aakhala anli ale ari vathi wa ahooleli a Fiifa , urukwayiyu Eugenio Figueredo ni Jeffrey Webb , wa Ilya Caimão ni ori muhooleli wa Konkakafi ( Konfeterasawu ya mphira wa Ameerika ya Norote , Senterale ni Karayipa ) , siiso ntoko parakwayiyo Nicolás Leoz , muhooleli ohinhye eKonfeterasawu wa Ameerika ya Suuli ( Konmebol ) , muhooleli a khalayi wa Konfeterasawu Parasileyira ya mphira José María Marín , mempuru wa komite axitokweene a Fiifa wa Isepwere Olimpiku Riyu 20216 , koxitarikenyu Eduardo Li , mempuru wa komite exekutivo wa Fiifa , ni Jack Warner , wa Trinidad ni Tobago , muhooleli vathi wa muhooleli a khalayi wa mutthukumano wa muhooleli ohinhye wa Konkakafi . <br><br> **Empréstimos/Adptações anotados:** <br><br> América do Norte: Ameerika ya Norote   Central e Caraíbas: Senterale ni Karayipa   membro: mempuru   Confederação : Konfeteresawu <br> comité da Fifa: komite wa Fiifa   paraguaio: parakwayiyo   costarriquenho: koxitarikenyu   Confederação de Futebol: Konfeterasawu ya mphira <br> América do Sul: Ameerika ya Suuli   comité: komite   Concacaf: Konkakafi   Fifa: Fiifa   Confederação Brasileira: Konfeterasawu Parasileyira <br> Olímpicos Rio: Olimpiku Riyu |

Figure 10: Screenshot of a spelling report. The report is organized into two columns: the first column lists the segment ID along with any potential translation issues (i.e., punctuations, source-target length ratio flag, number mismatch, loanwords not annotated, case mismatch, etc.). The second column displays the source text and its translation. Potential misspellings are highlighted within the translation. In the translation, potential misspellings are highlighted in yellow and red—yellow indicating that suggestions for corrections are available and red indicating that no suggestions exist. Additionally, the report lists all words that translators have annotated as loanwords from Portuguese, using the format *<donor sequence in Portuguese>:<recipient sequence in Emakhuwa>*

591

| Model | Size | Hyperparameters |
|---|---|---|
| byT5-base / afri-byT5-base | 580M | <ul><li>Max source length: 200</li><li>Max target length: 200</li><li>Batch size: 8</li><li>Beams: 4</li></ul> |
| mT5-base / afri-mT5-base | 580M | <ul><li>Max source length: 200</li><li>Max target length: 200</li><li>Batch size: 8</li><li>Beams: 4</li></ul> |
| mT0 | 580M | <ul><li>Max source length: 200</li><li>Max target length: 200</li><li>Batch size: 8</li><li>Beams: 4</li></ul> |
| NLLB-200-distilled-600M | 600M | <ul><li>Max steps: 60000</li></ul> |
| M2M100 | 418M | <ul><li>Max tokens: 1200</li><li>Layers: 12</li><li>Dropout: 0.3</li><li>Attention dropout: 0.1</li><li>Learning rate: 3e-05</li><li>Max update: 40000</li><li>Emakhuwa was mapped to Swahili (sw)</li></ul> |

Table 8: MT Models Configurations

# Enhancing Tuvan Language Resources through the FLORES Dataset

**Ali Kuzhuget**
Senior Member,
IEEE #100062617
iOS Developer
agisight@gmail.com

**Airana Mongush**
Machine Learning Founder
Algebras AI
aira@algebras.ai

**Nachyn-E. Oorzhak**
Data Scientist
moy.kot.tas.ool@gmail.com

## Abstract

FLORES is a benchmark dataset designed for evaluating machine translation systems, particularly for low-resource languages. This paper, conducted as a part of Open Language Data Initiative (OLDI) shared task, presents our contribution to expanding the FLORES dataset with high-quality translations from Russian to Tuvan, an endangered Turkic language. Our approach combined the linguistic expertise of native speakers to ensure both accuracy and cultural relevance in the translations. This project represents a significant step forward in supporting Tuvan as a low-resource language in the realm of natural language processing (NLP) and machine translation (MT).

## 1 Introduction

Tuvan is a Turkic language, written using the Cyrillic alphabet and spoken by approximately 258,000 people (as of 2020), according to Ethnologue (2024). It is one of two official languages, along with Russian, of the Republic of Tuva, which is located in South Central Siberia, Russia. Despite its historical and cultural significance, Tuvan is classified as vulnerable by UNESCO, making it a critical target for preservation and technological integration. The FLORES (Goyal et al., 2022), spearheaded by Meta, aims to enhance machine translation systems by providing high-quality, controlled datasets for under-resourced languages for evaluation purposes. This paper, as a part of OLDI shared task (Initiative, 2024b), details our efforts to contribute to this dataset by providing translations from Russian to Tuvan.

## 2 Related work

It is essential to provide a brief overview of the FLORES (Goyal et al., 2022) dataset for those unfamiliar with this resource. The FLORES dataset, introduced by Goyal et al. (2022), is a benchmark for evaluating machine translation models on low-resource languages, which was translated to over 200 languages. It is comprised of two sets: the

dev set contains 997 sentences and the devtest set includes 1012 sentences, that were sampled from Wikinews, Wikijunior and Wikivoyage. FLORES dataset is crucial in advancing NLP for languages like Tuvan by providing benchmark specifically designed to evaluate machine translation systems across a wide variety of languages. The NLLB project (NLLB Team et al., 2022) further exemplifies efforts to scale human-centered machine translation across diverse languages.

## 3 Language overview

### 3.1 Handling dialectal differences

As the Republic of Tuva is a federal subject of Russia, the majority of the population is bilingual, speaking both Russian and Tuvan. For that reason, Tuvan translators relied on the FLORES dataset in Russian when developing one in Tuvan. During the translation process, any variations in interpretation due to dialectal differences were resolved by defaulting to the Central dialect's interpretation. This approach ensured uniformity and consistency across the dataset, which is crucial for training machine translation models that need to generalize well across different contexts.

The language is taught in schools optionally, but one can get a higher education in the Tuvan language and Literature in one of the universities of the Republic. Although the number of youth that speaks the language fluently decreases, it is still widely used both in cities and rural areas.

### 3.2 Linguistic challenges

Tuvan is characterized by its complex phonological and grammatical structures, including vowel harmony and extensive use of suffixes. These features posed challenges in translation, particularly in ensuring that the meaning and tone of the original Russian texts were accurately conveyed in Tuvan. However, since our work was limited to written form, we did not address challenges related to vocal translation or spoken dialects.

593

## 4 Data collection

### 4.1 Expertise of translators

Our translation team for the FLORES dataset was led by four native Tuvan speakers who are also proficient in Russian. The team included professional linguists and language enthusiasts with formal education in the Tuvan language. Although Tuvan language education is currently facultative, several of our translators had attended schools where Tuvan was the primary medium of instruction. This deep linguistic knowledge was crucial in ensuring that translations were not only accurate but also culturally relevant and sensitive. The following team of translators put their utmost effort to make the FLORES dataset available in Tuvan.

- Mongush Salim (Моңгуш Салим)
- Oorzhak Lyudmila (Ооржак Людмила)
- Ongai-ool Choduraa (Оңгай-оол Чодураа)
- Kuzhuget Ali (Кужугет Али)

### 4.2 Translation guidelines and training

Before beginning the translation tasks, all translators were provided with comprehensive guidelines, prepared in Russian, detailing the translation process. These guidelines, which are included in Appendix A, covered key aspects such as:

1. Maintaining the tone and style of the original text.
2. Handling idiomatic expressions and culturally specific references.
3. Ensuring pragmatic accuracy, including the correct use of pronouns and proper nouns.

Translators were instructed to adhere strictly to these guidelines in order to ensure a high level of consistency and quality across the entire dataset.

### 4.3 Managing the translation workflow

The translation process was managed using a Telegram group, where tasks were assigned, and progress was tracked using project management tools. This system allowed for effective coordination among the translators and ensured that the project stayed on schedule. It is worth pointing out that the workflow included multiple rounds of review and feedback among the translators to check one another and refine the translations further.

## 5 Experimental validation

### 5.1 Contribution to the evaluation of translation models

Our work makes a significant impact on the evaluation part of the Tuvan translation models by creating a reliable benchmark for this task. As for the current model of September 2024 developed by our team, it was trained on the existing Tuvan-Russian corpus, which consisted of approximately 200,000 pairs (Kuzhuget and Choigan, 2024) of translations sourced mainly from Wikipedia and other early Tuvan language projects (Kuzhuget et al., 2023). By contributing to the FLORES (Goyal et al., 2022) dataset, we have provided a more structured and high-quality resource, developed and checked by professionals, that is better suited for evaluating machine translation models.

The new dataset allowed us to run experiments to compare the quality of the existing translation models, available in Tuvan (Claude Sonnet 3.5, Google Translate v2 API, tyvan.ru). The results of these experiments are demonstrated on the Table 1.

Claude Sonnet 3.5 shows the highest performance overall, with BLEU scores of 35.65 and 33.04 for the Tyv-Rus translation in dev and devtest, respectively. The ChrF2++ scores are also the highest, reflecting good contextual understanding of this model. Google Translate v2 performs well, particularly with Tyv-Rus, achieving a BLEU score of around 28 in both datasets with ChrF2++ scores being also strong. The tyvan.ru model tends to have lower scores and lags behind the two above mentioned models.

| Model | Dataset | Tyv-Rus | | Rus-Tyv | |
|---|---|---|---|---|---|
| | | BLEU | ChrF++ | BLEU | ChrF++ |
| tyvan.ru | dev | 16.94 | 43.94 | 13.12 | 45.92 |
| | devtest | 16.41 | 43.09 | 13.35 | 46.11 |
| Google Translate v2 | dev | 29.78 | 54.60 | 14.30 | 45.50 |
| | devtest | 27.16 | 52.87 | 15.58 | 46.18 |
| Claude Sonnet 3.5 | dev | 35.65 | 59.65 | 16.67 | 48.45 |
| | devtest | 33.04 | 57.41 | 17.09 | 49.08 |

Table 1: Scores of Russian-Tuvan translation models on the FLORES dataset.

### 5.2 Manual evaluation of machine translation with the FLORES dataset

#### Objective

The objective of the manual evaluation on the FLORES dataset is to evaluate and compare the perceived translation adequacy of three translation services: Google Translate v2, Claude Sonnet 3.5,

and tyvan.ru. For that purpose we asked five Tuvan native speakers (annotators) to assess the quality of the translations for adequacy by giving scores on a scale of 1 to 5, the higher the better.

## Data

31 sentences were translated from Russian to Tuvan using the above mentioned services. The data was taken from two sets of FLORES dataset:

- dev set: First 16 sentences from the dev set
- devtest set: Last 15 sentences from the devtest set

Annotators were given a table with the first column containing sentences in Russian and three other columns containing translated sentences to Tuvan by the services, they had no information of the service that provided the translation to make the evaluation unbiased.

## Results

Google Translate v2

- Median score: $\sim 4$
- Distribution: Skewed toward higher scores (4-5), with some variability and occasional lower ratings.

Claude Sonnet 3.5

- Median score: $\sim 4$
- Distribution: Consistently high (3-5), with some sentences receiving lower scores (down to 2).

tyvan.ru

- Median score: $\sim 3 - 4$
- Distribution: Most variability, with a wider range of scores (1-5). This service had the most mixed feedback, with some high and many low ratings.

## Key Insights

Google Translate v2 and Claude Sonnet 3.5 performed relatively well, with consistent high scores. tyvan.ru showed more variability in performance, reflecting either inconsistent translation quality or differences in how annotators perceived adequacy.

## Visuals

Box plots (Figure 1) were used to compare the distribution of scores across services in both dev and devtest sets. A histogram (Figure 2) illustrated the overall distribution of scores for each service, with distinct colors for easy comparison (yellow for Google Translate v2, blue for Claude Sonnet 3.5,

red for tyvan.ru). This experiment highlights differences in translation quality as perceived by human annotators, showing that while Google Translate v2 and Claude Sonnet 3.5 generally perform well, tyvan.ru's performance was less consistent across annotators.



Figure 1: Direct assessment adequacy scores per Annotator



Figure 2: Perceived translation score distribution by Service

## 6 Data sample

As far as the the FLORES dataset in Tuvan is concerned, it has the following characteristics:

- dev set: 997 sentences, 18.26 average number of words in a sentence;
- devtest set: 1012 sentences, 18.44 average number of words in a sentence.

A table 2 showcases the examples of translated sentences from Russian to Tuvan. English translations were provided as examples so that English speakers could understand the general meaning of these sentences. But it is important to highlight that our translators were based only on the Russian

| No. | English source | Russian translation | Tuvan translation |
| --- | --- | --- | --- |
| 1 | Tokyo will be the only city in Asia to have hosted the Summer Olympics twice. | Токио станет единственным городом Азии, который дважды принимал летние Олимпийские игры. | Токио Азияга чайгы олимпийжи оюннарны ийи катап эрттирген чаңгыс хоорай болуп арттып каар. |
| 2 | There were no reports of serious damage or casualties in Tonga, but there was a temporary power outage, which reportedly prevented the authorities from receiving the tsunami warning sent by the Pacific Tsunami Warning Center (PTWC). | Из Тонга не поступило сообщений о серьезных разрушениях или о пострадавших, но произошло временное отключение электроснабжения, что, по имеющимся данным, не позволило властям Тонга получить предупреждение о цунами, посланное Тихоокеанским центром предупреждения о цунами (PTWC). | Тонгадан шыңгыы үрегдээшкиннер азы когарааннар дугайында медээлер келбээн, ынчалза-даа электри хандырылгазы түр када хже берген, ол чүүл, амгы үеде бар медээлер-биле, Цунами дыңнадыр Ооожум-океанчы төптен (PTWC) цунами дугайында дыңнадыгны алыр арганы Тонганың чазаанга бербээн болуп турар. |
| 3 | The percentage of people with multidrug-resistant tuberculosis in the overall group of tuberculosis patients still appears low; 6000 out of 330,000 infected in South Africa. | Однако, процент людей с туберкулёзом с множественной лекарственной устойчивостью в целой группе больных туберкулезом все еще кажется низким; 6000 от общего числа 330 000 зараженных в ЮАР. | Туберкулез аарыг бөлүкте хөй санныг эмнер-биле эмнеттинмес туберкулезтуг кижилерниң хуузу ам-даа эвээш ышкаш сагындырар; амгы үеде ЮАР-да аарыг 330 000 кижиниң ниити санындан 6000 кижи. |
| 4 | In Japan, the first celebrations of cherry blossom viewing were arranged by the emperor only for himself and other members of the aristocracy. | В Японии первые празднования цветения сакуры устраивались императором только для себя и других членов аристократии при императорском дворе. | Японияга сакура частырының баштайгы байырлалдарын императорнуң чүгле бодунга болгаш императорнуң чанында өске-даа аристократчы кежигүннерге дээш эрттирип турган. |
| 5 | The airlines offering these services include: Air Canada, Delta Air Lines, Lufthansa - for flights departing from the USA or Canada, and WestJet. | Авиакомпании, предлагающие эти услуги включают: Air Canada, Delta Air Lines, Lufthansa — для рейсов, отправляющихся из США или Канады, и WestJet. | Ол ачы-дузаны чедирип турар авиакомпанияларже Air Canada, Delta Air Lines, Lufthansa – АКШ-тан азы Канададан чоруп турар рейстерге, база WestJet олар хамааржыр. |

Table 2: Examples of translated sentences from English to Russian and Tuvan.

FLORES dataset and did not consider the English part during the development of the FLORES in Tuvan.

## 7 tyvan.ru and language preservation

The development of a Tuvan AI translator has been instrumental in preserving and revitalizing the Tuvan language, classified as endangered by UNESCO. Born out of a personal commitment to language preservation, the project began as a response to the lack of online translation resources for Tuvan. The initiative gained momentum with contributions from David Dale, who recognized the urgent need for language preservation during his work on the Erzya language (Dale, 2022), and Ali Kuzhuget, who has spent over a decade developing Tuvan language resources (Kuzhuget and Choigan, 2024), including dictionaries, keyboards, and translations of major platforms.

tyvan.ru serves as the online hub for these efforts, offering a range of Tuvan language tools, including the AI translator. Since its launch, the translator has been used by over 80,000 individuals, showcasing the growing interest in and need for Tuvan language resources. The project also extends its impact through the "One Code - Different Languages" volunteer initiative, which supports other vulnerable low-resource languages, such as Bashkir, Tatar, Chuvash, and Mari.

## 8 Limitations

The main challenge that occurs when dealing with low-resource languages like Tuvan is the small number of professionals, who could correctly translate to a low-resource language, abiding by all the grammar rules and nuances of the language. This results in a rather long translation process. Another project we are engaged in is the Seed (Maillard et al., 2023) project in Tuvan. One of the key limitations we faced in the second project is the lack of Tuvan speakers who are bilingual in any language other than Russian. Due to the historical context of Tuvans being part of the Soviet Union and Russia, Tuvan speakers typically only speak Russian in addition to Tuvan. Consequently, our extended research group is required to translate the Seed dataset from English to Russian first, and then from Russian to Tuvan. This two-step translation process introduces additional complexity and potential for translation inaccuracies, but it is a necessary approach given the linguistic resources available.

## 9 Conclusion

Our contribution to the FLORES dataset (Goyal et al., 2022) represents a significant step forward in supporting Tuvan as a low-resource language in the field of natural language processing. By focusing on the Central dialect and leveraging human expertise, we have created a high-quality resource that will aid in the development of more accurate and culturally sensitive machine translation systems.

In addition, we are currently in the process of translating the Seed (Maillard et al., 2023) dataset from English to Russian, and subsequently from Russian to Tuvan. This effort further enhances the resources available for Tuvan, contributing to the development of multilingual datasets and promoting the digital presence of the language.

This work not only enhances the digital presence of the Tuvan language but also contributes to its preservation and promotion.

## 10 Acknowledgments

## References

David Dale. 2022. The first neural machine translation system for the Erzya language. In *Proceedings of the first workshop on NLP applications to field linguistics*, pages 45–53, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

David Dale. 2023. How to fine-tune a nllb-200 model for translating a new language, medium article. Accessed August 19, 2024.

Ethnologue. 2024. Ethnologue. Accessed August 19, 2024.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán,

and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Open Language Data Initiative. 2024a. Translation guidelines. Accessed August 19, 2024.

WMT24 Open Language Data Initiative. 2024b. The flores+ evaluation benchmark for multilingual machine translation. Accessed August 19, 2024.

Ali Kuzhuget and Ondar Choigan. 2024. Russian-tuvan parallel corpus, crowdsourcing organised by ali kuzhuget. Accessed August 19, 2024.

Ali Kuzhuget, Airana Mongush, and David Dale. 2023. The first Tyvan AI language project, tyvan.ru. Accessed August 19, 2024.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzman. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

# A Appendix. Translation guidelines

These guidelines were adapted and translated to Russian for the team of translators based on the OLDI translation guidelines (Initiative, 2024a).

**Version 1.01**
**Author:** Күжүгет А.А
**Date:** 2 марта 2024 г.

## A.1 Важное примечание:

Ваши переводы будут использоваться для обучения или оценки движков машинного перевода. Поэтому этот проект требует человеческого перевода.

## A.2 Общие рекомендации:

1. Контекст: Вы будете переводить предложения из разных источников. В некоторых случаях может быть предоставлена ссылка на исходный документ, чтобы дать вам больше контекста. Если она доступна, пожалуйста, обратитесь к ней.
2. Единицы измерения: Не переводите единицы измерения. Переводите их точно так, как указано в исходном содержании.
3. Сохранение тона: При переводе сохраняйте тон, используемый в исходном документе. Например, энциклопедический контент из источников вроде Википедии должен переводиться с использованием формального тона.
4. Плавность перевода: Предоставляйте плавные переводы, не отклоняясь слишком сильно от структуры исходного текста. Допускаются только необходимые изменения.
5. Точность: Не расширяйте или не заменяйте информацию по сравнению с тем, что присутствует в исходных документах. Не добавляйте никакой поясняющей или скобочной информации, определений и т.д.
6. Полнота перевода: Не игнорируйте любой значимый текст, который был в исходнике.
7. Выбор перевода: В случае нескольких возможных переводов, пожалуйста, выберите тот, который имеет наибольший смысл (например, для соответствия гендеру, культурной адаптации на целевом языке, уровня формальности и т.д.).

## A.3 Именованные сущности:

Именованные сущности - это люди, места, организации и т.д., которые обычно упоминаются с использованием собственного имени. Этот раздел содержит рекомендации о том, как обращаться с именованными сущностями:

1. Общепринятые названия: Если в целевом языке существует общепринятое название для именованной сущности, используйте его.
2. Транслитерация: Если общепринятое название отсутствует, используйте транслитерацию оригинального термина, если это возможно. Если транслитерация не будет широко понята в контексте, вы можете сохранить оригинальный термин.

## A.4 Идиоматические выражения:

Идиоматические выражения не должны переводиться дословно. Используйте эквивалентную идиому, если таковая существует. Если эквивалентная идиома отсутствует, используйте идиому схожего смысла. Если в целевом языке не существует похожих выражений, перефразируйте идиому так, чтобы значение было сохранено на целевом языке.

## A.5 Неоднозначные местоимения:

Когда переводимое местоимение является неоднозначным (например, может быть интерпретировано как он/она или его/ее), выбирайте гендерно-нейтральные местоимения (такие как они), если таковые существуют на целевом языке. Однако, когда местоимение в исходном тексте четко обозначено по гендеру, вы должны следовать исходному материалу и сохранять гендерную маркировку.

# Machine Translation Evaluation Benchmark for Wu Chinese: Workflow and Analysis

**Jay Hongjian Yu\***
University of Washington
`hjyu@uw.edu`

**Yiming Shi\***
East China Normal University
`10214507003@stu.ecnu.edu.cn`

**Zherui Zhou\***
Shanghai Normal University
`zheruizhou@outlook.com`

**Christopher Haberland**
University of Washington
`haberc@uw.edu`

## Abstract

Although the population of Wu speakers is the second largest among languages in China, it is a textually under-resourced language, creating significant challenges for building machine translation systems supporting Wu. In this paper, we describe our Wu Chinese contribution to the FLORES+[1] dataset to serve as a training corpus and evaluation benchmark for machine translation models and we demonstrate its orthographic compatibility with existing Wu data. Our contributions include: (1) an open-source, manually translated dataset, (2) full documentations on the process of dataset creation and validation experiments, (3) preliminary tools for Wu Chinese normalization and segmentation, and (4) benefits and limitations of our dataset, as well as implications for other under-resourced languages. The project codes are stored on Github.[2]

## 1 Introduction

Wu Chinese is a Sinitic language spoken in Shanghai, Zhejiang, parts of Jiangsu, Anhui, and Jiangxi of China. It represents a complex and internally divergent dialect group (Pan et al., 1991) with around 83 million speakers (Eberhard et al., 2024). Despite having a robust population of speakers, Wu Chinese has been facing a sharp decline in daily usage due to the promotion of Standard Chinese. Meanwhile, Wu Chinese lacks a widely-accepted writing system and is not commonly written by native speakers, which has relegated Wu to becoming

under-resourced with respect to text data. In this data-scarce context, machine translation of Wu is an extremely challenging task.

To assist in the development of language models in Wu Chinese, we construct a FLORES+ benchmark dataset for Wu machine translation models and conduct evaluations that validate its utility via a language identification task. FLORES+ is an extension of the initial FLORES-101 project by NLLB Team et al. (2022), aiming at expanding the coverage to more languages worldwide. FLORES features fully aligned data directly translated from English Wikimedia and is consequently ideal for multilingual translation systems. Currently, the FLORES+ benchmark covers 3 language varieties written in Hanzi: Mandarin Chinese (Standard Beijing), Mandarin Chinese (Taiwanese), and Yue Chinese (Hong Kong Cantonese). The addition of Wu Chinese would facilitate machine translation across these similar varieties.

After translating and proofreading the new Wu Chinese dataset, we validated its consistency and generalizability with respect to existing Mandarin, Wu, and Yue Wikimedia resources. We also devised measures to normalize and segment Wu Wikimedia data in order to enhance their fidelity and provide a standardized dataset upon which to conduct our evaluations. Finally, we discuss the results of the experiments and suggest future work on Wu Chinese.

## 2 Language Overview

### 2.1 Wu Sounds and Wu Writings

Wu Chinese is mutually unintelligible with other Sinitic languages such as Mandarin and Yue (Cantonese), but shares a common set of *Hanzi* (Chinese characters) with these varieties. A significant feature of Wu Chinese that differentiates it from other Sinitic languages is its three-way VOT contrast in the syllable-initial position and its glottalization of

---

*\*Equal contribution

[1]https://github.com/openlanguagedata/flores/
[2]https://github.com/HongjianYu/FLORES-WU

"checked tones" (Norman, 1988), inherited from Middle Chinese. A lot of Hanzi in Wu Chinese bear two ways of pronunciation: *Wendu* ("文读", literary reading) and *Baidu* ("白读", vernacular reading). Wendu is the borrowing of pronunciation from Northern and Jianghuai Mandarin dialects; Baidu is the indigenous pronunciation derived from antecedent tones. For instance, "学" reads [ɕɥæ] in Standard Mandarin; as for Wu, it reads [ɦoʔ] (Baidu) in "学堂" (the traditional word for "school") and [iɐʔ] (Wendu) in "学校" (the modern word for "school" borrowed from Mandarin). Although [ɕɥæ] and [iɐʔ] appear to be unrelated, the sound change is in fact systematically induced according to the phonotactic constraints of Wu Chinese. Wendu and Baidu can occur in one single word too, as "大学" ("university") reads [dɐ ɦoʔ] where [dɐ] is Wendu and [ɦoʔ] is Baidu.

With the above prerequisite knowledge, it is logical to believe that most syllabary utterances in Wu have had a Hanzi representation, since most Middle Chinese pronunciations have traceable documentations. However, through an evolution of sounds and lexicon over a thousand years, the *Benzi* ("本字", original character) of many sounds have been lost. To recover missing graphemes in Wu Chinese writings, 赵元任 (1956) proposed a guideline that for every Wu utterance:

1. If the Hanzi of the utterance is known, use that Hanzi;
2. Otherwise, use a known Hanzi of the same pronunciation in the target Wu variety.

This logic has constituted the overarching principle of modern Wu orthographies. In the second case where we have multiple Hanzi candidates, we adhere to the following rules based on 简明吴方言词典(闵家骥et al., 1986):

1. Use the Hanzi that historically appears in Ming-Qing literature;
2. When no historical usage is found, pick a Hanzi that best subscribes to the semantic meaning;
3. Otherwise, choose the Hanzi that most frequently appears in vulgar texts.

## 2.2 Dialectal Variations

In the past, the Suzhou dialect (a sub-dialect of Su-Hu-Jia within Northern Wu) has been the prestige form of Wu Chinese. Beginning from the late 20th century, the Shanghai dialect (spoken in the urban central area of Shanghai, also a sub-dialect of Su-

Hu-Jia) has served as the lingua franca of the surrounding regions (Chen and Gussenhoven, 2015) because of the city's population and economic importance. There have also been recent attempts by the community to create "Standard Wu", which closely relates to Taizhou, Shanghai, and Ningbo dialects[3]. That being said, none of the Wu varieties is influential enough to profoundly alter the pronunciations or lexicons of others, considering the hegemonic impact of Standard Mandarin.

As a result of multiple factors, Wu dialects have developed vastly divergent readings of the same Hanzi. This becomes especially problematic when no Hanzi is registered to an utterance, i.e. finding a Hanzi of the same pronunciation in the local dialect is necessary. To illustrate how complex the spelling variations can become, below is an *incomplete* list of the pronunciations and corresponding common spellings of the location/time preposition in 33 Wu dialects transcribed by 钱乃荣 (1992):

$$
\text{"Prep. of loc./time"} \begin{cases} \text{[ləʔ]} & \text{"勒", "了"} \\ \text{[lɐʔ]} & \text{"辣", "拉", "垃"} \\ \text{[liʔ]} & \text{"立"} \\ \text{[le], [læ]} & \text{"来"} \\ \text{[kɛ]} & \text{"该"} \\ \text{[tsɛ}^\text{e}\text{], [dʑɛ]} & \text{"在"} \end{cases}
$$

Inevitably, adopting one spelling standard means discarding the rest of equivalently common spellings. For instance, [lɐʔ lɐʔ] "辣辣" is the prevalent spelling of the location/time preposition (double syllable) in the Shanghai dialect, whereas [ləʔ lɐʔ] "勒拉" is commonly accepted in many other dialects and therefore more frequently encountered. We shall discuss the Chongming dialect which is the variety used in our dataset, and its corresponding orthography in the next section.

## 2.3 Relevant Resources

Before proceeding, we want to outline the resources available for natural language processing tasks related to Wu Chinese. The foundation of Wu Chinese studies was laid by 赵元任(Yuen Ren Chao) with his 现代吴语的研究(赵元任, 1956). Works by later scholars include 当代吴语研究(钱乃荣, 1992) and others. There is 简明吴方言词典(闵家骥et al., 1986), a dictionary that covers most Wu dialects with a light emphasis on the Shanghai dialect lexicon, as well as 上海话大词

---

[3]https://wuu.m.wikipedia.org/wiki/标准吴语/

典(钱乃荣, 2008) specifically for the Shanghai dialect. Thanks to community efforts, there is a Wu Wikipedia[4], a forum[5], and several online dictionaries[6][7] made by 吴语协会and 吴语学堂.

## 3 Methodology

The FLORES+ Wu dataset is directly translated from English into Wu Chinese. The target Wu variety is the Chongming dialect. The Chongming dialect (or more broadly the Shadi dialect) is a Wu dialect within the Su-Hu-Jia division, spoken in Chongming, Haimen and Qidong districts as well as in some areas of Zhangjiagang (张惠英, 2009). Although Chongming belongs to the Municipality of Shanghai, the dialect is distinctive from the urban variety on many aspects and is known for preserving many rare characteristics of Middle Chinese. There is also a large dedicated Chongming dictionary (张惠英et al., 2014), which unfortunately was not accessible to us during our efforts.

While the Chongming dialect is not the most used or researched Wu dialect, it is a representative one of Northern Wu along the dialect continuum, because it spreads in between Shanghai and Suzhou, where the two prestige forms of Northern Wu are used. Besides, the lexicon of the Chongming dialect has a significant overlap with other dialects for its preservation of archaic features common to most Wu dialects. In contrast, the Shanghai dialect is less generalizable to other Wu dialects due to its integration with other Chinese languages. However, the overlap between the Chongming dialect and Southern Wu dialects might be less prominent. As a result, our dataset may be less effective for Southern Wu linguistic tasks.

Data were equally distributed and translated by 2 native speakers of the Chongming dialect who have earned or are pursuing a university degree in English. Both translators grew up in Chongming; one went to Putuo, Shanghai for high school and college, and the other went to Fengxian, Shanghai for college. They mainly speak Wu at home, but also speak it with peers on occasion. They are exposed to the Shanghai urban dialect as well as other local varieties. All translated data were checked by a third independent Wu speaker.

The translators worked collaboratively on the task. They mainly used 简明吴方言词典(闵

---

[4]https://wuu.wikipedia.org/wiki/
[5]https://wu-chinese.com/bbs/
[6]http://wu-chinese.com/minidict/index.php
[7]https://www.wugniu.com/dict

家骥et al., 1986) and then 上海话大词典(钱乃荣, 2008) (see 2.3) if they were unable to recall the parallel Hanzi that represent the utterances. Despite the discrepancy between the Chongming and Shanghai lexicons, the dictionaries provided enough context to determine the appropriate orthography. For example, both "乃末" and "乃么" point to the same word [ne məʔ]; as both dictionaries only list "乃末", it was easy to make the choice. When the two dictionaries' orthographies diverged, 简明吴方言词典was prioritized. When there were phonetic distinctions between the Chongming and Shanghai dialect and the original character was indeterminable, we made sure that the selected Hanzi had aligned pronunciation. Noticeably, we did not use "勿" [vəʔ] but used "弗" [fəʔ] for the word of negation. We were also committed to maintaining a balanced language register, as the translated content is formal, though Wu Chinese is usually colloquial. We referred to the broadcasting-style Wu Chinese found in Shanghai and Suzhou to achieve the desired register. Beside the task of translation, the translators dedicated time to review Wu dictionary entries and online fora to grasp the overall construct of the written Wu landscape. During the proofreading process, when an alternative translation is suggested, the translator responsible for the line would ask for community guidance from the aforementioned fora. Proper wordings were selected according to the intuitive preferences of Wu native speakers from the community. In total, we have translated and verified the linguistic accuracy all 997 sentences in the dev set.

## 4 Data Samples

This sections lists the first 5 lines of translation along with their English counterparts.

1. 斯坦福医学院个科学家勒礼拜一公布一种可以按种类划分细胞个新个诊断家生个发明：一种可以用标准喷墨打印机大量生产，差弗多小到只有一美分一只个可印芯片。

   On Monday, scientists from the Stanford University School of Medicine announced the invention of a new diagnostic tool that can sort cells by type: a tiny printable chip that can be manufactured using standard inkjet printers for possibly about one U.S. cent each.

2. 首席研究者认为伊作兴可以勒低收入国家稍为早发现癌症、肺结核、艾滋病、疟疾个病人，伊拉个乳腺癌等疾病个治愈率是

只富裕个国家个一半。

Lead researchers say this may bring early detection of cancer, tuberculosis, HIV and malaria to patients in low-income countries, where the survival rates for illnesses such as breast cancer can be half those of richer countries.

3. JAS 39C 鹰狮战斗机勒当地辰光差弗多早晨九点半（中央时间两点半）撞向飞机跑道爆炸，葛商业航班个机场关闭。

The JAS 39C Gripen crashed onto a runway at around 9:30 am local time (0230 UTC) and exploded, closing the airport to commercial flights.

4. 瓣只飞行员拨认为是空军中队长迪罗克利特帕塔维（Dilokrit Pattavee）。

The pilot was identified as Squadron Leader Dilokrit Pattavee.

5. 当地个媒体报道，一只机场个救火车勒回应个辰光翻倒哉。

Local media reports an airport fire vehicle rolled over while responding.

# 5 Validation

## 5.1 Task

Language identification models have frequently been trained to corroborate translation datasets' correspondence with other texts in the target language. We train a language identification model on the FLORES+ Wu dataset to show its compatibility with other Wu Chinese text resources. The goal of our experimental setup is to train a model to correctly identify the language of an input sentence from among Mandarin, Wu, and Yue when prompted a sentence written in Hanzi.

Experiments are split into three parts. In each part, we trained three binary classification models: Mandarin-Wu, Mandarin-Yue, and Wu-Yue on their respective datasets, and a three-way classification model on all designated data. We recorded the performance of the models in terms of their accuracy on unseen data, collected separately by languages.

We broke down the experiments into distinct trials that reflect noteworthy characteristics of the training and evaluation data. In part 1, we conducted a 9:1 split on FLORES+ datasets to test their internal consistency. In part 2, we trained the model on Wikimedia data and tested on FLORES+ to showcase the compatibility of FLORES+ in its common use cases. In part 3, we reversed the training and testing data in part 2 to explore the extent of generalizability of FLORES+ given that it consists of small but parallel data.

## 5.2 Dataset Processing

We adopted two data sources, Wikimedia[8] and FLORES+.

We downloaded Wikimedia database XML dumps for Mandarin, Wu, and Yue. Since Mandarin and Yue dumps are significantly larger than Wu, only a portion of the data was extracted. After normalization, each dataset comprises 25,000 lines of texts.

FLORES+ datasets in use include the existing two Mandarin dev sets and Yue dev set, as well as the newly built Wu dev set. As a result, there are 1994 lines of sentences in Mandarin, 997 in Wu and 997 in Yue.

### 5.2.1 Normalization

The dumps were preprocessed with a simple filter removing Latin characters and template symbols. Because Mandarin and Yue Wikimedia were written in Traditional Chinese and Wu Wikimedia partially, we utilized OpenCC[9] which supports character-level and phrase-level conversion from Traditional to Simplified Chinese. OpenCC conversion was also applied to FLORES+ Taiwan Mandarin and Yue datasets.

For the lack of a standard orthography, Wu Wikimedia requires an additional step of normalization. For some characters and words that are often spelled differently, we replaced the constituents by our standard forms. However, the brute force method does not work for every character and word that need normalization. For example, "勒海" before normalization could be interpreted differently:

$$\text{勒海} \begin{cases} \text{勒嗨} & \text{"In there, over there"} \\ \text{勒海(浪)} & \text{"In the sea"} \end{cases}$$

The ambiguity of the language results in demands on more advanced normalization schemes, which are essential for language models to grasp semantic understandings.

### 5.2.2 Segmentation

Since Chinese languages do not depend on spaces to separate words, segmentation tools tailored to the respective languages are indispensable. For

---

[8]https://dumps.wikimedia.org/
[9]https://github.com/BYVoid/OpenCC/

| | cmn | wuu | yue | total |
|---|---|---|---|---|
| **cmn-wuu** | 0.995 | 0.990 | - | 0.993 |
| **cmn-yue** | 1.000 | - | 0.970 | 0.990 |
| **wuu-yue** | - | 0.990 | 0.990 | 0.990 |
| **c-w-y** | 0.995 | 0.990 | 0.979 | 0.990 |

Table 1: FastText classification precision rates by languages ($k$=1), trained with and evaluated on FLORES+ dev sets, 9:1 split. Rows represent in what languages the model is trained with; columns represent the language of the testing data. We use the ISO 639-3 codes for abbreviation: cmn (Mandarin), wuu (Wu), yue (Yue).

Mandarin, we used `jieba`[10], a popular open-source segmentation library with a prefix dictionary and a HMM-based model with Viterbi algorithm for unknown words; for Cantonese, we used `cantoseg`[11] which is built from `jieba` and reads in a merged corpus from `jieba` and `PyCantonese` (Lee et al., 2022).

As for Wu, We decided on adding an auxiliary dictionary to `jieba` for frequent words and phrases in Wu Chinese that are not present in Mandarin. We found this approach has also been used by Chen (2023). The entries in the auxiliary dictionary have been collected from 简明吴方言词典(闵家骥et al., 1986) for its wide coverage on Northern and Southern Wu dialects and 上海话大词典(钱乃荣, 2008) for its rich lexicon.

### 5.3 Model

We use fastText[12] text classification (Joulin et al., 2016) for all experiments. FastText is a CPU-based library for efficient learning of word representations and sentence classification. We tagged language labels to individual lines in every dataset and called the `supervised` command to train the models. When testing, fastText takes a $k$ parameter and returns both precision and recall rates within the first $k$ predicted labels. As only 2 or 3 distinct labels were present in each run, we used $k$=1 to compute the classification accuracy.

### 5.4 Results

In part 1, all four models demonstrate a high accuracy in classifying all languages. This validates the internal consistency of FLORES+ datasets including the new Wu Chinese addition.

[10]https://github.com/fxsjy/jieba/
[11]https://github.com/ayaka14732/cantoseg/
[12]https://fasttext.cc/

| | cmn | wuu | yue | total |
|---|---|---|---|---|
| **cmn-wuu** | 0.896 | 0.999 | - | 0.930 |
| **cmn-yue** | 0.968 | - | 0.987 | 0.975 |
| **wuu-yue** | - | 0.996 | 0.986 | 0.991 |
| **c-w-y** | 0.868 | 0.997 | 0.971 | 0.926 |

Table 2: FastText classification precision rates by languages ($k$=1), evaluated on FLORES+ dev sets, trained with Wikimedia texts. Rows and columns are the same as Table 1.

| | cmn | wuu | yue | total |
|---|---|---|---|---|
| **cmn-wuu** | 0.944 | 0.639 | - | 0.792 |
| **cmn-yue** | 0.949 | - | 0.796 | 0.872 |
| **wuu-yue** | - | 0.815 | 0.884 | 0.849 |
| **c-w-y** | 0.929 | 0.615 | 0.735 | 0.759 |

Table 3: FastText classification precision rates by languages ($k$=1), evaluated on Wikimedia texts, trained with FLORES+ dev sets. Rows and columns are the same as Table 1.

In part 2, we can observe a total accuracy over 90% for all four models. However, a drop in accuracy for the Mandarin-Wu model on Mandarin texts indicates false positives of Mandarin texts mislabelled as Wu. Alternatively speaking, training on Mandarin and Wu Wikimedia data allows the model to capture features of Wu and thus correctly label Wu data, but is less effective for recognizing Mandarin features.

In part 3, due to insufficient training data, the models exhibit tendency to misclassify Wu and Yue data as Mandarin. There is a more significant contraction in testing accuracy on Wu data than Yue. Meanwhile, the accuracy of Mandarin-Yue and Wu-Yue models is maintained at a relatively stable level.

Some typical misclassifications are listed below. These input data are Mandarin but mislabeled as Wu, presented in the segmented format. The corresponding Wu translations are provided as well. "cmn" denotes Mandarin data (mislabeled as Wu) and "wuu" denotes Wu data (correctly labeled).

1. (cmn) 人类 的 手 比脚 短 ， 指 （ 趾 ） 骨 更 直 。
   (wuu) 人个 手比脚 短 ， 趾骨 更 直 。
2. (cmn) 看到 有人 愿意 支持 我 ， 我 很 高兴 。
   (wuu) 我蛮 高兴 有人 愿意 支持 我个 。
3. (cmn) 越来越 多 超市 开始 提供 更 多样 化 的 即 食品 ， 部分 超市 甚至 提供

微波炉 或 其他 设备 以 食物 加热 。
(wuu) 超市 里个 现成 食品 种类 越来越
多 。 有眼 超市 甚至 提供 微波炉 或者
其他 方式 来 加热 食物 。

From these cases we can observe many common words in the two languages. There are nuanced differences in phrasing order and sentence structures but the presented orders and structures are generally acceptable in both languages and only subject to personal habits of the translators. Despite lexicon similarity, the model also seems to have difficulty in recognizing Wu constituents due to the relatively weak performance of the segmentation tool, evident in "我个" (1, 2), "我蛮" (2).

Overall, the FLORES+ Wu dataset is consistent and capable of evaluating models trained with common data after appropriate normalization and segmentation. However, its benchmarking quality might not be as good as the Mandarin and Yue dataset for several reasons.

The tokenization scheme could be further optimized with a better segmentation tool in use. The current manually configured word list for Wu in the auxiliary dictionary of `jieba` is relatively small compared to the pre-built Mandarin dictionary in `jieba` and the independently maintained Yue dictionary by Lee et al. (2022). As some syntactic structures of Wu have not been recognized by `jieba`, the models are unable to learn accurate representations of the constituents.

Although the spelling standard used in FLORES+ Wu dataset is relatively generalizable to other dialects, it fails to take account of many expressions in Southern Wu dialects which are a part of the Wikimedia data. Therefore, we suggest that the training and testing datasets should align in the range of dialects whenever possible.

Moreover, the tendency of Wu Chinese to be influenced by Mandarin poses problems, exemplified by our classifier mislabeling Wu data containing "了" as Mandarin, because this grammar particle is common in Mandarin but infrequent in older Wu data.

## 6 Conclusion

As for now, a contemporary, consistent, and organized corpus becomes crucial for high-quality Wu Chinese language models. Meanwhile, it is important for AI scientists and engineers to have an understanding in the linguistic properties of their training and testing data. We hope that our published dataset and code contributions provide a foundation for future efforts towards Wu Chinese machine translation and language modeling.

## References

Yiya Chen and Carlos Gussenhoven. 2015. Shanghai chinese. *Journal of the International Phonetic Association*, 45(3):321–337.

Yuanhao Chen. 2023. Improving tts for shanghainese: Addressing tone sandhi via word segmentation.

David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2024. Ethnologue: Languages of the world. twenty-seventh edition. dallas, texas: Sil international.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.

Jackson L. Lee, Litong Chen, Charles Lam, Chaak Ming Lau, and Tsz-Him Tsui. 2022. Pycantonese: Cantonese linguistics and nlp in python. *Proceedings of the 13th Language Resources and Evaluation Conference*.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Jerry Norman. 1988. *Chinese*. Cambridge University Press.

Wuyun Pan, S.F. Zhengzhang, R.J. You, and Lien Chinfa. 1991. An introduction to the wu dialects. *Journal of Chinese Linguistics Monograph Series*, (3):235–291.

张惠英. 2009. 崇明方言研究. 中国社会科学出版社.

张惠英, 顾晓东, and 王洪钟. 2014. 崇明方言大词典. 上海辞书出版社.

赵元任. 1956. 现代吴语的研究. 科学出版社.

钱乃荣. 1992. 当代吴语研究. 上海教育出版社.

钱乃荣, editor. 2008. 上海话大词典. 上海辞书出版社.

闵家骥, 范晓, 朱川, and 张嵩岳, editors. 1986. 简明吴方言词典. 上海辞书出版社.

# Open Language Data Initiative: Advancing Low-Resource Machine Translation for Karakalpak

**Mukhammadsaid Mamasaidov**
Tahrirchi
m.mamasaidov@tahrirchi.uz

**Abror Shopulatov**
Tahrirchi
a.shopolatov@tahrirchi.uz

## Abstract

This study presents several contributions for the Karakalpak language: a FLORES+ devtest dataset translated to Karakalpak, parallel corpora for Uzbek-Karakalpak, Russian-Karakalpak and English-Karakalpak of 100,000 pairs each and open-sourced fine-tuned neural models for translation across these languages. Our experiments compare different model variants and training approaches, demonstrating improvements over existing baselines. This work, conducted as part of the Open Language Data Initiative (OLDI) shared task, aims to advance machine translation capabilities for Karakalpak and contribute to expanding linguistic diversity in NLP technologies.

## 1 Introduction

The Karakalpak language, a member of the Turkic language family, is primarily spoken in the Republic of Karakalpakstan, an autonomous region within Uzbekistan, Central Asia. Current estimates suggest a native speaker population to be around 900,000 individuals (Ethnologue, 2024). Linguistically, Karakalpak is an agglutinative language which belongs to the Kipchak branch of the Turkic language family and shares close affinities with Kazakh and Nogai (Berdimuratov and Dáwletov, 1979).

As a low-resource language, Karakalpak presents significant challenges in the field of natural language processing, particularly in machine translation. Major translation platforms such as Google Translate (Google, 2024) currently do not offer support for this language as of the time of writing this paper, underscoring the need for dedicated research and development in this area.

This study, conducted as part of the Open Language Data Initiative (OLDI) shared task, presents fine-tuned neural models for Karakalpak translation, a fine-tuned version the No Language Left

Behind (NLLB) model (NLLB Team et al., 2022). In line with OLDI's goals of expanding language resources, we release several key contributions:

1. A FLORES+ devtest dataset (NLLB Team et al., 2022) translated to Karakalpak

2. Parallel corpora for Uzbek-Karakalpak, Russian-Karakalpak and English-Karakalpak of 100,000 pairs each [1]

3. Open-sourced fine-tuned neural models for translation across Uzbek, Russian, English and Karakalpak languages [2]

4. Scripts for Latin-Cyrillic transliteration for Karakalpak

Our research aims to advance the state of machine translation for Karakalpak, contributing to the broader OLDI objective of improving natural language processing capabilities for low-resource languages. This work demonstrates how shared tasks like OLDI can drive progress in expanding linguistic diversity in NLP technologies.

## 2 Related work

The field of machine translation for low-resource languages has experienced significant advancement with the advent of the No Language Left Behind (NLLB) model families. These innovative models demonstrate the capability to facilitate translation across more than 200 languages, leveraging extensive collections of online corpora.

Another notable multilingual translation model is MADLAD-400 (Kudugunta et al., 2024), which extends the capabilities of large language models to cover 400 languages, including many low-resource languages and Karakalpak. This model represents a significant step forward in expanding the reach of

---

[1] https://huggingface.co/datasets/tahrirchi/dilmash

[2] https://huggingface.co/collections/tahrirchi

| English | Karakalpak |
|---------|------------|
| According to Japan's nuclear agency, radioactive caesium and iodine has been identified at the plant. | Yaponiya yadro agentligi maǵlıwmatlarına kóre, stanciyada radioaktiv ceziy hám yod bar ekenligi anıqlanǵan. |
| The result of plotting analysis will be posted to a public website. | Syujet analiziniń nátiyjesi ǵalabalıq veb-saytqa jaylastırıladı. |
| The station's web site describes the show as "old school radio theater with a new and outrageous geeky spin!" | Stanciya veb-saytında show "jańa hám ádettegiden basqasha ájáyıp aylandıratuǵın eski mektep radio teatrı!" dep táriyiplenedi. |

Table 1: Examples from the FLORES+ dataset for English-Karakalpak language pair

machine translation to a broader range of linguistic communities.

In the specific context of Karakalpak machine translation, several notable efforts have been made. A prominent example is the Apertium platform (Forcada et al., 2011), a rule-based machine translation system designed for low-resource languages. Utilizing finite-state algebra and rule-based methodologies, Apertium has developed morphological analyzers and spell-checking tools for Karakalpak[3]. Furthermore, it has produced machine translation systems for language pairs for Uzbek-Karakalpak, Kazakh-Karakalpak, and Tatar-Karakalpak.

A recent contribution to the Karakalpak translation comes from the Turkic Interlingua (TIL) team (Mirzakhalov, 2021), who introduced a model specifically trained on Turkic languages and corpora, with Karakalpak included in its linguistic scope. This initiative not only enhances the translation capabilities for Karakalpak but also contributes to the broader landscape of Turkic language processing. Additionally, the team has made significant strides in corpus development, introducing parallel corpora for numerous Turkic language pairs, including those involving Karakalpak.

Moreover, a proprietary online translation service for Karakalpak exists at `https://from-to.uz/`. To provide a comprehensive evaluation of Karakalpak machine translation capabilities, we will assess this tool's performance using its API, comparing it with our proposed models. This comparison will offer insights into both open-source and commercial solutions for low-resource language translation.

To our best knowledge, these developments collectively represent important steps towards improving machine translation capabilities for Karakalpak and other low-resource languages within the Turkic language family.

As an additional benchmark, we will include Claude-3.5-sonnet, a commercial large language model (LLM) with multilingual capabilities. While not specifically designed for Karakalpak translation, Claude-3.5-sonnet represents the current state of general-purpose language models and can provide valuable insights into how well such models perform on low-resource language tasks.

## 3 Datasets

### 3.1 FLORES+ Devtest Dataset

This study introduces the Karakalpak FLORES+ devtest dataset, which comprises 1012 sentences translated from English to Karakalpak. The FLORES+ datasets, derived from Wikimedia content, have been widely employed in the evaluation of foundational models within the NLLB family.

This dataset was developed under the auspices of the Open Language Data Initiative (OLDI). Two annotators were responsible for the translation of a devtest split, with subsequent cross-verification to ensure accuracy. The Karakalpak translations adhere to the most recent iteration of the Latin script orthography (see Table 1).

The Karakalpak orthography has experienced multiple changes recently. Both Latin and Cyrillic scripts are utilized, with the Latin script, introduced in 1995, undergoing several revisions. Notable modifications occurred in 2009 and 2016, with the latter replacing digraphs with diacritic letters to overcome previous limitations. Conversion scripts for Cyrillic and older Latin versions to the current system are available on GitHub[4].

### 3.2 Training data

The training dataset comprises diverse parallel corpora sourced from multiple domains, including news articles, literary works, lexicographic resources, and educational materials. Specifically, the corpus encompasses on average across three languages:

---

[3] `https://github.com/apertium/apertium-uzb-kaa`

[4] `https://github.com/tahrirchi/kaa-scripts`

- 23% sentences from news sources

- 34% sentences from books (novels, non-fiction)

- 24% sentences from bilingual dictionaries

- 19% sentences from school textbooks

- 4,000 English-Karakalpak pairs from Gatitos Project (Jones et al., 2023)

In total, the dataset consists of 100,000 sentence pairs for Uzbek-Karakalpak, Russian-Karakalpak, and English-Karakalpak each, making 300,000 pairs in total. Since there were too few bitexts with English, we decided to create English-Karakalpak dataset by translating Russian sentences from the Russian-Karakalpak dataset to English using Claude 3.5 Sonnet (See Appendix A). To promote further research and development in this field, we have made these corpora publicly available.

### 3.3 Data Mining Process

For mining parallel sentences, we apply only local mining, when we are sure that parallel sentences are to be mined from the translations of the same book, document or article. For alignment, we use LaBSE embeddings, although Karakalpak is not a supported language in LaBSE. We found that due to similarities of Karakalpak to already included Uzbek and Kazakh languages, LaBSE performed well for aligning sentences so we skipped this step.

The sentence alignment method we use is identical to the one applied for Erzya, as described by (Dale, 2022). We utilize LaBSE (Language-agnostic BERT Sentence Embedding) (Feng et al., 2020) to generate embeddings for each sentence pair. To calculate the alignment score, we first determine the cosine similarity between these embeddings. We then adjust this similarity by multiplying it with a length ratio - specifically, the length of the shorter sentence divided by the length of the longer sentence.

Using dynamic programming, we identify the sequence of sentence pairs that maximizes the total similarity score. Finally, we apply a threshold to filter out low-scoring alignments.

## 4 Translation Experiments

### 4.1 Model Training

For our experiments, we utilized the nllb-200-distilled-600M model, which is a transformer-based neural machine translation model with an encoder-decoder architecture. This model comprises 12 layers and employs the following approach: the source and target languages are indicated by the first tokens of the encoder and decoder inputs, respectively. This architecture allows the model to process and translate between numerous language pairs. The training process for our experiments consisted of several key steps:

#### 4.1.1 Tokenizer Preparation

Initially, we trained a SentencePiece (Kudo and Richardson, 2018) tokenizer on an expanded mono-corpus of approximately 300,000 Karakalpak sentences with a total of 16,000 vocabulary length. We decided to train a separate tokenizer because we hypothesized that the intial vocabulary of the NLLB model was not suited for Karakalpak, as there were some non-ASCII characters in the writing system (see Table 2). We also provide an evaluation of a model without training a separate tokenizer and compare the model's performance with and without additional trained tokens.

$$ Á\ á \quad Ǵ\ ǵ \quad Í\ ı \quad Ń\ ń \quad Ó\ ó \quad Ú\ ú $$

Table 2: Non-ASCII letters from Karakalpak Latin alphabet.

#### 4.1.2 Vocabulary Expansion

Following tokenizer training, we augmented the model's vocabulary. This expansion resulted in a total of 269,399 tokens, representing an increase of 13,195 tokens from the original model configuration. We then resized the model's token embeddings and initialized the new embeddings by averaging the embeddings of their constituent subtokens from the original vocabulary.

#### 4.1.3 Model Variants

We developed three distinct model variants to evaluate the impact of additional tokens and training data composition:

1. **dilmash-raw**[5]: This model was trained exclusively on a our own parallel corpus comprising 300,000 sentence pairs in Uzbek, Russian, and English on the original nllb-200-600M.

2. **dilmash**: Same as **dilmash-raw**, but it is a fine-tuned model with additional tokens which

---

[5]dilmash [dil-mash] *n. (from Karakalpak)* an oral interpreter

| Model | en-kaa | ru-kaa | uz-kaa | kaa-en | kaa-ru | kaa-uz |
|---|---|---|---|---|---|---|
| madlad-400 | 2.68 / 22.48 | 2.01 / 19.93 | 1.31 / 16.81 | 28.42 / 53.06 | 16.95 / 41.12 | 10.26 / 38.75 |
| apertium-uzb-kaa | - | - | 12.26 / 42.27 | - | - | 5.61 / 35.82 |
| google-from-kaz | - | - | - | 20.95 / 44.63 | 13.55 / 36.91 | - |
| google-from-uzb | - | - | - | 21.40 / 45.50 | 13.78 / 37.64 | - |
| nllb-200-600M-from-kaz | - | - | - | 4.32 / 23.35 | 3.12 / 16.86 | 3.91 / 25.26 |
| nllb-200-600M-from-uzb | - | - | - | 8.89 / 32.26 | 5.82 / 26.33 | 4.83 / 29.68 |
| from-to.uz | - | - | 20.18 / 53.22 | - | - | 11.18 / 41.37 |
| claude-3.5-sonnet | 11.17 / 33.37 | 9.02 / 34.02 | 12.74 / 35.17 | **37.06 / 61.41** | **25.70 / 51.23** | **22.38 / 54.71** |
| dilmash-raw | 14.37 / **45.65** | 11.41 / **42.99** | 16.16 / 48.88 | 30.01 / 54.81 | 16.34 / 42.01 | 19.19 / 51.92 |
| dilmash | 12.31 / 42.22 | 10.72 / 40.29 | 16.13 / 48.42 | 28.75 / 53.70 | 15.69 / 41.58 | 18.52 / 51.03 |
| dilmash-TIL | **15.02** / 45.43 | **12.00** / 42.07 | 17.59 / 49.90 | 32.07 / 56.45 | 17.53 / 43.52 | 19.83 / 52.58 |

Table 3: Evaluation of several models on sacreBLEU/chrF++ across various language pairs with Karakalpak on FLORES+ devtest set.

were trained on a bigger Karakalpak monocorpus.

3. **dilmash-TIL**: This variant was trained on the same dataset and tokenizer configuration as the **dilmash**, but supplemented with a strategically sampled subset from the TIL corpus. The sampling strategy was as follows:

   - 20% of parallel datasets containing Uzbek or Kazakh
   - 5% of all other datasets in the TIL corpus

To maintain balance with the Karakalpak dataset, we imposed an upper limit of 300,000 sentence pairs on the TIL corpus sample for the **dilmash-TIL**. This constraint ensured that the Karakalpak data was not overwhelmed by the additional multilingual data, while still allowing for potential improvements in cross-lingual transfer and overall model performance.

With a batch size of 1024 and using the AdaFactor (Shazeer and Stern, 2018) optimizer, we trained each model variant for 3 epochs. We employed a learning rate of 1e-4 with a linear warmup over the first 10% steps, followed by a constant learning rate schedule. Weight decay was set to 0.01 to help prevent overfitting.

To maximize computational efficiency and enable training on larger batch sizes, we utilized DeepSpeed ZeRO Stage 3 (Rasley et al., 2020) for model parallelism across 16 GPUs. This configuration allowed us to effectively distribute the model parameters and optimize memory usage, facilitating faster training times.

### 4.2 Evaluation Metrics

To evaluate the performance of our translation models, we employ two widely used metrics in machine translation:

- sacreBLEU (Post, 2018)

- chrF++ (Popović, 2017)

sacreBLEU, a standardized BLEU implementation, calculates the similarity between the machine-generated translation and one or more reference translations based on n-gram precision. It addresses inconsistencies in tokenization and BLEU computation across different implementations. chrF++, an extension of the character n-gram F-score, computes the F-score of character n-grams and word unigrams, which is particularly useful for morphologically rich languages, like Karakalpak or Uzbek.

## 5 Results and Discussion

Our evaluation on the FLORES+ Karakalpak devtest reveals several interesting insights into the performance of various translation models. The results, presented in Table 3, demonstrate the effectiveness of our proposed models, dilmash, dilmash-raw, and dilmash-TIL, in comparison to existing approaches.

Notably, the dilmash-raw model, which was trained on the original nllb-200-600M without additional tokens, outperforms the dilmash model with expanded vocabulary in most language pairs. This result suggests that the initialization of new token embeddings may have introduced some challenges. Our hypothesis is that the new token embeddings weren't initialized optimally, and before the model could learn good values for them, they may have affected other model parameters. The limited amount of Karakalpak data alone might not have been sufficient for the model to fully compensate for this initial distortion.

The dilmash-TIL model, which incorporates additional multilingual data from the TIL corpus, consistently outperforms both dilmash and dilmash-

raw across all language pairs. This improvement is particularly notable in the **\*-kaa** directions, with gains of up to 2.71 BLEU points (en-kaa) compared to dilmash. These results underscore two important points: first, the potential of using related Turkic language data to enhance translation quality for low-resource languages like Karakalpak; and second, that the additional training data and epochs may have allowed the model to better utilize the expanded vocabulary, overcoming the initial challenges faced by the dilmash model. To provide a more qualitative assessment of our models' performance, we have included translation examples in Appendix B.

While expanding the vocabulary can potentially improve model performance, careful consideration must be given to the initialization of new embeddings and the amount of training data available. The success of the dilmash-TIL model suggests that incorporating data from related languages and allowing for longer training periods can help overcome these challenges, ultimately leading to improved translation quality.

Interestingly, the Claude-3.5-sonnet model demonstrates superior performance in the **kaa-\*** directions, surpassing our models by a significant margin. This suggests that large language models may have a particular advantage in understanding content in low-resource languages, possibly due to their extensive pretraining on diverse multilingual data.

The performance of other baseline models provides additional context. Google Translate when treating Karakalpak as Kazakh or Uzbek, achieves respectable results but falls short of our models and Claude-3.5-sonnet. The NLLB-200-600M model, despite not being originally trained on Karakalpak, shows some ability to transfer knowledge when treating Karakalpak as Uzbek rather than Kazakh. This aligns with linguistic expectations, given the closer relationship between Karakalpak and Uzbek (both in linguistic similarity and writing scripts).

## 6    Conclusion

Our key contributions in this work include:

1. Creation of a FLORES+ devtest dataset for Karakalpak.

2. Development of parallel corpora for Uzbek-Karakalpak, Russian-Karakalpak, and English-Karakalpak, each containing 100,000 sentence pairs.

3. Open-sourcing of fine-tuned neural models for translation across Uzbek, Russian, English, and Karakalpak languages.

4. Open-sourcing of scripts for Latin-Cyrillic transliteration for Karakalpak.

Looking ahead, we plan to explore data augmentation techniques to further enhance our models' performance. One promising approach is to leverage the capabilities of Claude-3.5-sonnet for back-translation, potentially expanding our training data with high-quality synthetic examples.

Additionally, we aim to expand our dataset by mining more data from a wider range of books and sources. This will not only increase the volume of our training data but also improve its diversity, potentially leading to more robust and versatile translation models.

## 7    Limitations

While our study presents some advancements in Karakalpak machine translation, several limitations should be noted. First, the relatively small size of our dataset, despite being substantial for a low-resource language, may limit the model's ability to generalize across diverse domains and linguistic contexts. Second, the reliance on machine translation for creating the English-Karakalpak dataset introduces potential biases and errors that could affect model performance. Additionally, our evaluation is primarily based on automatic metrics, which may not fully capture the nuances of translation quality, particularly for a morphologically rich language like Karakalpak. Future work should address these limitations through expanded data collection, human evaluation, and more diverse testing scenarios.

## 8    Acknowledgements

# References

E. Berdimuratov and A. Dáwletov. 1979. *Házirgi qaraqalpaq tili*. Qaraqalpqastan, Uzbekistan SSR. Textbook for philology students at higher education institutions (in Cyrillic).

David Dale. 2022. The first neural machine translation system for the Erzya language. In *Proceedings of the first workshop on NLP applications to field linguistics*, pages 45–53, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

Ethnologue. 2024. Karakalpak. Accessed: 2024-07-25.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2020. Language-agnostic bert sentence embedding. *arXiv preprint arXiv:2007.01852*.

Mikel L Forcada, Mireia Ginestí-Rosell, Jacob Nord-falk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine translation*, 25:127–144.

Google. 2024. Google translate. https://translate.google.com. Accessed: July 25, 2024.

Alexander Jones, Isaac Caswell, Orhan Firat, and Ishank Saxena. 2023. Gatitos: Using a new multilingual lexicon for low-resource machine translation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 371–405.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2024. Madlad-400: A multilingual and document-level large audited dataset. *Advances in Neural Information Processing Systems*, 36.

Jamshidbek Mirzakhalov. 2021. Turkic interlingua: a case study of machine translation in low-resource languages. Master's thesis, University of South Florida.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Maja Popović. 2017. chrf++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506.

Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR.

## A   Prompt for translating from Russian to English using Claude-3.5-sonnet

```
You are a professional translator specializing in Russian to English translations.
Your task is to translate the given Russian text into English with the highest level
of accuracy, preserving the original meaning and context. Use proper grammar,
punctuation, and idiomatic expressions appropriate for English speakers.
Do not include any additional explanations or commentary; provide only the translated text.

Russian: {sent}
English:
```

## B    Translation examples from dilmash-TIL

| eng | kaa | dilmash-til (eng → kaa) |
|---|---|---|
| The winter can be deceptively chilly: temperatures rarely go below freezing, but the wind and humidity combine to make it feel colder than what the thermometer says. | Qıs aldamshı dárejede suwıq bolıwı múmkin: temperatura kemnen-kem jaǵdaylarda muzlawdan tómenge túsedi, bıraq samal hám ıǵallıq termometr aytqanınan kóre suwıqlaw seziliw ushın birlesedi. | Qıs júdá suwıq bolıwı múmkin: temperaturanıń buzılıw dárejesinen az ushırasıwı múmkin, biraq samal hám ıallıqtıń birigiwinen termometrdiń aytqanına qaraganda suwıq boladı. |
| **kaa** | **eng** | **dilmash-til (kaa → eng)** |
| Transport joq ekenligi qáweteri oyındı topar tárepdarlarısız jabıq esikler artında ótkeriwge májbúr bolıw múmkinshiligin asırdı. | Fears of lack of transportation raised the possibility that the game would be forced to play behind closed doors without the team's supporters. | Concern over the lack of transportation increased the possibility of being forced to hold the game behind closed doors without group supporters. |
| **rus** | **kaa** | **dilmash-til (rus → kaa)** |
| Также полезно носить кольцо (только не слишком дорогое на вид). | Júzik taǵıw da paydalı bolıp tabıladı (júdá qımbat kórinetuǵın júzik emes). | Sonday-aq, dóńgelek kiyiwde paydalı boladı (tek júdá qımbat kórinetuın emes). |
| **kaa** | **rus** | **dilmash-til (kaa → rus)** |
| Ullı piramida Firaun Xufudı ullılaw ushın jaratılǵan hám kóplegen kishi piramidalar, qábirler hám ibadatcxanalar Xufunıń hayalları hám shańaraq aǵzaların ullılaw ushın qurılǵan. | Великая Пирамида была создана в честь фараона Хеопса, а много меньших пирамид, гробниц и храмов были построены в честь жён Хеопаса и его родственников. | Большая пирамида была создана для восхваления фараона Хуфа, а многие небольшие пирамиды, могилы и храмы были построены для восхваления жен Хуфа и членов его семьи. |
| **uzb** | **kaa** | **dilmash-til (uzb → kaa)** |
| Ayrim atomlar turg'un bo'lmagan yadroga ega, bu esa ularning kichkina turtki bilan yoki turtkisiz parchalanishga moyilligini anglatadi. | Ayırım atomlar turaqsız yadrolarǵa iye, yaǵnıy olar azǵantay yamasa hesh qanday túrtki bolmaǵan halda bóleklenedi. | Ayırım atomlar turaqsız yadroga iye, bul bolsa olardıń kishi túrtki menen yaki túrtkisiz bóliniwine beyim ekenligi anlatadı. |
| **kaa** | **uzb** | **dilmash-til (kaa → uzb)** |
| Keshesi 150 den 200 ge shekem nusqalar tayarlandı, házirde "Dunlap broadsides" dep ataladı. | Tun davomida 150 dan 200 ga qadar nusxalar tayyorlandi, ular hozirda "Danlep yon zambaraklari" deb ataladi. | Kechasi 150 dan 200 gacha nusxalar tayyorlandi, hozirda "Dunlap broadsides" deb ataladi. |

Table 4: Some translation examples of dilmash-TIL model on FLORES+ sentences.

# FLORES+ translation and machine translation evaluation for the Erzya language

**Isai Gordeev**
École Polytechnique
Paris, France
isai.gordeev@polytechnique.edu

**Sergey Kuldin**
Astana, Kazakhstan
sergey@kuldin.com

**David Dale**
Paris, France

## Abstract

This paper introduces a translation of the FLO-RES+ dataset into the endangered Erzya language, with the goal of evaluating machine translation between this language and any of the other 200 languages already included into FLORES+. This translation was carried out as a part of the Open Language Data shared task at WMT24. We also present a benchmark of existing translation models bases on this dataset and a new translation model that achieves the state-of-the-art quality of translation into Erzya from Russian and English.

## 1 Introduction

The Erzya language is the language of Erzya, one of the indigenous peoples of Russia. Despite its official status in the Republic of Mordovia (along with Russian and Moksha), the use of the Erzya language is limited mainly to everyday topics, and its representation in the digital space remains low. On the one hand, this situation contributes to the reduced status of language in society and inhibits its development and intergenerational transfer, and is one of the factors that make the language endangered (UNESCO, 2010). On the other hand, this makes it more difficult to develop natural language processing (NLP) technologies for Erzya, such as machine translation (MT), the availability of which could help increase the prestige of the language. Development of MT technologies for Erzya, in turn, is hampered by the lack of a generally recognized dataset for evaluating the translation quality.

In this article, we aim to close this gap by publishing the Erzya version of the FLORES+ dataset. This dataset, created as part of the No Language Left Behind project (NLLB Team et al., 2022) and later transferred to the community-run Open Language Data Initiative[1], became the de facto standard for evaluating MT quality for low-resource languages.

FLORES+ consists of two thousand sentences sampled from English texts of Wikimedia projects and translated into more than 200 languages. The emergence of the Erzya version of FLORES+ will allow researchers to evaluate the quality of MT between any of these languages and Erzya. This translation has been submitted to the Open Language Data shared task at the WMT24 conference[2].

The quality of the new FLORES+ translation has been validated by independent annotators (for a sample from the dataset) and with a set of automatic metrics of accuracy and fluency which were themselves validated against human judgements.

In addition, we present a new neural model for translating between Erzya and other languages (primarily English and Russian), created by fine-tuning an NLLB-200 model. It achieves a BLEU score of 22% for translation from Erzya to Russian and 17% in the opposite direction, which implies the translation quality suitable for practical applications. Along with our new model, we evaluated the translation quality of several other models that also support the Erzya language. The model from Yankovskaya et al. (2023), also based on NLLB, achieved the highest scores for translation from Erzya into Russian, and a Claude model (Anthropic, 2024), into English, whereas our model achieved the highest scores for translation into Erzya.

In total, the contributions[3] of this article are:

1. We release and describe the first complete translation of the FLORES+ dataset into the Erzya language and validate its quality.
2. We evaluate the performance of available MT systems for Erzya on this dataset.
3. We present a model for translation between Erzya and several high-resourced languages, a state-of-the-art for translating into Erzya.

---

[1] https://oldi.org

[2] https://www2.statmt.org/wmt24/open-data.html
[3] Our code, data, and models will be made publicly available at https://github.com/slone-nlp/myv-nmt.

## 2 Related work

### 2.1 The FLORES+ dataset

FLORES+ is the next stage in the evolution of the FLORES-200 dataset (NLLB Team et al., 2022), which, in turn, grew out of FLORES-101 (Goyal et al., 2022). The dataset is based on 3001 English sentences taken from three sources: Wikinews (international news), Wikijunior (non-fiction literature for children), and Wikivoyage (travel tips). It is divided into three roughly equal parts (dev, devtest, and test), of which the first two (2005 sentences) are included in FLORES+. The sentences were translated from English into 203 other languages by professional translators; among these languages are Russian and three Finno-Ugric ones (related to Erzya): Finnish, Estonian, and Hungarian.

A small subset of FLORES-200 (the first 250 sentences from the devtest subset, news domain) were translated by Yankovskaya et al. (2023) into 9 low-resource Finno-Ugric languages, including Erzya, and used to evaluate the quality of the machine translation system from this article.

As far as we know, no other multiway parallel datasets including the Erzya language have been published (one exception is the dataset of the Tatoeba project[4], however, it currently contains less than 100 Erzya sentences).

### 2.2 Erzya datasets

The available corpora of the Erzya language (especially parallel ones) are not numerous. Rueter and Tyers (2018) presented the Erzya corpus with morphosyntactic markup, including the translation of several hundred sentences into English and Finnish. Arkhangelskiy (2019) has compiled a web corpus of the Erzya language (available for download); there is also a corpus of the literary Erzya language, avaliable only for search[5]. Another corpus of the Erzya language, also searchable, is described by Rueter (2024a).

Medium-scale parallel datasets of Erzya with other languages have been considered only in recent papers on neural machine translation for Erzya: Dale (2022) and Yankovskaya et al. (2023).

### 2.3 Machine translation for Erzya

The Apertium platform implements a machine translation system between Erzya and related Moksha and Finnish languages[6], but this work gives an impression of being incomplete. In Dale (2022), one of the first machine translation systems for the Erzya language was created, based on a parallel Russian-Erzya corpus consisting of dictionaries and automatically aligned sentences from books and web publications (a total of 77K parallel pairs of sentences, words and phrases, as well as 333K sentences in Erzya without translation). Yankovskaya et al. (2023) collected parallel and monolingual datasets for 20 low-resource Finno-Ugric languages and trained a neural model for their translation, but the parallel part of their dataset remained unpublished. In both works, the new languages were added to the pre-trained multilingual translation models (mBART-50 (Tang et al., 2020) and NLLB-200 (NLLB Team et al., 2022), respectively) by adding new tokens to the model vocabulary and fine-tuning it with parallel and back-translated data.

The third known paper on neural machine translation for Erzya, Alnajjar et al. (2023), used the rule-based Apertium system to generate a synthetic Erzya-Moksha corpus, and also fine-tuned an NLLB-200 model on this data.

Finally, Mordovian State University announced the development of an Erzya-Russian MT system[7], but no publication on this topic has appeared yet.

### 2.4 Multilingual MT systems

In addition to MT systems explicitly trained with Erzya parallel data, some models might have learned this language from monolingual texts or parallel texts unintentionally found in web corpora. One could also hope that multilingual models can achieve some understanding of the Erzya language based on the grammar and vocabulary of other Finno-Ugric languages, as well as on the vocabulary borrowed by Erzya from Russian. Therefore, we are considering several public and proprietary systems that do not always explicitly include Erzya, but may still be suitable for its translation.

One of such systems is the NLLB-200 family of models, which is still the leader among open models in terms of the translation quality and language coverage trade-off. Their training data did not include the Erzya language (except perhaps a small number of web texts mistakenly classified as Russian), so the NLLB ability to translate Erzya is limited by the knowledge transferred from other languages. The

---

[4]https://tatoeba.org/
[5]https://erzya.web-corpora.net/

[6]https://github.com/apertium/apertium-myv-mdf and https://github.com/apertium/apertium-myv-fin
[7]See the announcement on msru.ru.

MADLAD-400 models (Kudugunta et al., 2024) were also not trained with Erzya parallel data, but used the monolingual web corpora collected in this paper (including Erzya) for unsupervised training and back-translation.

Finally, we consider three proprietary systems: Google Translate, which has recently added 110 new languages (Caswell, 2024) powered by the PALM-2 (Google, 2023) large language model (LLM), and the Claude (Anthropic, 2024) and GPT-4o (OpenAI, 2023) LLMs. All the three systems were trained with multilingual web data and with a large amount of parallel data (including, possibly, parallel texts for the Erzya language). Unfortunately, the technical reports for these systems give only very brief descriptions of their language coverage.

## 3 About the Erzya language

The Erzya language (myv) is one of the largest Finno-Ugric languages in the world and belongs to the Mordvinic branch of the Finno-Ugric group of the Uralic language family. Linguists distinguish five main dialects: Central, Western (Insar), north-western (Alatyr), southeastern (Sura) and Shoksha (isolated in the northwestern regions). They differ mostly in their phonetic, and, to a lesser extent, morphological features. Our translation was based on the literary standard of the Erzya language (built primarily upon the Central dialect). Most modern Erzya is written using the Russian Cyrillic alphabet, although there are several Latin script proposals.[8]

Phonetically, the Erzya language contrasts palatalized and plain consonants and features vowel harmony. Grammatically, it is an agglutinative language with extensive systems of declension (including 12 noun cases, possesive suffixes and definitiveness marking) and conjugation (including 7 moods and marking the person and number of subject and object). The word order is SVO, and postpositions are widely used. Lexically, most words have Finno-Ugric origins,[9] with a significant number of Russian and Turkic loanwords.

In the XX century, the number of Erzya was approaching one million people (according to the 1970 census, 1,263 thousand people, along with Moksha). The dispersed settlement of the people

led to accelerated assimilation and the decrease in the number of native speakers. Thus, the Republic of Mordovia, where Erzya and Moksha are the titular nation, and their languages are co-official with Russian, hosts only about 30% of all Erzya, with the rest settled in Samara, Orenburg, Nizhny Novgorod, Penza, Saratov regions and other regions of Russia.

The status of the official language allows Erzya to function in the public space: there are newspapers and magazines, TV shows, theater, and popular music in this language. In addition, being a state language, Erzya is studied for 1-2 hours per week as an elective or mandatory lesson in many schools in the Republic of Mordovia, and in the settlements with a predominantly Erzya population, even as the first language.

Nevertheless, the Erzya language is poorly represented in the digital space. A few areas where it nevertheless functions are mentioned below:

- Several documentaries and feature films, music videos and video blogs;[10]
- The Wikipedia in Erzya with 7877 articles;[11]
- The Erzya interface of vk.com;[12]
- A few websites, thematic groups, and channels on social networks and messengers.

The problem of transferring the Erzya language to the younger generation is pressing: most children only understand, but almost do not speak their native language. We hope that translation of FLORES into Erzya will facilitate the development of machine translation for this language, which, in turn, could spur other technologies, such as speech synthesis, text and image generation models, contributing to the preservation, popularization and development of the language.

## 4 Translation of FLORES+

We translated FLORES+ from from Russian into the Erzya language. The translation was carried out by two native speaker volunteers who are also teachers of the language and writers (one with a doctoral degree in philology). The 250 translated sentences from Yankovskaya et al. (2023) were also included, but only after a thorough revision. All 2009 translations were revised by one of the native translators and a linguist with a profound expertise in the language. In addition, the translations were

---

[8]Table 2 features an example of a sentence in Erzya alongside with its Latin transliteration and translation.

[9]There are numerous, but not always easily recognisable cognates with other Uralic languages, such as Finnish and Estonian: for example, "keď / käsi / käsi / hand", "koto / kuusi / kuus / six", or "ĕräms / elää / elama / to live".

[10]E.g., Azor and Кода эри эрзянь морось movies.

[11]https://myv.wikipedia.org

[12]A news article about the interface.

| Neologisms | Examples in Erzya, Russian and English |
|---|---|
| Валдокаямо (luminosity, from валдо=light and каямо=output) | Весе **валдокаямось** ды чарамось сайневить тештень Россби числанть муеманзо кис, конась сюлмазь плазмань потоконть марто. |
| | Совокупность **светимости** и вращения используется для определения числа Россби звезды, связанного с потоком плазмы. |
| | The **luminosity** and rotation are used together to determine a star's Rossby number, which is related to plasma flow. |
| Тевконёв (document, from тев=business and конёв=paper) | Ломантненень, конат арсить теемс сымень полавтомань операция омбо масторсо, эряви парсте ванкшномс, улезт сынст марто мекев самонень эрявикс **тевконёвост**. |
| | Люди, планирующие пройти операцию по смене пола за границей, должны убедиться, что у них при себе есть действительные **документы** для обратного пути. |
| | Voyagers planning sex reassignment surgery abroad must ensure they're carrying valid **documents** for the return trip. |
| Инедавол (hurricane, from ине=great and давол=storm) | Раськень **инедаволонь** кунщкакуронть коряс те шкас Джерри а канды кодамояк зыян мода лангс. |
| | Согласно Национальному **ураганному** центру, на данный момент Джерри не представляет никакой угрозы на суше. |
| | The National **Hurricane** Center (NHC) says that at this point Jerry poses no threat to land. |
| Озавтозетне (inmates, "the imprisoned") | Чокшне ланга 10:00 ды 11:00 шканть юткcо пандонь шканть коряс **озавтозетне** тейсть кирвазтема вальмалост. |
| | Между 10:00 и 11:00 вечера по горному времени **заключенные** устроили пожар во дворе. |
| | Between 10:00-11:00 pm MDT, a fire was started by the **inmates** in the yard. |
| Кортницятне (negotiators, "talkers") | **Кортницятне** варчизь витемс тевенть, ансяк озавтозетнень вешемаст зть чарькодеве. |
| | **Переговорщики** попытались исправить ситуацию, но требования заключённых не ясны. |
| | **Negotiators** tried to rectify the situation, but the prisoners' demands are not clear. |
| Превмаксый (advisor, "intellect-giver") | 1960 иетнестэ Бзежинский ульнесь Джон Ф. Кеннединь ваксcо **превмаксыекс**, мейле Линдон Б.Джонсононь администрациясо. |
| | В 1960-х гг. Бжезинский занимал должность **советника** при Джоне Ф. Кеннеди, затем в администрации Линдона Б. Джонсона. |
| | Throughout 1960s, Brzezinski worked for John F. Kennedy as his **advisor** and then the Lyndon B. Johnson administration. |

Table 1: Examples of lexical neologisms created according to the word-formation models of the Erzya language (the top 3) and semantic neologisms created by assigning a new meaning to already known words (the bottom 3).

scored with automatic quality metrics (Section 5), which helped identify several omissions and typos.

The successful translation of texts on topics that are uncommon for Erzya allows us to assess the capabilities of this language positively. However, we should note the difficulties in translating special terminology in various domains (such as science, politics, and sports). In such cases, we used lexical and semantic neologisms to solve the problem of insufficient vocabulary (see the examples in Table 1). For some sentences, to avoid distorting their meaning, we had to preserve the Russian terminology, (e.g. Table 2). In some cases, the neologisms are translations of one part of a complex word, for example, "пельсфинал/pelśfinal" ("полуфинал" in Russian, "semi final" in English).

It would be difficult to directly evaluate how acceptable are these neologisms to the native speakers. But a human evaluation by two independent native speakers (Section 8) resulted in all sampled translations annotated as at least "acceptable", and the majority, as "good". This suggests that the new words are not perceived as serious problems to meaning preservation or fluency.

## 5 Automatic validation of data quality

To validate the quality of the newly translated Erzya FLORES dataset experimentally, we applied several automatic metrics of translation accuracy and fluency. To demonstrate the validity of the metrics themselves and to establish their baseline values, we needed human judgements of translation quality on some other dataset, and, fortunately, we had one.

**Data**. The baseline data consists of 1500 Erzya-Russian sentence pairs in the dev subset from Dale (2022), automatically aligned from various parallel documents. The sentence pairs were manually annotated by a proficient Erzya speaker for accuracy and fluency, with 0 standing for "unacceptable", 0.5 for "barely acceptable", and 1, for "good". The problems with meaning were mostly results of overly loose translations or incorrect alignment, whereas most of the fluency problems were caused by too literal translation from Russian (often by simply adding Erzya suffixes to the Russian words).

Our **simple metrics** are rel_sim (computed as the edit distance between the source and the target, divided by the maximum of the length of the two and then subtracted from 1) and len_ratio (the ratio

| Examples in Erzya (Cyrillic), Erzya (transliterated to Latin), Russian and English |
|---|
| Докладсонть ули малав эрьва аспектэнь пшти критика малав неень исполнительной властень Ираконь коряс политиканть коряс ды виев тердема сеске полавтомс улиця курсонть. |
| Dokladsonť uli malav ěŕva aspektěń pšti kritika malav neeń ispolniteĺnoj vlasteń Irakoń koräs politikanť koräs dy viev terdema seske polavtoms ulicä kursonť. |
| В докладе содержится резкая критика почти каждого аспекта нынешней политики исполнительной власти в отношении Ирака и настоятельный призыв к незамедлительной смене курса. |
| The Report is highly critical of almost every aspect of the present policy of the Executive towards Iraq and it urges an immediate change of direction. |

Table 2: An example of a translation with multiple loanwords from Russian (the loanwords are highlighted). Note that 4 out of these 6 words (aspektěń, kritika, politikanť, and kursonť), while being borrowed into Erzya from Russian, have ultimately Greek or Latin origins and are recognisable internationally.

of the character lengths of the source and target, the shortest of the two to the longest). We also include here the LID_rus metric: a probability, according to the GlotLID model (Kargaran et al., 2023), that the Erzya sentence is in fact Russian. We expect len_ratio to correlate with omissions; rel_sim and LID_rus are expected to correlate with fluency issues. The rest of the metrics, described below, target translation accuracy.

**Dictionary-based metrics** are computed by lemmatizing the words in a sentence pair and computing the proportion of words on the one side that has a counterpart on the other side that can be matched using a dictionary. WR_rus computes the proportion of Russian words whose translations can be found in the Erzya sentence, and WR_myv shows the opposite (they don't necessarily match because the Russian and Erzya sentences may have a different number of words).

**Model-based metrics** include two cosine similarities of sentence embeddings: LaBSE uses a LaBSE model (Feng et al., 2022) distilled for Erzya[13], and enc_sim uses the mean token embeddings from the encoder of the NLLB-based MT model described in Section 6. The latter model is also used for computing Ppl: the mean cross-entropy loss (perplexity) of the model for translating the Erzya sentence to Russian and in reverse. The Att metric is based on the encoder-decoder attention maps for this model: we average the cross-attentions to each encoder token over the layers and the heads, add up across the decoder tokens, and average across all the encoder tokens, except the first one (language code) and the last one (end of sentence), as they are expected to serve as "attention sinks". This metric is also averaged across two translation directions.

**Correlations**. To evaluate the quality of the metrics, we report their Spearman correlation with the quality annotation labels in the two last rows of Table 3. As we expected, rel_sim and LID_rus are predictive of fluency problems, and all dictionary- and model-based metrics correlate with accuracy.

**Data comparison**. The top two rows of Table 3 report the average values of our automatic metrics on our dev and devtest splits of FLORES. The next four rows report their values on the baseline data, depending on the presence or absence of fluency and meaning preservation problems. According to most metrics, the FLORES translations are similar or even better than the "good" (problem-free) subset of the baseline data. The only exception is WR_rus which for FLORES is slightly below the baseline; this might be explained by the difficulty of the FLORES domains, where many words are not yet covered by the Erzya dictionaries.

## 6 A new MT model for Erzya

Our preliminary exploration suggested that few existing models are capable of producing reliable Erzya sentences, so we tried training our own translation model, focused on translation into Erzya.

### 6.1 Datasets

**"Natural" data**. To train our model, we used the same monolingual Erzya and parallel Erzya-Russian datasets as Dale (2022). In addition, we collected and aligned some parallel news articles[14], and included several new translated books[15] and pairs of words and phrases from the Russian-Erzya dictionaries at Emerald (Rueter, 2024a,b) and Panlex (Kamholz et al., 2014). The books and articles have been aligned at the sentence level using the

---

[14]http://e-mordovia.ru
[15]A physics textbook, two books of Alexander Doronin, and a translation of The Little Prince. All copyright holders gave consent to use the texts as training data.

| Dataset | rel_sim | len_ratio | LID_rus | WR_rus | WR_myv | LaBSE | enc_sim | Ppl | Att |
|---|---|---|---|---|---|---|---|---|---|
| FLORES dev | 0.28 | 0.90 | 0.07 | 0.55 | 0.64 | 0.89 | 0.82 | 1.49 | 0.32 |
| FLORES devtest | 0.28 | 0.90 | 0.06 | 0.55 | 0.64 | 0.89 | 0.83 | 1.51 | 0.32 |
| BL (good) | 0.27 | 0.83 | 0.06 | 0.59 | 0.62 | 0.86 | 0.84 | 2.51 | 0.28 |
| BL (fluency problems) | 0.38 | 0.87 | 0.14 | 0.70 | 0.74 | 0.94 | 0.90 | 1.73 | 0.28 |
| BL (meaning problems) | 0.23 | 0.79 | 0.07 | 0.43 | 0.44 | 0.66 | 0.69 | 3.56 | 0.25 |
| BL (both problems) | 0.23 | 0.90 | 0.13 | 0.35 | 0.37 | 0.65 | 0.69 | 3.64 | 0.23 |
| BL, corr. with meaning | 0.23 | 0.01 | -0.01 | 0.38 | 0.38 | 0.48 | 0.43 | -0.39 | 0.30 |
| BL, corr. with fluency | -0.30 | -0.17 | -0.35 | -0.17 | -0.19 | -0.33 | -0.26 | 0.26 | 0.05 |

Table 3: The average values of the automatic metrics on our FLORES translation (top 2 rows); their average values on the baseline data grouped by quality (next 4 rows), and their Spearman correlations with human quality labels on the baseline data (the last 2 rows).

algorithm from Dale (2022). We filtered out the pairs for which the Levenshtein distance was less than 20% of the text length, since their inclusion could lead to the model learning to copy the source words too often instead of translating them. We also dropped the pairs with more than 55% difference in length, as they were likely incorrect as translations. The volume of the cleaned Erzya-Russian dataset was 174460 pairs of sentences, phrases or words.[16]

**Back-translation**. To take advantage of the monolingual Erzya texts, we translated 200K Erzya sentences with the previous version of our model (trained only on "natural" data): 50% into Russian and 50% into 13 other languages previously represented in NLLB-200[17]. After filtering by string edit distances and length ratios, 176462 texts remained. To enhance the model's ability to translate from Erzya into languages other than Russian, we also translated 30K Russian sentences from the Erzya-Russian parallel dataset into the same 13 other languages using an NLLB-200-600M model.

During the training, the pairs of texts from one natural and two synthetic sources were randomly selected in the following proportion: 70% from the natural data (in both directions), 25% from the data translated from Erzya (only in the opposite direction, into Erzya), and 5% from the data translated from Russian to other languages (in both directions). We normalized 100% of the texts on the target side and 80% of the texts on the source side[18] using the

algorithm from NLLB Team et al. (2022)[19].

## 6.2 Model training

When training the model, we followed an approach similar to Tars et al. (2022) and Dale (2022). As a basic model, we took NLLB-200 with 600 million parameters, enriched its dictionary with new tokens for the Erzya language, further trained embeddings for these tokens, and then further trained the entire model on parallel Russian-Erzya data, as well as on data obtained by reverse translation.

**Vocabulary update**. To better represent Erzya words in the model, we trained a new Sentencepiece tokenizer (Kudo and Richardson, 2018) on the Erzya side of the training dataset. Most of its tokens (6208 out of 8192) were missing from the NLLB vocabulary, so we added them there. The corresponding embeddings of each new token were initialized by the mean of the embeddings of the "old" tokens into which the new token could be decomposed. We also added a new language code to the tokenizer: myv_Cyrl.

**Fine-tuning**. Since the embeddings of the new tokens were initialized with a naive method, fine-tuning of the whole model with these parameters could introduce undesirable disturbances into other parameters. To avoid this, for the first 45K training steps, we updated only the embedding matrix, "freezing" the rest of the model parameters, and used an additional loss function (with a weight of 100): the mean squared error between the original and current embeddings of the "old" tokens. We used a single GPU, a batch of size 6, 4 gradient accumulation steps, and an Adafactor optimizer (Shazeer and Stern, 2018) with the learning rate that linearly warmed up from zero to $10^{-4}$ during the first 3000 steps and then stayed constant. After updating the token embeddings this way, we

---

[16]The collected parallel dataset (at least, its part that is free of copyright restrictions) will be made publicly available in our repository.

[17]Arabic, English, Estonian, Finnish, French, German, Kazakh, Mandarin, Mongolian, Spanish, Turkish, Ukrainian, and Uzbek. This choice was motivated by the international importance of the languages, by their connections to the post-Soviet region where most Erzya live, and by an attempt to represent diverse language families.

[18]We kept 20% of the source texts unnormalized to better prepare the model for potential downstream use cases when the input is not normalized.

[19]The normalization code was adopted from the Sentence-SplitClean class in the Stopes repository.

continued training the whole model (with all the parameters unfrozen), for 220K more steps.

# 7 Comparing MT systems for Erzya

We used our FLORES+ translation (the dev subset) to evaluate the current quality of MT between Erzya on the one side and Russian and English on the other. We evaluated the translation quality with BLEU (Papineni et al., 2002).

## 7.1 The systems

Here we describe each of the systems that we tried including in our benchmark. For all models except the LLMs, we use beam search with the beam size of 5, and keep all other inference parameters at their default values, unless otherwise specified.

**Ours**. We used the model described in Section 6. For this model (and for NLLB), we normalized the input texts with the NLLB normalization algorithm.

**SLONE** (Dale, 2022). We use the myv-mul and mul-myv models from Dale (2022) to translate from and to Erzya, respectively. We use beam size of 5 and repetition penalty of 5, like in the original work (but we do not use reranking with the LID model).

**SMUGRI** (Yankovskaya et al., 2023). We used their latest model, also based on the NLLB-200 (with 1.3B parameters)[20], for translation in all 4 directions. We used the fairseq-interactive interface with beam size of 5, penalty of 100 for the unk token, and a maximum of 4 consecutive tokens repeating, and ran the model in the fp16 precision.

**NLLB**. We tried to use the NLLB-200-600M model to translate from Erzya into English and Russian. Since NLLB requires specifying the source language, and Erzya is not included in it, we indicated Estonian as the source language: it is genetically related to Erzya, and also, like Erzya, has experienced some lexical influence of Russian.

We used the MADLAD-400-3B-MT model (Kudugunta et al., 2024) in two configurations: as is (**MADLAD**), and fine-tuned with our Erzya-Russian dataset (without other languages) for 60K steps using a new Erzya token (**MADLAD-ft**). For both of them we used the HuggingFace package with beam size of 3.

---

[20]The paper mentions M2M-100 as a base model. But the model that we used, https://huggingface.co/tartuNLP/smugri3_14-finno-ugric-nmt, is, apparently, a newer version, and it has the language code formats and the size of NLLB-200-1.3B. According to the model card, it is currently powering https://translate.ut.ee.

| System | ru-myv | en-myv | myv-ru | myv-en |
|---|---|---|---|---|
| NLLB | - | - | 3.23 | 1.05 |
| SLONE | 8.11 | 4.99 | 15.12 | 11.70 |
| SMUGRI | 11.46 | 6.58 | 28.44 | **21.12** |
| MADLAD | 1.10 | 1.06 | 18.48 | 13.87 |
| MADLAD-ft | 15.50 | - | 25.50 | - |
| Claude | 14.13 | 7.09 | **34.68** | 20.18 |
| GPT-4o | 3.49 | 1.24 | 11.07 | 8.73 |
| Ours | **17.09** | **7.35** | 22.06 | 16.42 |

Table 4: BLEU scores for the evaluated MT systems on the dev subset of FLORES+.

**Google Translate**. We have tried several configurations of the Google Translate API: the general/nmt and general/translation-llm models in the Advanced v3 interface. Both models proved unable to translate from Erzya into Russian or English either when specifying related languages (Estonian, Finnish, Udmurt, Mari) as the source language, or when automatically detecting the language (most often it was defined as Udmurt, Mari, or Komi). In most cases, the models simply transliterated the Erzya text into the Latin alphabet, without trying to translate most of the words. Based on these results, we excluded Google Translate from the benchmark.

**Claude** (Anthropic, 2024) and **GPT-4**. We used the API of the Claude 3.5 Sonnet and GPT-4o models, respectively, to obtain the translations. For some texts, we obtained the necessary translations only after several iterations of inference without changing the prompt, because in the case of GPT, the prompt did not pass the jailbreak protection, and in the case of Claude, the prompt could produce empty translations, which were corrected upon retranslation. An example of our prompt is given in Appendix A.

## 7.2 Evaluation results

The results of MT evaluation with BLEU are in Table 4. As we had hoped, our model achieved the best quality of translation into Erzya from Russian and English. The Claude model won the first prize for translating from Erzya to Russian, and the model from Yankovskaya et al. (2023), from Erzya to English. In addition, the MADLAD model demonstrated promising results, with rather good understanding of Erzya even before fine-tuning.

# 8 Human validation of human and machine translation

While the BLEU scores reported above may help ranking translation systems, they cannot tell how good the translation is from the human point of

view. To shed some light on this, we engaged two new native Erzya speakers (not involved into the FLORES translation) to evaluate the translations of 30 randomly chosen FLORES sentences from Russian into Erzya. We showed them the human translations and machine translations by SLONE, Claude, and our new system, without displaying the system names and in a random order, to reduce the potential bias. The annotators were asked to rate each translation on a 1-5 point scale from Dale (2022), with the label 3 for "acceptable", 4 for "good", and 5 for "great" translations. Their full guidelines are given in Appendix B.

| System | Annotator 1 | Annotator 2 | $\kappa$ | $\rho$ |
|--------|-------------|-------------|----------|--------|
| Human  | 4.37(0.14)  | 4.33(0.12)  | -0.03    | -0.02  |
| Claude | 4.17(0.14)  | 4.1(0.13)   | 0.24     | 0.41   |
| Our MT | 3.9(0.19)   | 3.87(0.17)  | 0.29     | 0.54   |
| SMUGRI | 3.47(0.21)  | 3.57(0.21)  | 0.27     | 0.53   |

Table 5: Mean assigned scores (with standard errors in brackets), Cohen's kappa $\kappa$ and Spearman correlation $\rho$ of the two annotators' labels.

Table 5 reports the mean scores assigned by the annotators to each system, and their agreement scores. The inter-annotator agreement is fair for the MT systems, but there is none for human translations, indicating that more fine-grained annotation schemes might be needed in the future. Nevertheless, both annotators assign the highest (and similar) average scores to the human translations, reaffirming their quality.

Our MT system takes the third place in the human ranking, after the human translations and Claude. The reason for this low position is 3 "stupid" translation errors (out of the 30): two undertranslations and one cyclical hallucination. We hope that in the future, simple modifications of the decoding algorithm would help avoid such errors.

## 9 Conclusion

Although the *endangered* and *low-resourced* statuses for a language are by no means equivalent (Hämäläinen, 2021), they reinforce each other in a vicious circle. Lack of resources for a language lowers its prestige, which reduces the number of active speakers, which, in turn, disincentivises creation and maintenance of the language resources. As an example, the endangered Erzya language, with its few hundred thousand speakers and a modest Wikipedia community, did not make it to the FLORES-200 dataset and the NLLB-200 models —

but if it did, it could have a positive impact at least on the Erzya Wikipedia.

Such projects as the Open Language Data Initiative shared task give one more chance to such languages. By releasing a FLORES+ translation into Erzya and using it to benchmark the few existing MT models that support it, we hope to help rolling the vicious circle in the opposite, virtuous direction. An Erzya version of FLORES might open a way to include this language into other NLP evaluation datasets, such as FLORES-based FLEURS (Conneau et al., 2023) for speech translation and Belebele (Bandarkar et al., 2024) for machine reading comprehension. And the presence of the language in such datasets, we hope, might motivate the researchers to include it in foundation models, which, in turn, might influence the development of practical applications, supporting the speakers of the language and increasing its status and chances of survival.

More directly, we are hoping that the emergency of the Erzya version of FLORES will facilitate improvements in machine translation research and applications for this language.

Some possible areas of future research based on the Erzya translation of FLORES+ include:

- Translating new MT training datasets, such as the Seed (Maillard et al., 2023)), into Erzya.
- Creating an automatic reference-free metric of translation quality for Erzya that would highly correlate with human judgments.
- Setting up a system of active learning that would help collect human translations for the sentences at which MT most likely fails.
- Creating an MT system suitable for translating content (such as books) into the Erzya language, minimizing the necessary revision and post-correction by human translators.

## 10 Acknowledgements

## References

Khalid Alnajjar, Mika Hämäläinen, and Jack Rueter. 2023. Bootstrapping Moksha-Erzya neural machine translation from rule-based apertium. In *Proceedings of the Joint 3rd International Conference on Natural Language Processing for Digital Humanities and 8th International Workshop on Computational Linguistics for Uralic Languages*, pages 213–218, Tokyo, Japan. Association for Computational Linguistics.

Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf. Accessed: 2024-08-19.

Timofey Arkhangelskiy. 2019. Corpora of social media in minority Uralic languages. In *Proceedings of the Fifth International Workshop on Computational Linguistics for Uralic Languages*, pages 125–140, Tartu, Estonia. Association for Computational Linguistics.

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. 2024. The belebele benchmark: a parallel reading comprehension dataset in 122 language variants. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 749–775, Bangkok, Thailand. Association for Computational Linguistics.

Isaac Caswell. 2024. 110 new languages are coming to google translate. https://blog.google/products/translate/google-translate-new-languages-2024/. Accessed: 2024-08-19.

Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2023. Fleurs: Few-shot learning evaluation of universal representations of speech. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 798–805. IEEE.

David Dale. 2022. The first neural machine translation system for the Erzya language. In *Proceedings of the first workshop on NLP applications to field linguistics*, pages 45–53, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Google. 2023. Palm 2 technical report. https://ai.google/static/documents/palm2techreport.pdf. Accessed: 2024-08-19.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Mika Hämäläinen. 2021. Endangered languages are not low-resourced! *arXiv preprint arXiv:2103.09567*.

David Kamholz, Jonathan Pool, and Susan Colowick. 2014. PanLex: Building a resource for panlingual lexical translation. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3145–3150, Reykjavik, Iceland. European Language Resources Association (ELRA).

Amir Hossein Kargaran, Ayyoob Imani, François Yvon, and Hinrich Schuetze. 2023. GlotLID: Language identification for low-resource languages. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6155–6218, Singapore. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2024. Madlad-400: A multilingual and document-level large audited dataset. *Advances in Neural Information Processing Systems*, 36.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzman. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

NLLB Team, Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2207.04672*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia,

Pennsylvania, USA. Association for Computational Linguistics.

Jack Rueter. 2024a. On searchable mordvin corpora at the language bank of finland, emerald. *Journal of Data Mining & Digital Humanities*, (V. The contribution of corpora).

Jack Rueter and Francis Tyers. 2018. Towards an open-source universal-dependency treebank for Erzya. In *Proceedings of the Fourth International Workshop on Computational Linguistics of Uralic Languages*, pages 106–118, Helsinki, Finland. Association for Computational Linguistics.

Jack Michael Rueter. 2024b. erzya-bidix.

Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR.

Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. Multilingual translation with extensible multilingual pretraining and finetuning. *arXiv preprint arXiv:2008.00401*.

Maali Tars, Taido Purason, and Andre Tättar. 2022. Teaching unseen low-resource languages to large translation models. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 375–380, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

UNESCO. 2010. Unesco atlas of the world's languages in danger (pdf).

Lisa Yankovskaya, Maali Tars, Andre Tättar, and Mark Fishel. 2023. Machine translation for low-resource Finno-Ugric languages. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 762–771, Tórshavn, Faroe Islands. University of Tartu Library.

## A  Example prompts

We used the following prompt format for translating with Claude: *"You are an AI assistant responsible for translating phrases between Russian and Erzya-Mordvin. Your task is to translate a given sentence from Erzya-Mordvin into Russian. Ensure that the translation remains specific to Erzya-Mordvin, avoiding confusion with Komi, Moksha, or other Finno-Ugric languages to prevent false cognates. Verify that each word in the original Erzya-Mordvin sentence has a corresponding translated word in Russian, maintaining the accuracy and completeness of the content. The final translated sentence should retain a similar word count without omitting any parts of the original*

*text. Output only the final translated result in Russian.*
*{The source sentence}. Output only the translated result."*. With GPT, we used a similar format, with the first paragraph fed as the system prompt, and the source sentence as the user prompt.

## B  Annotation guidelines

For human evaluation of human and machine translation in Section 8, we provided the annotators with a short guideline text in Russian. Below is its translation into English.

*We ask you to rate the translations from Russian to Erzya on a scale of 1-5 points:*

*5 points - perfect translation (the meaning and style are fully preserved, the grammar and word choice are correct, the text looks natural);*

*4 points - good translation (the meaning is fully or almost completely preserved, the style and choice of words are acceptable for the target language);*

*3 points - acceptable translation (the core meaning is preserved; mistakes in word choice and grammar do not interfere with understanding; most of the text is fluent and in the target language);*

*2 points - poor translation (the text is mostly understandable and mostly in the target language, but there are serious errors in meaning preservation, grammar or word choice);*

*1 point - unsuitable translation (most of the text is in the wrong language, or nonsense, or has little in common with the original text).*

*If at least one word is incorrectly translated, the resulting score should not be 5; the choice between 1 and 4 is at your discretion.*

*If a word is an overly literary term or a neologism, but its meaning is clear, it does not lower the score. However, if the usage of an unusual word is unclear or it changes the original meaning, lower the score.*

# Spanish Corpus and Provenance with Computer-Aided Translation for the WMT24 OLDI Shared Task

**Jose Cols**
Department of Linguistics, University of Washington
jcols@uw.edu

## Abstract

This paper presents the SEED-CAT submission to the WMT24 Open Language Data Initiative shared task. We detail our data collection method, which involves a computer-aided translation tool developed explicitly for translating SEED corpora. We release a professionally translated Spanish corpus and a provenance dataset documenting the translation process. The quality of the data was validated on the FLORES+ benchmark with English-Spanish neural machine translation models, achieving an average chrF++ score of 34.9.

## 1 Introduction

In recent years, the NLP community has made significant strides in reducing the data gap for hundreds of languages (Tiedemann, 2012; Bañón et al., 2020; Federmann et al., 2022; NLLB Team et al., 2022). Nonetheless, finding parallel corpora for machine translation and other NLP applications remains challenging for many language pairs (Haddow et al., 2022; Ranathunga et al., 2023). The WMT24 Open Language Data Initiative shared task aims to continue expanding language coverage with contributions from communities of native speakers.

This work describes our data collection method to expand the SEED dataset (NLLB Team et al., 2022; Maillard et al., 2023) with the Spanish language. Specifically, we focus on Latin American Spanish varieties to match the existing coverage of this language in the FLORES+ benchmark (NLLB Team et al., 2022). While the Spanish language benefits from the availability of multiple parallel corpora datasets (Aulamo et al., 2020), the majority of this corpora features other well-resourced languages such as English and French, and translation directions of regional significance to other languages like Asturian and Quechua remains a challenge (Oliver et al., 2023; Ahmed et al., 2023).

The multilingual alignment of the SEED dataset (Doumbouya et al., 2023) allows for the addition of a single corpus to enable dozens of translation directions into low-resource languages. Including Spanish represents an essential step toward incorporating other low-resource languages where finding English translators is challenging, as was the case for Ligurian, where half the data was translated from Italian (NLLB Team et al., 2022).

Considering the impact that high-quality parallel corpora can have on machine translation performance (Maillard et al., 2023), this work aims to facilitate extending the SEED dataset while supporting quality improvements in Spanish machine translation. Our main contributions are:

1. The expansion of the SEED dataset with professional translations of Latin American Spanish, created by native speakers, along with a neural machine translation baseline.

2. The open-source release of SEED-CAT, a web computer-aided translation tool explicitly designed to assist human translators in the translation of SEED files.

3. The automated generation and public release of a provenance dataset documenting the creation of each Spanish translation.

## 2 Background

**Language overview**  According to a 2022 report from the Cervantes Institute,[1] there are more than 496 million native Spanish speakers in the world. Speakers are mainly concentrated in the Americas and the Iberian Peninsula, with Mexico having the largest population.

Spanish is an Ibero-Romance language that developed from Latin on the Iberian Peninsula. Thanks to its global expansion, this language has evolved into several dialectal variations. An example of this variation is the absence of the informal

---

[1] https://cvc.cervantes.es

second-person plural 'vosotros' in Latin American Spanish, where most varieties use the pronominal form 'ustedes' to address speakers in both formal and informal contexts (Hualde et al., 2012).

Although there are social, phonological, and lexical variations, Spanish retains a fundamental cohesiveness (Hualde et al., 2012). The Royal Spanish Academy and the Association of Academies of the Spanish Language collaborate to publish a unified set of orthography, dictionaries, and other language resources. The Spanish writing system is based on the Latin script, with the addition of the character ⟨ñ⟩ forming an alphabet of 27 letters (Hernández Gómez, 2015). This script is represented in our collected data.

**Seed dataset**   The SEED dataset (NLLB Team et al., 2022; Maillard et al., 2023), currently managed by the Open Language Data Initiative (OLDI),[2] contains 6,193 parallel sentences in English along with professional translations into 40 low-resource languages. The English sentences were originally sampled from thousands of Wikipedia articles across various categories such as arts, history, mathematics, people, and technology, offering diverse content from notable topics (NLLB Team et al., 2022). In this work, we use the English corpus (eng_Latn) from the SEED dataset as the source text for the Spanish translations.

**FLORES+ benchmark**   FLORES+ is an evaluation benchmark for machine translation with support for 212 languages based on the initial FLORES-101 dataset (Goyal et al., 2022) and its recent expansions (NLLB Team et al., 2022; Doumbouya et al., 2023). The collection of this data involved a rigorous and iterative quality assurance process with professional translators, pre-defined standards, post-editing, and automatic quality assessments. We rely on this benchmark to assess the quality of the Spanish translations.

**Computer-aided translation**   CAT, or computer-aided translation, refers to software tools, such as word processing, translation memory (TM), and terminology management, used by human translators to assist the translation process (Bowker and Fisher, 2010). Studies have shown that these tools can enhance the productivity and translation quality of human translators (Federico et al., 2012; Koehn, 2009). While machine translation differs from other CAT tools, as it is the primary

driver of the translation (Bowker and Fisher, 2010), modern CAT suites often include machine translation as a key feature. According to a user survey study involving 736 translators (Zaretskaya et al., 2017), machine translation ranked as the third most commonly used functionality, following translation memory and terminology management. CAT users from that study and other usability surveys (Alotaibi, 2020; Vargas-Sierra, 2019) also reported dissatisfaction with the ease of use and learnability of these systems, highlighting the importance of user-friendly interfaces for computer-assisted translation.

## 3   Data Collection

**Seed-CAT**   Various commercial CAT solutions exist, with SDL Trados, Memsource, and Wordfast being recognized as popular options by different research (Alotaibi, 2020; Picton et al., 2017). Apart from requiring purchasing a license, these systems use custom file formats that may not be compatible with other tools, leading to interoperability issues in translation projects. Using general-purpose software can also result in unaligned parallel sentences due to translators re-ordering the files, a problem highlighted by Doumbouya et al. (2023) in their review of the original NLLB-SEED dataset (NLLB Team et al., 2022). Furthermore, commercial CAT systems often integrate machine translation models, such as Google Translate and DeepL, that restrict the use of their outputs for training other models.

Recognizing these challenges, we release SEED-CAT,[3] an open-source web application specifically designed to assist human translators in translating SEED dataset files. This application was at the center of our data collection efforts and was designed with the three core principles.

- The user interface and features are optimized for usability, device compatibility, and seamless integration with the SEED dataset. The list of languages and corpora is fetched at runtime from the dataset's repository, and metadata is displayed alongside each sentence (Figure 1).

- The system architecture facilitates application deployment, as it does not require configuring databases or user accounts. Data persistency is achieved via IndexedDB,[4] a transactional

---

[2] https://oldi.org/

[3] https://github.com/josecols/seed-cat
[4] https://www.w3.org/TR/IndexedDB/

625

database for object storage in web browsers.

- The application data model adheres to the W3C PROV-DM (Missier and Moreau, 2013) recommendation for data provenance, adding an additional layer of transparency to the translation creation process.

SEED-CAT integrates a focused set of features, such as machine translation and terminology consultation. Machine translation is supported for local[5] and remote inference, with the latter being recommended for broader device compatibility. The machine translation feature relies on the `facebook/nllb-200-distilled-600M` model (NLLB Team et al., 2022) to generate outputs. Likewise, terminology consultation in English is enabled by WordNet (Miller, 1995).

Users can also compare translations using text differencing (Myers, 2023), with word-based comparison for Latin-based languages and character-based comparison for other scripts. Additionally, part-of-speech color highlighting for English words can be toggled based on user preference. This feature relies on the Brill tagger (Brill, 1992) and a Treebank tokenizer (Marcus et al., 1993), both implemented using the `natural` library.[6]

**Sourcing translators** We recruited a team of ten freelance translators who were individually sourced through Fiverr, an online marketplace for digital services. We relied on the platform's reputation system and freelancer profiles to identify potential candidates. The final translators were selected based on specific criteria: native Latin American Spanish speakers, a minimum of two years of translation experience on the platform, at least 500 completed projects, and a brief English conversation assessment. The median translator had nine years of experience, 1,900 completed projects, and 779 reviews. In addition, we sourced an independent freelance translator with a degree in Applied Languages who underwent a similar vetting process. Additional background information on all translators is reported in Table 1.

**Compensation** Each translator determined their compensation separately based on the number of English words in their assigned task. The tasks were divided into two stages: *translation* and *review*, which had different compensation rates. The

| Category | Detail | % |
|---|---|---|
| Education | Master's degree | 9.1 |
| | Bachelor's degree | 54.5 |
| | Course or certificate | 27.3 |
| | No formal training | 9.1 |
| Country | Argentina | 9.1 |
| | Chile | 9.1 |
| | Colombia | 9.1 |
| | Mexico | 18.2 |
| | Panama | 9.1 |
| | Venezuela | 45.5 |

Table 1: Percentage distribution of participants by educational background in translation and country of origin.

median compensation per translated word was 0.017 US dollars, with an average of 0.022, and the median compensation per reviewed word was 0.012 US dollars. Translators were also given a user guide to the SEED-CAT application, along with data samples and translation guidelines, to help them assess the complexity of the task when determining their rates.

**Translation workflow** The translation process was divided into two stages: first, translating all sentences from English to Spanish, and second, reviewing every sentence to ensure accuracy and quality. Both phases were carried out by the team of translators using the SEED-CAT application.

A team of ten translators completed the initial translation phase in 16 days. The translators worked on a contiguous set of segments for better contextual reference, with an average task size of 593 sentences. Each translator received a unique URL to access SEED-CAT, which automatically configured their browser with the target language file (`spa_Latn`), sentence range, and user identifier. Translators did not need to handle administrative tasks such as creating user accounts or managing assignments. When they opened the application URL, they were prompted to review and acknowledge the translation guidelines, and then they were directed to their first assigned segment for translation.

The review phase was carried out by three translators. Each reviewed segments that were initially translated by others, finishing the task in five days and proofreading an average of 2,064 sentences. The translators received a specific URL to open SEED-CAT in `review` mode. In this mode, the application automatically loads and deserializes the

---

[5] https://github.com/xenova/transformers.js
[6] https://github.com/NaturalNode/natural/

Figure 1: SEED-CAT's user interface with two resizable panels to display the original sentence and the translation editor. Users can select different languages, track progress, review guidelines, or access other actions, such as importing/exporting provenance graphs using the top navigation bar. Translators can also open the source document and generate machine translations.

provenance information collected up to that point, enabling the consolidation of the translation and review history of a sentence into a single PROV-JSON file (Huynh et al., 2013). A total of 686 translations were copy-edited, with most corrections involving mistranslations, syntactic and lexical refinements, and grammatical issues such as verb agreement. Additionally, the decimal and thousand separators were standardized following established Spanish orthographic norms (Real Academia Española, 2010).

**Dataset sample**    The final dataset contains 6,193 Spanish sentences (152,664 words) professionally translated by eleven native speakers from six countries in Latin America. Table 2 provides a brief excerpt of the parallel sentences.

**Provenance dataset**    According to the World Wide Web Consortium (W3C) (Groth and Moreau, 2013), "provenance is information about entities, activities, and people involved in producing a piece of data or thing, which can be used to form assessments about its quality, reliability or trustworthiness."

During the translation process, the SEED-CAT application automatically recorded provenance information on how activities such as

EditTranslation, MachineTranslate, and QueryWordNet were used to generate, invalidate, and revise translations. This information can be serialized into JSON files (Huynh et al., 2013), enabling data sharing with its complete history. Users can also import these files and modify the data entities while maintaining the provenance's integrity.

In the PROV data model (Missier and Moreau, 2013), entities, activities, and agents are linked through relations. These links can be used to create a directed graph to visualize dependencies and data interactions. Appendix E provides examples of these graphs from the Spanish dataset. This dataset containing 6,193 PROV-JSON files is released as part of our SEED contribution.

**System usability scale**    The system usability scale (SUS) (Brooke, 1996) is a standardized 10-item questionnaire for assessing perceived usability. Users rate each statement of the survey on a scale from 1 to 5, enabling the calculation of the SUS score, which ranges from 0 to 100. Substantial research has found that this score is a reliable metric of perceived system usability (Lewis, 2018). We administered the SUS questionnaire to the eleven translators involved in the project to evaluate the SEED-CAT application, obtaining an SUS score of

| # | English | Spanish |
|---|---------|---------|
| 663 | For Gibbon, "The decline of Rome was the natural and inevitable effect of immoderate greatness. | Para Gibbon, "El declive de Roma fue el efecto natural e inevitable de la grandeza excesiva. |
| 2079 | By 1843 Richard Hoe developed the rotary press, and in 1844 Samuel Morse sent the first public telegraph message. | En 1843 Richard Hoe inventó la prensa rotativa, y en 1844 Samuel Morse envió el primer mensaje público por telégrafo. |
| 5500 | But mental ideas or judgments are true or false, so how then can mental states (ideas or judgments) be natural processes? | Pero las ideas o juicios mentales son verdaderos o falsos, entonces, ¿cómo pueden los estados mentales (ideas o juicios) ser procesos naturales? |

Table 2: Sample sentences from the spa_Latn dataset with English source text and corresponding translations.

82.95. Appendix C summarizes the participants' responses.

## 4 Experimental Validation

Following the shared task's recommendation for experimental validation, we trained four bilingual machine translation models on the 6,193 newly collected Spanish sentences and evaluated their performance on the FLORES+ benchmark.

### 4.1 Data

**Italic experiments** To validate our model training setup, we reproduced the bilingual results reported by Maillard et al. (2023) for bidirectional translations between English and three Italic languages: Friulan (fur_Latn), Venetian (vec_Latn), and Ligurian (lij_Latn). We focus on these languages due to their linguistic relation to Spanish and their complete data availability on the SEED dataset and the FLORES+ benchmark.

**Spanish experiments** For our bidirectional English and Spanish (spa_Latn) machine translation models, we divided the collected data described in Section 3 into two versions: one before and one after the translation review process. This allowed us to analyze the scoring effect of a more streamlined review process based solely on proofreading and copy-editing in contrast to the iterative quality assurance pipeline implemented in FLORES-200 (NLLB Team et al., 2022). Table 3 summarizes the employed corpora with their corresponding source, size, and split.

| Language | Split | Lines | Source |
|----------|-------|-------|--------|
| eng, fur lij, vec | train | 6193 | Seed |
| | valid | 997 | FLORES+ |
| | test | 1012 | FLORES+ |
| spa | **train** | **6193** | **This work** |
| | valid | 997 | FLORES+ |
| | test | 1012 | FLORES+ |

Table 3: Corpora used in model experiments. Our contribution is highlighted in bold font.

### 4.2 Tokenization

We trained a SentencePiece model (Kudo and Richardson, 2018) on the train split for each language pair using a joined vocabulary of 8k tokens and byte-pair encoding (BPE) (Sennrich et al., 2016) for subword segmentation. In total, we trained three tokenizers for the bilingual Italic experiments and two tokenizers for the Spanish experiments, one for each version of the translated dataset.

### 4.3 Models

The machine translation models in this work are implemented with the fairseq toolkit (Ott et al., 2019) using the transformer architecture (Vaswani et al., 2017). Modifications are also made to match the bilingual model configurations in Maillard et al. (2023) for comparison purposes. The resulting architecture consists of 8 attention heads, 6 encoder and decoder layers, each with 4096-dimensional feedforward networks. We trained each model with an inverse square root learning rate of 0.001 and 400 warm-up updates. Training is conducted on a cloud virtual machine with an

NVIDIA L4 24GB GPU and an image preloaded with Debian 11, Python 3.10, PyTorch 1.13, and CUDA 11.3. Data preparation, model training, and evaluation recipes are available.[7]

**Italic models** We train two models per language pair, one for each direction (eng_Latn ↔ xxx_Latn) between English and the three selected Italic languages from the Seed dataset. We use the dev split of the FLORES+ benchmark for validation and the highest BLEU score (Papineni et al., 2002) on this split as the checkpoint selection criterion. Training is stopped when the validation BLEU score fails to improve after 10,000 gradient updates, and the selected checkpoint is used to calculate the chrF++ scores (Popović, 2015) on the FLORES+ devtest split. These scores provide a baseline for guiding our training parameters until achieving performance on par with the results reported by Maillard et al. (2023). This enables comparing the metrics from our Spanish models and assessing the quality of our spa_Latn data contribution.

**Spanish models** Using only the newly collected Spanish data, we trained four models, two for each version of the spa_Latn dataset. Model training and architecture parameters were defined during the Italic experiments and remained constant for these models. For validation, we used the dev split from the FLORES+ dataset. Training was conducted for 2,000 epochs (averaging a total runtime of 12h 31m), with the best checkpoint selected based on the highest validation BLEU score.

### 4.4 Results

**Italic experiments** We evaluated all model hypotheses using the sacrebleu tool (Post, 2018) against the devtest split of the FLORES+ benchmark. Table 4 compares the original performance of bilingual machine translation models reported by Maillard et al. (2023) with our reproduction attempts, which employed a similar model architecture and training routine. Our reproduced models nearly match the average chrF++ score for the English-to-Italic direction, falling short by 0.2 points while showing an improvement of 2.8 chrF++ points for the Italic-to-English direction.

**Spanish experiments** We achieved a chrF++ score of 35.0 for English-to-Spanish translation

| Language | Original | | Reproduction | |
|---|---|---|---|---|
| | eng→ | →eng | eng→ | →eng |
| fur_Latn | 35.4 | 35.6 | 35.7 | 36.8 |
| lij_Latn | 34.1 | 32.1 | 33.4 | 36.0 |
| vec_Latn | 33.5 | 32.3 | 33.2 | 35.5 |
| *Average* | 34.3 | 33.3 | 34.1 | 36.1 |

Table 4: Performance comparison (chrF++) between original (Maillard et al., 2023) and our reproduced bilingual models for three Italic languages (fur_Latn, lij_Latn, vec_Latn).

and 34.7 for the reverse direction by training exclusively on the collected spa_Latn data. The average score of 34.9 is comparable to the 35.1 mean obtained by the Italic models trained on existing SEED corpora. This result suggests that the new Spanish training data is representative of the spa_Latn data in the FLORES+ benchmark.

Analyzing the effect of the translation review process, we observed an average improvement of 0.3 chrF++ points. Specifically, the English-to-Spanish model trained on the reviewed data improved from 34.5 to 35.0, while the reverse direction decreased slightly from 34.8 to 34.7. Figure 2 breaks down the performance of the four bilingual Spanish models.



Figure 2: Performance (chrF++) of the eng_Latn↔ spa_Latn bilingual models trained on two versions of Spanish data: before (*Initial*) and after the translation review process (*Reviewed*).

## 5 Discussions

**Seed English corpus** Each line in the English corpus is an excerpt from a Wikipedia article,

which may consist of complete sentences or fragments. Translators identified two primary challenges when working with this data: incomplete sentences and a lack of context due to changes in the original article. For example, segment 5540, "By way of example, they provide two proofs of the irrationality of ." is missing an object at the end. Translations of such sentences inevitably reflect the original issues.

Furthermore, Wikipedia articles support versioning,[8] so including the date of compilation in the metadata or augmenting the dataset with provenance information could enable the correct context retrieval at the time of translation.

**Seed-CAT** The SEED-CAT application facilitated the translation of the English corpus, the review of translations, and the generation of the provenance dataset. Translators rated its usability highly, with an "A" grade based on the Sauro–Lewis curved grading scale (Lewis and Sauro, 2018). Notably, the system's perceived learnability, identified in Alotaibi (2020) and Zaretskaya et al. (2017) as a key area for improvement in other CAT systems, scored the highest in our study, with an average of 93.2. This result underscores our efforts in user-centered design to make participation by language communities more accessible.

**Translation workflow** During the review phase, translators proofread and copy-edited the entire spa_Latn dataset, modifying 686 sentences and 1,815 words. Although the cost of this phase was lower than the initial translation, they were still comparable. Given the marginal improvement in the post-review model's metric and the significant impact of high-quality parallel sentences on machine translation performance (Maillard et al., 2023), teams should consider allocating review resources toward generating more translations. In our case, this approach could have generated 3,498 additional translations of similar average length.

## 6 Conclusions

This work presented the SEED-CAT application and its role in expanding the SEED dataset with a professionally translated Spanish corpus. By integrating a provenance data model and its serialization in SEED-CAT, we automatically obtained a

---

[8] https://en.wikipedia.org/wiki/Help:Page_history

JSON dataset detailing the origin of each translation. Our experimental machine translation validation on the FLORES+ benchmark demonstrates that the collected Spanish data is of high quality, achieving on-par performance with other established language pairs in the SEED dataset. The excellent grade from the system usability scale survey suggests that the SEED-CAT application has the potential to facilitate the inclusion of additional languages in future efforts.

## Limitations

To effectively support the expansion of the SEED dataset, localizing the SEED-CAT user interface is essential. Future translation projects may involve translators who work in languages other than English. Identifying these relevant languages and implementing the UI localization requires further work. Similarly, the machine translation feature is constrained by the availability of open models and their supported translation directions.

While the provenance dataset includes timestamps for when activities are performed, this information is not a reliable source for measuring the time taken to translate a sentence or other similar metrics. Users may experience interruptions, and the system does not track user engagement or attention.

Latin American Spanish varieties exhibit dialectal divisions that affect morphosyntactic features such as word order and verb tense (Hualde et al., 2012). Our data collection methodology does not distinguish between these variations. However, the specific variety spoken by each translator who participated in the project is detailed in Appendix D.

## Acknowledgements

## References

Nouman Ahmed, Natalia Flechas Manrique, and Antonije Petrović. 2023. Enhancing Spanish-Quechua machine translation with pre-trained models and diverse data sources: LCT-EHU at AmericasNLP shared task. In *Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP)*, pages 156–162, Toronto, Canada. Association for Computational Linguistics.

Hind M. Alotaibi. 2020. Computer-assisted translation tools: An evaluation of their usability among arab translators. *Applied sciences*, 10(18):6295–.

Mikko Aulamo, Umut Sulubacak, Sami Virpioja, and Jörg Tiedemann. 2020. OpusTools and parallel corpus diagnostics. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3782–3789, Marseille, France. European Language Resources Association.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

Lynne Bowker and Desmond Fisher. 2010. Computer-aided translation. *Handbook of translation studies*, 1:60–65.

Eric Brill. 1992. A simple rule-based part of speech tagger. In *Third Conference on Applied Natural Language Processing*, pages 152–155, Trento, Italy. Association for Computational Linguistics.

J Brooke. 1996. Sus: A "quick and dirty" usability scale. *Usability Evaluation in INdustry/Taylor and Francis*.

Moussa Doumbouya, Baba Mamadi Diané, Solo Farabado Cissé, Djibrila Diané, Abdoulaye Sow, Séré Moussa Doumbouya, Daouda Bangoura, Fodé Moriba Bayo, Ibrahima Sory 2. Condé, Kalo Mory Diané, Chris Piech, and Christopher Manning. 2023. Machine translation for nko: Tools, corpora, and baseline results. In *Proceedings of the Eighth Conference on Machine Translation*, pages 312–343, Singapore. Association for Computational Linguistics.

Marcello Federico, Alessandro Cattelan, and Marco Trombetti. 2012. Measuring user productivity in machine translation enhanced computer assisted translation. In *Proceedings of the 10th Conference of the Association for Machine Translation in the Americas: Research Papers*, San Diego, California, USA. Association for Machine Translation in the Americas.

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online. Association for Computational Linguistics.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10.

Paul Groth and Luc Moreau. 2013. PROV-overview. W3C note, W3C. Https://www.w3.org/TR/2013/NOTE-prov-overview-20130430/.

Barry Haddow, Rachel Bawden, Antonio Valerio Miceli Barone, Jindřich Helcl, and Alexandra Birch. 2022. Survey of low-resource machine translation. *Computational Linguistics*, 48(3):673–732.

Elena Hernández Gómez. 2015. Diccionario panhispánico de dudas.

Jose Ignacio Hualde, Antxon. Olarrea, and Erin. O'Rourke. 2012. *The handbook of Hispanic linguistics*, 1st ed. edition. Blackwell handbooks in linguistics. Wiley-Blackwell, Chichester, West Sussex [England] ;.

Trung Dong Huynh, Michael O. Jewell, Amir Sezavar Keshavarz, Danius T. Michaelides, Huanjia Yang, and Luc Moreau. 2013. The prov-json serialization. Project report, W3C.

Philipp Koehn. 2009. A process study of computer-aided translation. *Machine translation*, 23(4):241–263.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

James Lewis. 2018. The system usability scale: Past, present, and future. *International Journal of Human-Computer Interaction*, pages 1–14.

James R. Lewis and Jeff Sauro. 2018. Item benchmarks for the system usability scale. *J. Usability Studies*, 13(3):158–167.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzmán. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330.

George A Miller. 1995. Wordnet: A lexical database for english. *Communications of the ACM*, 38(11):39–41.

Paolo Missier and Luc Moreau. 2013. PROV-dm: The PROV data model. W3C recommendation, W3C. Https://www.w3.org/TR/2013/REC-prov-dm-20130430/.

Eugene W. Myers. 2023. Ano(nd) difference algorithm and its variations. *Algorithmica*, 1(1–4):251–266.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Antoni Oliver, Mercè Vàzquez, Marta Coll-Florit, Sergi Álvarez, Víctor Suárez, Claudi Aventín-Boya, Cristina Valdés, Mar Font, and Alejandro Pardos. 2023. TAN-IBE: Neural machine translation for the romance languages of the Iberian peninsula. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 495–496, Tampere, Finland. European Association for Machine Translation.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, page 311–318, USA. Association for Computational Linguistics.

Aurélie Picton, Emmanuel Planas, and Amélie Josselin-Leray. 2017. Monitoring the use of newly integrated resources into cat tools: A prototype. In *Trends in e-tools and resources for translators and interpreters, Approaches to Translation Studies*, 45, pages 109–136. Brill Editions.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Surangika Ranathunga, En-Shiun Annie Lee, Marjana Prifti Skenduli, Ravi Shekhar, Mehreen Alam, and Rishemjit Kaur. 2023. Neural machine translation for low-resource languages: A survey. *ACM Comput. Surv.*, 55(11).

Real Academia Española. 2010. *Ortografía de la lengua española*. Espasa, Madrid.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Chelo Vargas-Sierra. 2019. Usability evaluation of a translation memory system. *Quaderns de Filologia - Estudis Lingüístics*, 24:119.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Anna Zaretskaya, Gloria Corpas Pastor, and Míriam Seghiri. 2017. Chapter 2: User perspective on translation tools: Findings of a user survey. In *Trends in E-Tools and Resources for Translators and Interpreters*, pages 37–56. Brill.

## A    Performance of bilingual models

Table 5 summarizes the performance on the FLORES+ benchmark of the bilingual machine translation models using both BLEU and chrF++ metrics.

## B    Machine translation examples

Table 6 presents three sample machine translations generated by the eng-spa model trained solely on the 6,193 reviewed translations of the spa_Latn data.

## C    System usability scale

Table 7 details the responses of each translator to the system usability scale (SUS) survey. The columns correspond to each numbered statement as they appear in the standard questionnaire (Brooke, 1996), while the rows represent the translators in no particular order. The table also summarizes the average score per participant, the score per question, and the total SUS score.

| Model | BLEU | | chrF++ |
|---|---|---|---|
| | valid | test | test |
| eng→fur | 10.3 | 10.4 | 35.7 |
| eng←fur | 10.8 | 10.0 | 36.8 |
| eng→lij | 7.5 | 8.0 | 33.4 |
| eng←lij | 9.6 | 9.3 | 36.0 |
| eng→vec | 7.0 | 6.3 | 33.2 |
| eng←vec | 9.9 | 9.4 | 35.5 |
| *Average* | 9.2 | 8.9 | 35.1 |
| spa←eng | 8.4 | 8.1 | 35.0 |
| spa→eng | 7.2 | 7.2 | 34.7 |
| *Average* | 7.8 | 7.7 | 34.9 |

Table 5: Performance of the bilingual models evaluated using automatic metrics on the `valid` and `test` splits.

## D Spanish varieties

Table 8 relates each translator identifier in the provenance dataset with their specific Latin American Spanish variety.

## E Provenance graphs

Figures 3 and 4 depict the provenance graphs of two translations. The translation process for each sentence can vary significantly, leading to graphs of different complexity. These graphs were generated using the Python `prov` package.[9]

---

[9] https://github.com/trungdong/prov

| # | English | Spanish |
|---|---------|---------|
| 663 | This behavior oftentimes results in rifts between the leaders and the rest of the team. | Este comportamiento de los resultadosmes inestables entre los líderes de los líderes y el resto del equipo. |
| 702 | European influence and colonialism began in the 15th century, as Portuguese explorer Vasco da Gama found the Cape Route from Europe to India. | La influencia europea y el colonialismo comenzó en el siglo XV, como Portugués Verés Vasco Gama fundó la Cautela de Europa desde Europa. |
| 1009 | Japanese work culture is more hierarchical and formal that what Westerners may be used to. | La cultura japonés es más jerárquica y que se puede utilizar en los occidentales. |

Table 6: Sample machine translations from the eng-spa bilingual model. The English source sentences are drawn from the devtest split of the FLORES+ benchmark.

| User | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | *Score* |
|------|---|---|---|---|---|---|---|---|---|----|---------|
| 1 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 100.00 |
| 2 | 2.00 | 3.00 | 3.00 | 3.00 | 3.00 | 3.00 | 3.00 | 3.00 | 3.00 | 2.00 | 70.00 |
| 3 | 3.00 | 4.00 | 4.00 | 4.00 | 3.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 95.00 |
| 4 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 100.00 |
| 5 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 100.00 |
| 6 | 4.00 | 1.00 | 4.00 | 1.00 | 4.00 | 1.00 | 4.00 | 4.00 | 4.00 | 4.00 | 77.50 |
| 7 | 2.00 | 3.00 | 4.00 | 4.00 | 2.00 | 2.00 | 4.00 | 3.00 | 3.00 | 4.00 | 77.50 |
| 8 | 1.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 4.00 | 1.00 | 4.00 | 4.00 | 85.00 |
| 9 | 2.00 | 3.00 | 3.00 | 4.00 | 2.00 | 2.00 | 4.00 | 4.00 | 4.00 | 4.00 | 80.00 |
| 10 | 2.00 | 2.00 | 2.00 | 3.00 | 2.00 | 1.00 | 3.00 | 2.00 | 2.00 | 4.00 | 57.50 |
| 11 | 2.00 | 3.00 | 4.00 | 4.00 | 1.00 | 4.00 | 2.00 | 1.00 | 3.00 | 4.00 | 70.00 |
| *Score* | 68.18 | 79.55 | 90.91 | 88.64 | 75.00 | 75.00 | 90.91 | 77.27 | 88.64 | 95.45 | **82.95** |

Table 7: System usability scale scores for each translator (normalized).

| Translator ID | Variety | Glottocode |
|---------------|---------|-----------|
| 14a33724-59b6-45f3-b056-f9d384e48a59 | Caribbean Spanish | cari1288 |
| 2460a2a5-1a59-4e0a-afff-a83be7af3872 | Caribbean Spanish | cari1288 |
| d67b54ab-6325-47be-b578-02f4b7ba942c | Chilean Spanish | chil1286 |
| 599ec44e-1b13-4f0c-a71f-296bbf0f2c6a | Mexican Spanish | mexi1248 |
| ef29b2b9-ecc8-4766-95a7-40b794d0053f | Mexican Spanish | mexi1248 |
| 548b0e62-71a4-448c-ab47-96f58f81a935 | Rioplatense Spanish | riop1234 |
| 237fa953-c66e-4d5c-9f5a-919b171766be | Venezuelan Spanish | vene1262 |
| 142058e1-0375-4b16-bcc3-655af871ff1c | Venezuelan Spanish | vene1262 |
| 8fa01aed-835b-4912-b648-c86ae67e3599 | Venezuelan Spanish | vene1262 |
| 250663c9-8d8e-43da-a116-840b8cf39cf4 | Venezuelan Spanish | vene1262 |
| e730a639-0928-4801-a97b-f070e661dff9 | Venezuelan Spanish | vene1262 |

Table 8: Translator identifiers and their corresponding Latin American Spanish varieties with Glottocodes.

Figure 3: This provenance graph represents a simple workflow in which the translator consulted the original English text and translated it into Spanish in a single, continuous edit.



Figure 4: This provenance graph represents a workflow that begins with an initial machine translation, followed by two rounds of copy-editing.

# Efficient Terminology Integration for LLM-based Translation in Specialized Domains

**Sejoon Kim, Mingi Sung, Jeonghwan Lee, Hyunkuk Lim, Jorge Froilan Gimenez Perez**
PwC Korea GenAI Team, Seoul, South Korea
*{sejoon.s.kim, mingi.sung, jeonghwan.lee, hyunkuk.lim, gimenez.jorge}@pwc.com*

## Abstract

Traditional machine translation methods typically involve training models directly on large parallel corpora, with limited emphasis on specialized terminology. However, In specialized fields such as patent, finance, or biomedical domains, terminology is crucial for translation, with many terms that needs to be translated following agreed-upon conventions. In this paper we introduce a methodology that efficiently trains models with a smaller amount of data while preserving the accuracy of terminology translation. We achieve this through a systematic process of term extraction and glossary creation using the Trie Tree algorithm, followed by data reconstruction to teach the LLM how to integrate these specialized terms. This methodology enhances the model's ability to handle specialized terminology and ensures high-quality translations, particularly in fields where term consistency is crucial. Our approach has demonstrated exceptional performance, achieving the highest translation score among participants in the WMT patent task to date, showcasing its effectiveness and broad applicability in specialized translation domains where general methods often fall short.

## 1 Introduction

Conventional approaches to machine translation typically rely on training models using extensive parallel corpora, with little focus on specialized vocabulary. While this can be an effective approach in general, it demands large amounts of data and may lead to inconsistent translations of technical or domain-specific terminology. This challenge is particularly acute in specialized fields, where precise terminology usage is crucial and high-quality training data is often scarce (Skianis et al., 2020; Ghazvininejad et al., 2023; Zhang et al., 2023). Datasets for training models in these specialized domains are usually limited, and even when they exist, many are private due to security concerns.



Figure 1: Training method in terminology-based LLM translation

Consequently, certain industries lag behind in the advancement of deep learning-based translation. This disparity is even more pronounced for less commonly spoken languages, where specialized translation capabilities are significantly underdeveloped, resulting in an unequal distribution of progress in neural machine translation.

Numerous approaches have been explored to integrate terminology constraints into Neural Machine Translation (NMT) systems, aiming to improve domain-specific translation quality. Recent research on terminology-based machine translation has shifted towards incorporating constraints during the training phase, which eliminates the computational overhead during inference and enhances translation quality. Dinu et al. (2019) introduced a method where NMT models are trained with augmented datasets that include terminology constraints as inline annotations, allowing the model to learn the appropriate use of these terms during training. Building on this, Ailem et al. (2021) proposed further enhancements by using token masking and a modified cross-entropy loss function, which biases the model towards generating constraint terms more effectively. Additionally, the use of large language models for post-translation re-

finement has been explored to improve terminology recall, demonstrating the evolving nature of terminology integration in NMT (Bogoychev and Chen, 2023; Ghazvininejad et al., 2023; Moslem et al., 2023). These training-based approaches have demonstrated significant improvements in both BLEU scores and terminology usage rates compared to decoding-time methods, indicating their effectiveness in satisfying lexical constraints without compromising translation quality.

In this paper, we propose a fine-tuning approach to resolve the domain-specific terminology mismatch problem using only a small dataset. Our approach focuses on extracting a glossary from the existing training datasets and fine-tuning the model to integrate these terms effectively into translations. First, we train a terminology extraction model to generate a glossary from existing training datasets, which we integrate into our trie data structure (Bodon and Rónyai, 2003). We then extract domain-specific terms from the source sentences using the tree structure and pass them along with the source texts to instruct our Large Language Model (LLM) to effectively incorporate specialized terminology into translations. This approach ensures high-quality and consistent results in specialized fields. Figure 1 illustrates how our approach differs from traditional fine-tuning methods. This targeted refinement process enhances the model's capacity to manage specialized terminology, thereby maximizing the utility of the original training data and significantly improving translation accuracy and consistency. Our methodology has proven to be exceptionally effective, particularly in specialized translation tasks, where general translation methods often struggle to maintain accuracy and consistency. Notably, our approach achieved the highest translation score among all participants in the WMT patent task, underscoring its superior performance and broad applicability across various specialized translation domains. Through this systematic and targeted strategy, we ensure that our translations are not only accurate but also contextually relevant, thereby providing a reliable solution for specialized translation needs.

## 2 Methodology

In this section, we describe the methodology employed in developing a domain-specific terminology-based LLM translation system, focusing on three key processes: (1) the creation of a

terminology glossary, (2) the identification of terms within the source text, and (3) the application of these terms during the translation process using LLM prompts and sLLM fine-tuning.

### 2.1 Construction of the Terminology Glossary: Terminology Aligner

**System Message:**
I will now show you source sentences in Japanese and target sentences in Korean. Your task is to extract and pair key terms from both the original and translation texts. Maintain the exact form of the terms without modification.
Please follow these instructions for extracting term pairs:

- Extract term pairs that are closely related to patents.
- Only extract nouns.
- The extracted term pairs will be used to create a Japanese-Korean glossary.
- Return the results in the form of a Python dictionary, as shown in the example.
- However, if the exact same term appears more than once include it only once.

**Example 1:**

src_sentence = それぞれについて官能評を行った結果を表４２に示す。

tgt_sentence = 각각에 대하여 관능 평가를 행한 결과를 표 42에 나타낸다.

result = {"官能評": "관능 평가"}

**Example 2:**

src_sentence = 各種の特許や技術標準化にする問題が討された。

tgt_sentence = 각종 특허권과 기술 표준화에 관한 문제가 검토되었다.

result = {"特許": "특허권", "技術標準化": "기술 표준화"}

Figure 2: Instructions for Term Extraction

To enable the translation model to produce accurate translations that incorporate specialized terminology, we first construct a "Terminology Pair Dictionary," aligning key terms between the source and target languages. We achieve this by fine-tuning

the Mistral Nemo model, creating a Terminology Aligner model whose primary objective is to extract pairs of key terms from both the original and translated texts.

For our training data, we leverage the GPT-3.5 API to generate synthetic data by crafting prompts that instruct the API to extract key term pairs from existing Japanese-Korean translation pairs in our dataset, along with the system prompt shown in Figure 2. From the 1,000,000 training samples provided by the organizers, we randomly select 1,000 examples to fine-tune Mistral Nemo for a single epoch. We adopt this conservative approach, recognizing that Mistral Nemo already possesses a robust grasp of both Korean and Japanese and is capable of performing various tasks, including the one at hand. Our goal is to specialize the model for our particular task without compromising its broader capabilities or confusing it with unrelated tasks.

Furthermore, when the entire dataset was used for fine-tuning, the model frequently extracts non-essential term pairs or entire sentences as pairs, indicating overfitting. By carefully selecting the amount of data and limiting the number of training epochs, we ensure that it extracts only the most relevant, domain-specific term pairs and effectively fine-tune the Terminology Aligner model.

## 2.2 Term Identification in the Source Text: Trie-Tree Algorithm



Figure 3: Overall process of term extraction to translation

The next step in our methodology involves identifying and extracting specialized terms from the source text that must be accurately translated using the glossary we constructed. To account for industries where there is often a high volume of technical terms and the need for efficient text scanning, we implement the Trie Tree data structure to extract the domain-specific terms.

The Trie Tree is particularly well-suited for this task due to its efficiency in string searching and matching. The algorithm operates by placing a cursor at the first Unicode character of the text, while another cursor points to the root of the tree. As the text cursor advances through each character, the tree cursor checks for corresponding child nodes. If a match is found, the tree cursor moves to the next node; if not, it resets to the root. When the cursor reaches a node marked as a 'term,' the term is identified, and its position is recorded. This allows us to quickly retrieve the term's translation and include it in the LLM prompt, ensuring that all relevant terms in the text are accurately and efficiently identified. The process is visually illustrated in Figure 3, which describes the step-by-step progression of the Trie Tree algorithm from text scanning to term retrieval and integration into the LLM prompt.

## 2.3 Application of Terms in Translation: LLM Prompting and sLLM Fine-Tuning



**System Message:**
You are a professional translator. You are especially familiar with specialized patent knowledge and terms in chemistry, electricity, mechanical engineering, and physics, as well as general everyday terms. Translate the following Japanese source text into Korean.

- Refer to the word pairs in the glossary when you translate.
- Do not translate the glossary itself.
- Do not include anything but translation result only.
- If a term in the glossary has multiple possible translations separated by '|', choose the most appropriate one.
- The translation result must be written in a single line. There must be no newline character at the end.

**Glossary:**
{セレノール化合物 : 셀레놀 화합물,
端部 : 끝부분 | 단부 | 모서리 ,
絶膜 : 절연막,
送信回路 : 송신 회로 | 전송 회로}

Figure 4: Instructions for Term Extraction

The final phase of our methodology involves the use of the extracted terminology during the translation process. To do this, we first extract term pairs from all translation pairs in our dataset using the created tree structure. These extracted term pairs are then combined with each original translation pair and the system message in Figure 4 to create an instruction-based training dataset to fine-tune our translation model.

Similar to our fine-tuning process with the Terminology Aligner, we observed that both the amount of data and the number of training epochs significantly influence the quality of the translation output, particularly in terms of how natural the translations sound. Interestingly, when working with smaller datasets, the model tends to produce more natural, conversational translations. However, as the dataset size increases, the model increasingly adheres to the original sentence structure, resulting in a more formal and literal style of translation.

To balance these tendencies, we use approximately 1,000 data points for training and limit the training to three epochs, with a temperature setting of 0.1. This configuration allows the model to generate translations that were both accurate and natural, making an effective use of the specialized terminology while maintaining a high level of fluency and readability.

## 3 Experimental Results and Application

|   | Team | BLEU (mecab) | RIBES |
|---|------|--------------|-------|
| 1 | GenAI | **70.60** | **0.939073** |
| 2 | Chatgpt (w/ glossary) | 69.00 | 0.929945 |
| 3 | sakura | 68.00 | 0.926839 |
| 4 | Bering Lab | 66.25 | 0.925226 |
| 5 | ryan | 65.74 | 0.922837 |
| 6 | goku20 | 64.30 | 0.922486 |
| 7 | ORGANIZER | 62.43 | 0.915266 |
| 8 | tpt_wat | 61.00 | 0.918436 |
| 9 | Chatgpt (w/o glossary) | 59.90 | 0.908637 |

Table 1: BLEU (mecab) and RIBES scores for the Japanese-to-Korean translation task

Our proposed methodology has been rigorously tested and evaluated within the framework of the WMT patent task, where it achieves the highest translation score to date among all participants. This success demonstrates the effectiveness of our approach in handling domain-specific translations,

particularly in maintaining consistency in terminology.

In addition to the translation results generated by our model, we submitted two additional translations using ChatGPT. The first result, labeled 'ChatGPT (w/ glossary)' in Tables 1 and 2, was obtained by replacing our model with ChatGPT while keeping the system prompt and glossary identical to our methodology. The second result was generated using ChatGPT alone without any additional inputs.

Several interesting findings emerged: for the Japanese-to-Korean translation task, ChatGPT without the glossary scores lower than other models in the patent translation domain. However, the score significantly improves when our glossary is provided. This demonstrates that the integration of a terminology glossary substantially enhances translation performance, regardless of the underlying model's capabilities. By comparing ChatGPT with and without the glossary, it becomes evident that our system effectively boosts translation quality through efficient terminology integration. Our specialized language model, trained specifically to use the glossary, outperforms ChatGPT even with the glossary. Upon reviewing the outputs, we notice that ChatGPT sometimes fails to correctly apply terms inside the glossay and occasionally uses Japanese terms instead of their Korean equivalents in the Japanese-to-Korean translation.

These findings highlight that our model can be effectively trained with a small dataset, achieving high-quality translations while remaining a smaller, more efficient model. Beyond patent translation, our methodology can be extended to specialized fields such as legal and financial translation where accurate term alignment is critical, providing a robust solution where general translation methods may fall short.

## 4 Discussion

### 4.1 Advantages of Our Methodology Over Traditional Approaches

The effectiveness of our methodology is further underscored by several key advantages it holds over traditional approaches:

### 4.1.1 Focused Learning on Domain-Specific Terms

Traditional models typically assign equal importance to all words in the training data, which can result in inconsistent translations of specialized terms

| Rank | Team | BLEU | | | RIBES | | |
|---|---|---|---|---|---|---|---|
| | | juman | kytea | mecab | juman | kytea | mecab |
| 1 | **GenAI** | **67.00** | **67.40** | **66.90** | **0.924474** | **0.919657** | **0.923416** |
| 2 | Chatgpt (w/ glossary) | 62.20 | 62.50 | 61.90 | 0.916385 | 0.912133 | 0.914275 |
| 3 | Chatgpt (w/o glossary) | 61.60 | 62.50 | 61.50 | 0.912482 | 0.907932 | 0.911476 |
| 4 | EHR | 53.83 | 55.83 | 54.23 | 0.907358 | 0.903857 | 0.905654 |
| 5 | sarah | 53.59 | 55.68 | 53.94 | 0.903211 | 0.900313 | 0.902430 |
| 6 | KNU_Hyundai | 53.56 | 55.68 | 54.02 | 0.901627 | 0.900091 | 0.901877 |
| 7 | TMU | 52.85 | 54.92 | 53.24 | 0.906113 | 0.903179 | 0.906320 |
| 8 | Bering Lab | 52.74 | 54.55 | 53.15 | 0.902984 | 0.898627 | 0.902621 |
| 9 | ORGANIZER | 52.02 | 53.93 | 51.99 | 0.897348 | 0.896897 | 0.898316 |
| 10 | sakura | 51.90 | 54.10 | 52.30 | 0.899781 | 0.896489 | 0.898412 |

Table 2: BLEU and RIBES scores for the Korean-to-Japanese translation task

across different contexts. Our methodology addresses this by prioritizing domain-specific terms, ensuring they are recognized and used consistently in relevant translations.

### 4.1.2 Efficient Data Utilization through Terminology Extraction

Traditional methods often require large volumes of data to achieve satisfactory performance, particularly in specialized domains. Our method optimizes the use of training data by focusing on key term pairs and creating a dedicated glossary, enabling more efficient learning even with a smaller dataset.

### 4.1.3 Enhanced Translation Consistency and Accuracy

A common challenge with traditional translation methods is inconsistency in translating specialized terms, especially when these terms have multiple possible translations depending on context. Our approach mitigates this by ensuring the model is trained with a consistent set of term translations derived from the glossary.

### 4.1.4 Improved Model Generalization

Traditional models trained on large corpora may overfit to specific sentence structures or styles present in the training data, leading to poor generalization to new texts. Our approach incorporates the glossary into training, acting as a regularizing factor that improves generalization to new texts within the same domain.

### 4.1.5 Customizability for Different Domains

Our methodology allows for greater flexibility in adapting the model to different specialized fields. By updating the glossary with new terms relevant

to a particular domain, the model can be quickly tailored to perform well without extensive retraining.

## 5 Conclusion

Our terminology-based LLM translation methodology represents a significant advancement in the field of machine translation, particularly for specialized domains requiring precise and consistent term usage. By constructing a terminology glossary using the Terminology Aligner, implementing an efficient term identification process with a Trie Tree algorithm, and fine-tuning the translation process using LLM prompts, we present a system that not only improves translation accuracy but also maintains a high level of naturalness in the output. Our approach has proven successful in terms of performance, operational cost, and training data efficiency, showing great promise for a wide range of professional translation applications.

## References

Melissa Ailem, Jingshu Liu, and Raheel Qader. 2021. Encouraging neural machine translation to satisfy terminology constraints. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1450–1455, Online. Association for Computational Linguistics.

F. Bodon and L. Rónyai. 2003. Trie: An alternative data structure for data mining algorithms. *Mathematical and Computer Modelling*, 38(7):739–751. Hungarian Applied Mathematics.

Nikolay Bogoychev and Pinzhen Chen. 2023. Terminology-aware translation with constrained decoding and large language model prompting. In *Proceedings of the Eighth Conference on Machine*

*Translation*, pages 890–896, Singapore. Association for Computational Linguistics.

Georgiana Dinu, Prashant Mathur, Marcello Federico, and Yaser Al-Onaizan. 2019. Training neural machine translation to apply terminology constraints. *CoRR*, abs/1906.01105.

Marjan Ghazvininejad, Hila Gonen, and Luke Zettlemoyer. 2023. Dictionary-based phrase-level prompting of large language models for machine translation. *Preprint*, arXiv:2302.07856.

Yasmin Moslem, Gianfranco Romani, Mahdi Molaei, John D. Kelleher, Rejwanul Haque, and Andy Way. 2023. Domain terminology integration into machine translation: Leveraging large language models. In *Proceedings of the Eighth Conference on Machine Translation*, pages 902–911, Singapore. Association for Computational Linguistics.

Konstantinos Skianis, Yann Briand, and Florent Desgrippes. 2020. Evaluation of machine translation methods applied to medical terminologies. In *Proceedings of the 11th International Workshop on Health Text Mining and Information Analysis*, pages 59–69, Online. Association for Computational Linguistics.

Huaao Zhang, Qiang Wang, Bo Qin, Zelin Shi, Haibo Wang, and Ming Chen. 2023. Understanding and improving the robustness of terminology constraints in neural machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6029–6042, Toronto, Canada. Association for Computational Linguistics.

# A Appendix

In this appendix, we provide additional details on the training procedures, model configurations, and methodologies employed in our approach for efficient terminology integration in LLM-based translation within specialized domains.

## A.1 Training Details and Additional Information

For both the translation task and the terminology extraction task, we used the **Mistral-Nemo-Instruct-2407** model as our base language model. This model was selected due to its strong capability in following instructions, including tasks such as translation and terminology extraction.

### A.1.1 Training Details

**Model Configuration and Training Parameters**

| | |
|---|---|
| **Base Model** | mistralai/Mistral-Nemo-Instruct-2407 |
| **LoRA Adapter Settings** | |
| Alpha | 8 |
| Rank | 8 |
| Dropout Rate | 0.1 |
| Target Modules | ["q_proj", "k_proj", "v_proj", "o_proj", "gate_proj", "down_proj", "up_proj"] |
| **Learning Rate** | 1e-5 |
| **Optimizer** | AdamW |
| **Learning Rate Scheduler** | Linear |
| **Warmup Ratio** | 0.01 |
| **Epochs** | 1 |
| **Batch Size** | 4 |
| **Gradient Checkpointing** | Enabled |

LoRA was applied to mitigate GPU memory limitations and prevent catastrophic forgetting.

## A.2 Fine-Tuning Challenges and Considerations

Our preliminary experiments indicated that using smaller datasets for fine-tuning resulted in more effective performance for both the terminology aligner and the translation model. Based on these observations, we concluded that a smaller dataset was sufficient to format the model's outputs appropriately and guide it to produce task-specific responses without deviating from the desired content.

The Mistral-Nemo model already exhibited strong abilities in instruction following, including translation and terminology extraction. Therefore, extensive fine-tuning was unnecessary and could potentially degrade performance. Training with larger datasets led to overfitting, where the model's training loss decreased, but the actual translation quality did not improve. In some cases, the model exhibited issues like repetitive outputs. We attempted to mitigate overfitting by increasing dropout rates and weight decay. However, these adjustments did not yield significant improvements in our experiments.

# Rakuten's Participation in WMT 2024 Patent Translation Task

**Ohnmar Htun** and **Alberto Poncelas**
Rakuten Institute of Technology
Rakuten Group, Inc.
{ohnmar.htun alberto.poncelas}@rakuten.com

## Abstract

This paper introduces our machine translation system (team *sakura*), developed for the 2024 WMT Patent Translation Task. Our system focuses on translations between Japanese-English, Japanese-Korean, and Japanese-Chinese. As large language models have shown good results for various natural language processing tasks, we have adopted the *RakutenAI-7B-chat* model, which has demonstrated effectiveness in English and Japanese. We fine-tune this model with patent-domain parallel texts and translate using multiple prompts.

## 1 Introduction

Machine Translation (MT) systems are becoming increasingly important in the translation industry. While generic MT models are good at translating common phrases into everyday language, they often struggle with specialized domains unless they have been specifically tuned for those areas. Patent documents are an example of this specialized content.

The patent translation shared task[1] at Conference on Machine Translation (WMT) 2024 aims to bring together Natural Language Processing (NLP) researchers to assess and explore innovative methods for translating patents, specifically between Japanese (Ja) and English (En), Korean (Ko) or Chinese (Zh), and vice versa.

Recently, significant advancements have been made in the field of NLP due to the development of Large Language Models (LLMs). Unlike encoder-decoder models, which are typically created to perform a single task such as machine translation, LLMs are designed for multiple NLP purposes. As a result, LLMs are often pre-trained on larger and more diverse texts, which helps improve the model's language understanding. In our work, we propose using an LLM fine-tuned with parallel data to perform accurate translations in the patent domain.

An LLM that has been specifically adapted to multiple NLP tasks in both English and Japanese is the *RakutenAI-7B* (Rakuten Group, Inc. et al., 2024) model. It has been pre-trained on a large volume of data, and its tokenizer has been optimized for the character-per-token rate in Japanese, making it ideal for complex tasks such as Japanese translation.

We participated in the patent translation (*sakura* team) shared task. In our proposal, we fine-tune the LLM with patent-domain bilingual data to build a multilingual model that achieves high-quality translations in multiple language directions. In addition, we produce translations using multiple prompts to further boost performance.

## 2 Related Work

LLM has been explored in the patent industry for tasks such as claim generation (Jiang et al., 2024), Question-Answer, or Classification (Bai et al., 2024).

Regarding the patent translation, previous participants in the JPO shared task have explored various methodologies, including training encoder-decoder models as suggested by Park and Lee (2021), utilizing Transformer-based NMT model (Vaswani et al., 2017) with ensemble decoding (Susanto et al., 2019) and adapting pre-trained models such as BART (Lewis et al., 2020) mBART (Liu et al., 2020) with patent-specific data (Wang and Htun, 2020; Kim and Komachi, 2021).

## 3 Task Description

The shared task consists of translating a set of sentences from patent publications in the En ↔ Ja, Ko ↔ Ja and Zh ↔ Ja language directions. The

---

[1] https://www2.statmt.org/wmt24/patent-task.html

text belongs to the domains of Chemistry, Electricity, Mechanical Engineering or Physics.

These sentences, are organized as different test set according to the year the patents were published:

- *test-n1*: Published between 2011 and 2013 (same test sets used in the past years).

- *test-n2*: Published between 2016 and 2017 (not available for Ko-Ja).

- *test-n3*: Published between 2016 and 2017 (but target sentences were manually created by translating source sentences).

- *test-n4*: Published between 2019 and 2020.

- *test-2022*: The union of the previous n1 to n4 sets.

The *test-n1* to *test-n4* vary in size from 2K to 5K sentences depending on the language. The only exception is *test-n3*, which was created manually and contains between 200 and 700 sentences. The total size of these tests, i.e. *test-2022*, ranges from 7K to 10K sentences.

## 3.1 Evaluation

In order to determine the performance of our model, we submit the translation of the test sets mentioned above. The results of the different tasks are published in https://lotus.kuee.kyoto-u.ac.jp/WAT/evaluation/index.html.

The translations are tokenized using Juman (Kurohashi and Kawahara, 2009), KyTea[2], Mecab (Kudo, 2005) or Moses tokenizer[3]. The website presents multiple evaluation metrics for evaluating the translation. In this paper we present only the BLEU (Papineni et al., 2002) scores. The other metrics are correlated with BLEU. We refer to their website for the rest of the metrics.

## 3.2 Training Data

The organizers of the shared task also provide the JPO Patent Corpus (JPC) for training. This is a dataset built by the Japan Patent Office[4] consisting of sets of 1M parallel sentences for each language pair (English-Japanese, Chinese-Japanese and Korean-Japanese).

Figure 1: Performance of the fine-tuned model on the dev set using different beam sizes for decoding.

The data also include a dev set in the same domain with around 2K sentences each.

## 4 Experimental Settings

For our experiments, we fine-tune the *RakutenAI-7B-chat*[5] model, which has been optimized for the English and Japanese languages. However, it has not been explicitly adapted for other languages such as Korean and Chinese.

We use the JPO data described in Section 3.2 for this fine-tuning and do not incorporate any additional data other than what has been provided by the shared task organizers. The training process involves 200K steps with a batch size of 8. We fine-tune the model using the prompt "*Translate the following English text to Japanese:*" appending the sentence to be translated and replacing the source and target languages as needed for each language direction.

### 4.1 Influence of Beam Size

For decoding, we chose a beam size of three. While larger beam sizes involve considering more candidate translations, this does not always result in better performance. We tested our fine-tuned model with beam sizes of 1, 3, 5, and 7 on the development set. The results, measured using CHRF (Popović, 2015) metric, are shown in Figure 1. Although there is no single optimal beam size, our findings indicate that increasing the beam size beyond three does not lead to significant improvements and in some cases it may even degrade performance.

| Test | Direction | BLEU | Δ |
|------|-----------|------|---|
| *test-2022* | En → Ja | 53.4 | +4.5 |
| | Ja → En | 50.1 | +5.6 |
| *test-n1* | En → Ja | 51.1 | +5.8 |
| | Ja → En | 49.3 | +5.2 |
| *test-n2* | En → Ja | 46.3 | +5.7 |
| | Ja → En | 43.9 | +6.2 |
| *test-n3* | En → Ja | 54.9 | +7.4 |
| | Ja → En | 43.1 | +8.1 |
| *test-n4* | En → Ja | 62.1 | +1.6 |
| | Ja → En | 59.7 | +4.8 |

Table 1: BLEU scores for Japanese-English translation (using Moses tokenizer for Ja → En and kytea for En → Ja). The column Δ indicates the difference between the scores of our model and those of the organizers.

| Test | Direction | BLEU | Δ |
|------|-----------|------|---|
| *test-2022* | Zh → Ja | 56.6 | +5.5 |
| | Ja → Zh | 46.2 | +1.5 |
| *test-n1* | Zh → Ja | 53.4 | +6.7 |
| | Ja → Zh | 41.7 | +2.6 |
| *test-n2* | Zh → Ja | 51.3 | +5.3 |
| | Ja → Zh | 40.6 | +1.5 |
| *test-n3* | Zh → Ja | 21.8 | +4.0 |
| | Ja → Zh | 27.0 | +3.2 |
| *test-n4* | Zh → Ja | 68.7 | +3.7 |
| | Ja → Zh | 58.7 | +1.2 |

Table 2: BLEU scores for Japanese-Chinese translation (using Kytea tokenizer). The column Δ indicates the difference between the scores of our model and those of the organizers.

| Test | Direction | BLEU | Δ |
|------|-----------|------|---|
| *test-2022* | Ko → Ja | 74.3 | +0.4 |
| | Ja → Ko | 75.4 | +2.6 |
| *test-n1* | Ko → Ja | 73.3 | +1.6 |
| | Ja → Ko | 72.9 | +2.2 |
| *test-n3* | Ko → Ja | 52.3 | +0.3 |
| | Ja → Ko | 68.0 | +5.6 |
| *test-n4* | Ko → Ja | 77.3 | +0.2 |
| | Ja → Ko | 78.4 | +3.7 |

Table 3: BLEU scores for Japanese-Korean translation (using Mecab tokenizer). The column Δ indicates the difference between the scores of our model and those of the organizers.

## 4.2 Influence of the Prompt

At decoding time, we perform multiple translations using different variations of the prompt. The prompts are the following:

- *Translate the following English text to Japanese* (same as training data)

- *Translate the following English sentence to Japanese:* (replace "text" with "sentence")

- *Translate the following text to Japanese:* (omit the source language)

- *Translate the text to Japanese:* (above prompt rephrased)

- *Translate the following English patent text to Japanese:* (explicitly indicate that is a patent text)

We use LASER (Heffernan et al., 2022) scores compared to the source to retrieve the best translation. Although all of them are similar, there are small nuances that can increase the quality by around 0.5-1 BLEU points.

## 5 Results

In this section we present the translation performance achieve by our model on the different language directions.

## 5.1 Japanese-English

First, Table 1 illustrates the performance of our model on English-Japanese translation. We observe that our model achieves the best results for

this language pair compared to the model of the organizers, with an average improvement of 5 BLEU points for English-to-Japanese and 6 BLEU points for Japanese-to-English. Furthermore, it shows greater improvements in this pair compared to the other language pairs. This success can be attributed to the fact that our model was pre-trained on these two languages, benefiting from higher exposure.

## 5.2 Japanese-Chinese

Table 2 presents the results for Chinese-Japanese translation. Although improvements are observed across all test sets, there is a notable disparity between the language directions. While Zh → Ja shows an improvement of 5 BLEU points, for the reverse direction there is an improvement of 1.5 BLEU points.

## 5.3 Japanese-Korean

Lastly, in Table 3 we show the performance of the Japanese-Korean translation. For this language pair we achieve smaller improvements when compared to the baseline of the organizers.

## 6 Conclusion

In this paper, we described our MT model developed for the 2024 WMT Patent Translation Task, specifically for English-Japanese, Japanese-Korean, and Japanese-Chinese translations. The ranking system has evaluated participating teams every year from 2016 to 2024. Our model achieved first place in 20 out of the 28 tasks without using external data. Our approach involves fine-tuning the "*RakutenAI-7B-chat*" model using sentences from the patent domain and decoding with multiple prompts. Although this model was originally pre-trained only on English and Japanese data, fine-tuning with Korean and Chinese text has led to good translation performance, surpassing the models submitted in previous years for the same task.

## References

Zilong Bai, Ruiji Zhang, Linqing Chen, Qijun Cai, Yuan Zhong, Cong Wang Yan Fang, Jie Fang, Jing Sun, Weikuan Wang, Lizhi Zhou, et al. 2024. PatentGPT: A Large Language Model for Intellectual Property. *arXiv preprint arXiv:2404.18255*.

Kevin Heffernan, Onur Çelebi, and Holger Schwenk. 2022. Bitext Mining Using Distilled Sentence Representations for Low-Resource Languages. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2101–2112, Abu Dhabi, United Arab Emirates.

Lekang Jiang, Caiqi Zhang, Pascal A Scherz, and Stephan Goetz. 2024. Can Large Language Models Generate High-quality Patent Claims? *Preprint*, arXiv:2406.19465.

Hwichan Kim and Mamoru Komachi. 2021. TMU NMT system with Japanese BART for the patent task of WAT 2021. In *Proceedings of the 8th Workshop on Asian Translation*, pages 133–137, Bangkok, Thailand.

Taku Kudo. 2005. Mecab: Yet another part-of-speech and morphological analyzer. *https://taku910.github.io/mecab/*.

Sadao Kurohashi and Daisuke Kawahara. 2009. Japanese Morphological Analysis System JUMAN 7.0 Users Manual. *http://nlp.ist.i.kyoto-u.ac.jp*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual Denoising Pre-training for Neural Machine Translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA.

Heesoo Park and Dongjun Lee. 2021. Bering Lab's Submissions on WAT 2021 Shared Task. In *Proceedings of the 8th Workshop on Asian Translation*, pages 141–145, Bangkok, Thailand.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal.

Rakuten Group, Inc., Aaron Levine, Connie Huang, Chenguang Wang, Eduardo Batista, Ewa Szymanska, Hongyi Ding, Hou Wei Chou, Jean-François Pessiot, Johanes Effendi, Justin Chiu, Kai Torben Ohlhus, Karan Chopra, Keiji Shinzato, Koji Murakami, Lee Xiong, Lei Chen, Maki Kubota, Maksim Tkachenko, Miroku Lee, Naoki Takahashi, Prathyusha Jwalapuram, Ryutaro Tatsushima, Saurabh Jain, Sunil Kumar Yadav, Ting Cai, Wei-Te Chen, Yandi Xia, Yuki Nakayama, and Yutaka Higashiyama. 2024. RakutenAI-7B: Extending Large Language Models for Japanese. *Preprint*, arXiv:2403.15484.

Raymond Hendy Susanto, Ohnmar Htun, and Liling Tan. 2019. Sarah's Participation in WAT 2019. In *Proceedings of the 6th Workshop on Asian Translation*, pages 152–158, Hong Kong, China.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *Advances in Neural Information Processing Systems*, volume 30.

Dongzhe Wang and Ohnmar Htun. 2020. Goku's Participation in WAT 2020. In *Proceedings of the 7th Workshop on Asian Translation*, pages 135–141, Suzhou, China.

# The SETU-ADAPT Submissions for WMT24 Biomedical Shared Task

**Antonio Castaldo**[a*]**, Maria Zafar**[*]**, Prashanth Nayak**[b]**,**
**Rejwanul Haque**, **Andy Way**[c]**, Johanna Monti**[a]

[a]UNIOR NLP Research Group, University of Naples "L'Orientale", Naples, Italy
South East Technological University, Carlow, Ireland
[b]KantanAI, Dublin, Ireland
[c]ADAPT Centre, Dublin City University, Dublin, Ireland
antonio.castaldo@phd.unipi.it,c00304029@setu.ie,pnayak@kantanai.io
rejwanul.haque@setu.ie,andy.way@adaptcentre.ie,jmonti@unior.it

## Abstract

This paper presents SETU-ADAPT's submissions to the WMT 2024 Biomedical Shared Task, where we participated for the language pairs English-to-French and English-to-German. Our approach focused on fine-tuning Large Language Models (LLMs), using in-domain and synthetic data, employing different data retrieval strategies. We introduce a novel MT framework, involving three autonomous agents: a Translator Agent, an Evaluator Agent and a Reviewer Agent. We present our findings and report the quality of the outputs.

## 1 Introduction

Translating texts in the biomedical domain presents unique challenges that sets it apart from general domain translation tasks. The domain is characterised by the use of specialised terminology, fixed expressions and relative data scarcity. In recent times, LLMs (Brown et al., 2020; Przystupa and Abdul-Mageed, 2019) have become the go-to systems for building Machine Translation (MT) systems, due to their impressive performance in generating accurate translations across diverse domains. Precisely, the ability to fine-tune these models on new data, adapting them to the specialised terminology used in the biomedical domains, makes them particularly suitable for our task.

In our experiments, we built our MT systems using Llama-3 (Dubey et al., 2024) and No Language Left Behind (NLLB) (Costa-jussà et al., 2022), based on the high performance reported in their relevant publications. We further design and develop strategies to address data scarcity and improve the quality of the outputs. Our first approach

involves back-translation (Xu et al., 2019), where we leverage monolingual data and translate them back into the source language, thus generating synthetic data to be combined with the original dataset. This approach is widely recognized as an effective method to overcome the challenges caused by the translation of low-resource languages and specific domains. Another data augmentation method that we adopt is based on terminology-aware mining (Haque et al., 2020), where we extract a terminology list from our training data and use it to mine semantically similar sentences from the general domain corpus. We further experimented using few-shot prompting, where we provided the model with a few translation samples retrieved through semantic search based on the source sentence. Finally, we propose an innovative MT system powered by `GPT-4o` (OpenAI et al., 2024) that employs an agentic workflow (Wang et al., 2024). This system follows a collaborative framework, where three LLM-based agents work together autonomously to produce translations.

The paper is organised as follows. We present an overview of our proposed systems in Section 3. We describe our datasets and our data augmentation strategies in Section 4 and Section 5. We introduce our last system, involving LLM-based autonomous agents in Section 6. We present the results of our evaluation in Section 7 and draw our conclusions in Section 8.

## 2 Related Work

The main difficulties found in biomedical MT have been the highly specialised domain, the lack of relevant data, and the importance of using the correct terminology. To address the issues caused by

---

*Both authors are equal contributors to this work.

domain-specific terminology, Choi et al. (2022) adopted a soft-constrained translation approach, where terminology constraints retrieved from the training corpus are provided to the MT system as a suggestion rather than a hard constraint. Soft-constrained decoding appears to be a promising solution to drive the systems to include the necessary terminology in the output while preserving the model's fluency and flexibility in the translation.

Ballier et al. (2022) trained different systems on a selection of texts from WMT, Khresmoi (Dušek et al., 2017) and UFAL (Bojar et al., 2017) datasets, comparing the results. Interestingly, they find that mBART-50 (Tang et al., 2021), despite producing fluent grammatical sentences, fails at translating consistently domain-specific terminology. Their study suggests that this well-known model may not be adequate for the task of biomedical translation, especially in the context of translating biomedical abstracts where its small context window may cause inaccurate translations.

Manchanda and Bhagwat (2022) confirms previous studies that showed how fine-tuning any model from a general domain to a specialised one, as is the case with clinical and biomedical texts, improves the translation quality in most cases. Their study introduces a novel approach, based on combining general-purpose and domain-specific datasets for fine-tuning while applying a higher learning rate to the general domain data. Their experiments demonstrate how this combined fine-tuning approach may improve translation quality in both domains.

In the last few years, we have seen a general surge of LLMs applied to MT (Hendy et al., 2023; He et al., 2024). Several studies have been conducted with a high degree of success on the application of LLMs for the translation of biomedical texts. The first study of this sort was published by Han et al. (2022), where they compare MT models of different sizes to investigate the applicability of Kaplan's scaling laws (Kaplan et al., 2020) in biomedical translation. Their findings confirmed that larger general-purpose models consistently outperform smaller models, even when the latter are fine-tuned on domain-specific data. Interestingly, the performance gap narrows significantly when the training data for smaller models is meticulously curated, bringing their efficacy close to that of the NLLB model. The efficacy of LLMs in translating biomedical data is further confirmed by several recent studies (Jahan et al., 2024; Keles et al., 2024;

García-Ferrero et al., 2024).

Finally, we underline one of the latest research directions in the study of LLMs: multi-agentic workflows. Agents are instances of LLMs, each with a tailored system prompt that defines their behaviour, adhering to specific criteria and output requirements. Usually, agents also have access to external features, such as memory mechanism (Zhang et al., 2024), retrieval-augmented generation (Gao et al., 2024), and tool use (Qu et al., 2024). In the experiment conducted by Liang et al. (2024), the authors exemplify this approach with a novel translation framework called Multi-Agent Debate (MAD). Their system is based on a guided interaction between multiple agents who engage in a debate to determine the most effective translation for a given source text. A designated judge agent oversees this process and ultimately decides on the final solution. This iterative strategy allows successive agents to refine the initial translation hypothesis, progressively improving translation quality. They achieve good performance with the models gpt-3.5-turbo and gpt-4 (Brown et al., 2020).

## 3 Systems Overview

We submit five MT systems for evaluation, each employing different approaches to biomedical translation. These systems range from traditional fine-tuning on in-domain data to various data augmentation approaches and the use of LLMs, prompt engineering, and multi-agent workflows.

Table 1 provides an overview of the five systems submitted for evaluation. System 1 utilizes the NLLB model fine-tuned with terminology mining techniques, applied in both directions (see §5.1). System 2 also uses NLLB, but we fine-tune it on both in-domain and synthetic data. For this system, we augmented the training data with an additional 5,000 backtranslated sentences to address data scarcity (see §5.2). System 3 uses a combination of agents powered by NLLB and GPT, who are tasked with post-editing and refining the NLLB outputs to make them more fluent and effective. For System 4, we select the smallest checkpoint of the most recent models developed by Meta AI, called LLama-3-8B (Dubey et al., 2024). This system uses parameter-efficient fine-tuning on the in-domain data, and the output is improved with three fuzzy matches prepended to the prompt. Finally, our last submission, System 5 uses a multi-agent crew powered by GPT4-o (OpenAI et al., 2024).

| System | Model | In-Domain FT | Backtranslation | Terminology Mining | Agents | FSP |
|---|---|---|---|---|---|---|
| 1 | NLLB | ✓ | | ✓ | | |
| 2 | NLLB | ✓ | ✓ | | | |
| 3 | NLLB | ✓ | | | ✓ | |
| 4 | LLama-3 | ✓ | | | | ✓ |
| 5 | GPT-4o | | | | ✓ | ✓ |

Table 1: Overview of our submitted systems. The checkmark (✓) indicates the presence of a feature. FSP stands for Few-Shot Prompting.

The multi-agents workflow is described in depth in the relevant section (see §6).

## 4 Dataset Selection

In this section, we describe the composition of the datasets used for our experiments. We curated a selection of parallel sentences from the corpora provided by the shared task organisers, including part of the Biomedical Translation repository and the UFAL Medical[1] corpus. This resulted in two datasets: 11,190 parallel sentences for English-German and 13,032 for English-French. We investigated synthetically increasing the training data by employing different data augmentation techniques for the English-to-German language pair. We provide an overview of the dataset selection in Table 2.

| Dataset | EN-DE | EN-FR |
|---|---|---|
| Original | 11,190 | 13,032 |
| + Term. Mining | 14,583 | NA |
| + Backtranslation | 16,190 | NA |

Table 2: Overview of datasets.

## 5 Data Augmentation

In this section, we describe the different approaches we have used to augment the datasets used for our MT systems. We adopt back-translation, terminology mining, and fuzzy matches.

### 5.1 Terminology Mining

We perform terminology mining on English-to-German language pairs. We extract biomedical terms from the training data using the pre-trained named entity recognition (NER) model d4data/biomedical-ner-all. This model is designed to identify biomedical entities within the text, such as diseases, disorders, and therapeutic

procedures, providing a confidence score and the specific unit being identified. The implementation utilises the pipeline function from the Hugging Face Transformers[2] library (Wolf et al., 2020), configured for the task of token classification.

First, the NER model iterates over every term in the dataset, obtaining a list of identified entities. We then filter them, collecting only those labeled as *B-Disease-disorder* or *B-Therapeutic-procedure*, provided that the model's confidence score for the entity exceeds 0.98 and the length of the identified word is greater than five characters. Entities meeting these criteria are then printed for verification and appended to the list of extracted terms. This approach ensures that only relevant biomedical terms are extracted from the dataset, focusing specifically on diseases, disorders, and therapeutic procedures. The terminology mining process yielded a total of 14 biomedical terms, that we used to collect 3,393 sentences containing at least one biomedical term.

### 5.2 Backtranslation

We adopt back-translation for the English-German language pair, to address data scarcity with semantically similar sentences, extracted with a pre-trained sentence embedding model, and then backtranslated with NLLB. We initially filter the EMEA monolingual dataset (Calzolari et al., 2012), selecting only sentences that exceed 100 characters in length, and further limit our selection to the first 1,000 entries for easier processing. We encode the text data using the pre-trained sentence embedding model multi-qa-mpnet-base-dot-v1 from the Sentence-Transformers library. The sentence embeddings are stored as a new column in the dataset, on which we perform semantic search, using FAISS index (Douze et al., 2024) for more efficient computation. The index is then queried to retrieve the top 5 most similar samples from the original dataset. We collect a total of 5,000 sentences, aggregated

and sorted by similarity scores in descending order.

This methodology enables the efficient retrieval of contextually relevant sentences from large datasets. We make use of the resulting sentences, backtranslating them to English using the baseline `NLLB-200-600M` and leading to the creation of a synthetic dataset, that is in the same domain. The synthetic dataset is added to the original dataset to fine-tune the baseline model.

## 5.3 Fuzzy Matches

Fuzzy matches are human translated segments, stored in parallel datasets. Drawing on findings from Moslem et al. (2023), we incorporate semantically-similar fuzzy matches in a three-shot prompting scenario. This approach leverages the model's in-context learning ability (Brown et al., 2020) to further improve the quality of the MT outputs. A wide range of academic literature has demonstrated that incorporating fuzzy matches in a few-shot scenario may improve the model's understanding of domain-specific terminology and fixed expressions (Castaldo and Monti, 2024; Moslem et al., 2022; Knowles et al., 2018).

To extract fuzzy matches, we employ semantic search on sentence embeddings generated by the `all-MiniLM-L6-v2` model. The embeddings are stored in a flat index created with the FAISS[3] library, from which we retrieve the three most similar sentences. After extracting the fuzzy matches for our input sentence, we prepend them to a minimalist prompt that directly maps the source language to the target language. We incorporate fuzzy matches in System 4 and 5, achieving substantial improvements over the zero-shot baseline. We present an overview of the prompt templates used in this study in Table 3, with the following annotations: ♦ shows the presence of a line break, [src] stands for source language, [tgt] stands for target language, and [input] stands for the text to be translated.

| Prompt Type | Template |
|---|---|
| Zero-Shot | [src]: [input] ♦ [tgt]: |
| Few-Shots | [src]: [source$_1$] ♦ [tgt]: [target$_1$] ♦ ... [src]: [source$_k$] ♦ [tgt]: [target$_k$] ♦ [src]: [input] ♦ [tgt]: |

Table 3: Overview of the prompt templates used in this study.

## 6 Multi-Agents Workflow

We design a team composed of three autonomous agents that collaborate to simulate a translation agency with the goal of refining an initial translation hypothesis from multiple perspectives. The process begins with the creation of our agent crew, using the CrewAI library[4].

The first agent, the Translator Agent, is tasked with translating a given sentence. Following this, the Evaluator Agent assesses the translation based on fluency and accuracy. This assessment is quantified with a numerical quality metric that ranges from 0 to 100, where 100 signifies a translation that is both perfectly fluent and accurate.

If the translation receives a score below 80, the Reviewer Agent intervenes to review the initial hypothesis, aiming to improve its accuracy. This iterative process repeats until the Evaluator Agent awards a quality score greater than 80, indicating a successful translation. We provide the reference code used for this experiment in the relevant GitHub repository.[5]

## 7 Evaluation

This section discusses the results that we obtained from our experiments. Table 4 shows the results obtained by evaluating our models on the validation set, using BLEU (Papineni et al., 2002), ChrF (Popović, 2015) and the COMET model `wmt22-comet-da` (Rei et al., 2022a). Our quality estimation is based on the reference-free COMET model `wmt22-cometkiwi-da` (Rei et al., 2022b). Our systems achieve good results for both language pairs, and the data augmentation approaches visibly improve the translation outputs, as documented in the evaluation of Systems 1, 2 and 3. Terminology mining seems particularly effective, improving the BLEU score of our first system significantly above the others.

In order to confirm the results of our automatic evaluation and to allow for a more precise comparison of the different systems used for the primary language pair of our study, we include a manual evaluation on a small sample of translations, conducted by two professional translators in the EN-DE language pair, for which we adopt the MQM-DQF framework (Burchardt, 2013; Lommel and

---

Melby, 2018).

The results of the MQM evaluation reveal that System 5 produces the fewest overall errors, with the majority of these errors falling under the Fluency category. In contrast, the other systems exhibit a higher concentration of errors in the Accuracy category. In System 1, where Terminology Mining was applied, fewer terminology-related errors were detected, further confirming the effectiveness of this strategy. Additionally, we find that in System 3 the involvement of `GPT-4o` agents in post-editing led to a reduction in accuracy-related errors.

| System | BLEU | ChrF | COMET | QE |
|---|---|---|---|---|
| **English-to-German** | | | | |
| Baseline | 2.97 | 26.01 | 70.09 | 0.62 |
| System 1 | 25.29 | 60.13 | 79.50 | 0.58 |
| System 2 | 23.80 | 58.89 | 78.33 | 0.56 |
| System 3 | 23.97 | 59.22 | 78.93 | 0.63 |
| System 4 | 22.95 | 58.90 | 84.32 | 0.73 |
| System 5 | 25.24 | 63.01 | 86.13 | 0.77 |
| **English-to-French** | | | | |
| System 3 | 29.18 | 57.01 | 75.70 | 0.65 |

Table 4: Experiment Results for Different Systems

## 8 Conclusions

This study presents the approaches we have adopted to address the challenges caused by biomedical translation, specifically the need for consistent translation of domain-specific terminology and the lack of in-domain parallel data.

By adopting data augmentation techniques, we found that our models improved consistently in translating biomedical terminology, achieving better results in our evaluation. Terminology mining proved particularly effective, resulting in our best overall submission. We also explored the use of backtranslation, but we found that its effectiveness may be limited in fine-tuning LLMs. We speculate that it may require a different ratio of original to synthetic data used during training, or a different weighting. Our experiments with fuzzy matches demonstrated the potential to use in-context learning to improve MT quality and adapt LLMs to domain-specific terminology.

Finally, we introduced a novel MT workflow based on the collaboration of three autonomous LLM-based agents. This approach offers an innovative way to refine an initial translation hypothesis from multiple perspectives, potentially leading to more accurate outputs.

## 9 Limitations

We acknowledge that several aspects of our study have room for improvement. First, the evaluation was conducted on a relatively small dataset of 50 biomedical abstracts, limiting the objectivity of the results. Second, while data augmentation helped improve performance, the training data could be expanded by incorporating larger corpora and potentially leading to better quality. Additionally, the models employed in this study may not represent the best performing MT systems by the time of publication, requiring further experiments with more recent models to validate our findings. Finally, manual evaluation was only conducted for a single language pair, limiting the scope of our analysis.

## 10 Acknowledgements

## References

Ballier, N., Yunès, J.-b., Wisniewski, G., Zhu, L. (2022). The SPECTRANS System Description for the WMT22 Biomedical TaskIn Koehn, P., Barrault, L., Bojar, O., Bougares, F., Chatterjee, R., Costa-jussà, M. R., Federmann, C., Fishel, M., Fraser, A., Freitag, M. et al., editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 895–900, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Bojar, O., Haddow, B., Marecek, D., Sudarikov, R., Tamchyna, A. (2017). Report on building translation systems for public health domain. UFAL Medical Corpus.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A. et al. (2020). Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, pages 1877–1901, Red Hook, NY, USA. Curran Associates Inc.

Burchardt, A. (2013). Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.

Calzolari, N., Choukri, K., Declerck, T., Doğan, M. U., Maegaard, B., Mariani, J., Moreno, A., Odijk, J. (2012). Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12).

Castaldo, A. (2024). Prompting Large Language Models for Idiomatic Translation. In *Proceedings of the First Workshop on Creative-text Translation and Technology*, pages 37–44, Sheffield, UK. Accepted: 2024-06-19T21:00:05Z.

Choi, Y., Shin, J., Ryu, Y. (2022). SRT's Neural Machine Translation System for WMT22 Biomedical Translation TaskIn Koehn, P., Barrault, L., Bojar, O., Bougares, F., Chatterjee, R., Costa-jussà, M. R., Federmann, C., Fishel, M., Fraser, A., Freitag, M. et al., editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 901–907, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Costa-jussà, M. R., Cross, J., Çelebi, O., Elbayad, M., Heafield, K., Heffernan, K., Kalbassi, E., Lam, J., Licht, D., Maillard, J. et al. (2022). No Language Left Behind: Scaling Human-Centered Machine Translation. arXiv:2207.04672 [cs].

Douze, M., Guzhva, A., Deng, C., Johnson, J., Szilvasy, G., Mazaré, P.-E., Lomeli, M., Hosseini, L. (2024). The Faiss library. arXiv:2401.08281 [cs].

Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A. et al. (2024). The Llama 3 Herd of Models. arXiv:2407.21783 [cs].

Dušek, O., Hajič, J., Hlaváčová, J., Libovický, J., Pecina, P., Tamchyna, A. (2017). Khresmoi summary translation test data 2.0. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Wang, M. (2024). Retrieval-Augmented Generation for Large Language Models: A Survey. arXiv:2312.10997 [cs].

García-Ferrero, I., Agerri, R., Atutxa Salazar, A., Cabrio, E., de la Iglesia, I., Lavelli, A., Magnini, B., Molinet, B., Ramirez-Romero, J., Rigau, G. et al. (2024). MedMT5: An Open-Source Multilingual Text-to-Text LLM for the Medical DomainIn Calzolari, N., Kan, M.-Y., Hoste, V., Lenci, A., Sakti, S., editors, *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 11165–11177, Torino, Italia. ELRA and ICCL.

Han, L., Erofeev, G., Sorokina, I., Gladkoff, S. (2022). Examining Large Pre-Trained Language Models for Machine Translation: What You Don't Know about ItIn Koehn, P., Barrault, L., Bojar, O., Bougares, F., Chatterjee, R., Costa-jussà, M. R., Federmann,

C., Fishel, M., Fraser, A., Freitag, M. et al., editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 908–919, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Haque, R., Moslem, Y. (2020). Terminology-Aware Sentence Mining for NMT Domain Adaptation: ADAPT's Submission to the Adap-MT 2020 English-to-Hindi AI Translation Shared TaskIn Sharma, D. M., Ekbal, A., Arora, K., Naskar, S. K., Ganguly, D, L, S., Mamidi, R., Arora, S., Mishra, P., editors, *Proceedings of the 17th International Conference on Natural Language Processing (ICON): Adap-MT 2020 Shared Task*, pages 17–23, Patna, India. NLP Association of India (NLPAI).

He, Z., Liang, T., Jiao, W., Zhang, Z., Yang, Y., Wang, R., Tu, Z., Shi, S. (2024). Exploring Human-Like Translation Strategy with Large Language Models. *Transactions of the Association for Computational Linguistics*, 12:229–246.

Hendy, A., Abdelrehim, M., Sharaf, A., Raunak, V., Gabr, M., Matsushita, H., Kim, Y. J., Afify, M. (2023). How Good Are GPT Models at Machine Translation? A Comprehensive Evaluation. arXiv: 2302.09210.

Jahan, I., Laskar, M. T. R., Peng, C. (2024). A comprehensive evaluation of large Language models on benchmark biomedical text processing tasks. *Computers in Biology and Medicine*, 171:108189.

Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J. (2020). Scaling Laws for Neural Language Models. arXiv:2001.08361 [cs, stat].

Keles, B., Gunay, M. (2024). LLMs-in-the-loop Part-1: Expert Small AI Models for Bio-Medical Text Translation. arXiv:2407.12126 [cs].

Knowles, R., Ortega, J. (2018). A Comparison of Machine Translation Paradigms for Use in Black-Box Fuzzy-Match RepairIn Astudillo, R., Graça, J., editors, *Proceedings of the AMTA 2018 Workshop on Translation Quality Estimation and Automatic Post-Editing*, pages 249–255, Boston, MA. Association for Machine Translation in the Americas.

Liang, T., He, Z., Jiao, W., Wang, X., Wang, R., Yang, Y., Tu, Z. (2024). Encouraging Divergent Thinking in Large Language Models through Multi-Agent Debate. arXiv:2305.19118 [cs].

Lommel, A. (2018). Tutorial: MQM-DQF: A Good Marriage (Translation Quality for the 21st Century)In Campbell, J., Yanishevsky, A., Doyon, J., editors, *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 2: User Track)*, Boston, MA. Association for Machine Translation in the Americas.

Manchanda, S. (2022). Optum's Submission to WMT22 Biomedical Translation TasksIn Koehn, P., Barrault, L., Bojar, O., Bougares, F., Chatterjee, R.,

652

Costa-jussà, M. R., Federmann, C., Fishel, M., Fraser, A., Freitag, M. et al., editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 925–929, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Moslem, Y., Haque, R., Kelleher, J. (2022). Domain-Specific Text Generation for Machine TranslationIn Duh, K., editors, *Proceedings of the 15th biennial conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 14–30, Orlando, USA. Association for Machine Translation in the Americas.

Moslem, Y., Haque, R., Kelleher, J. D. (2023). Adaptive Machine Translation with Large Language ModelsIn Nurminen, M., Brenner, J., Koponen, M., Latomaa, S., Mikhailov, M., Schierl, F., Ranasinghe, T., Vanmassenhove, E., Vidal, S. A., Aranberri, N. et al., editors, *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 227–237, Tampere, Finland. European Association for Machine Translation.

OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S. et al. (2024). GPT-4 Technical Report. arXiv:2303.08774 [cs].

Papineni, K., Roukos, S., Ward, T. (2002). Bleu: a Method for Automatic Evaluation of Machine TranslationIn Isabelle, P., Charniak, E., editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Popović, M. (2015). chrF: character n-gram F-score for automatic MT evaluationIn Bojar, O., Chatterjee, R., Federmann, C., Haddow, B., Hokamp, C., Huck, M., Logacheva, V., editors, *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Przystupa, M. (2019). Neural Machine Translation of Low-Resource and Similar Languages with Back-translationIn Bojar, O., Chatterjee, R., Federmann, C., Fishel, M., Graham, Y., Haddow, B., Huck, M., Yepes, A. J., Koehn, P., Martins, A. et al., editors, *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 224–235, Florence, Italy. Association for Computational Linguistics.

Qu, C., Dai, S., Wei, X., Cai, H., Wang, S., Yin, D., Xu, J. (2024). Tool Learning with Large Language Models: A Survey. arXiv:2405.17935 [cs].

Rei, R., C. de Souza, J. G., Alves, D., Zerva, C., Farinha, A. C., Glushkova, T., Lavie, A., Coheur, L. (2022a). COMET-22: Unbabel-IST 2022 Submission for the Metrics Shared TaskIn Koehn, P., Barrault, L., Bojar, O., Bougares, F., Chatterjee, R., Costa-jussà, M. R., Federmann, C., Fishel, M., Fraser, A., Freitag, M.

et al., editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Rei, R., Treviso, M., Guerreiro, N. M., Zerva, C., Farinha, A. C., Maroti, C., C. de Souza, J. G., Glushkova, T., Alves, D., Coheur, L. et al. (2022b). CometKiwi: IST-Unbabel 2022 Submission for the Quality Estimation Shared TaskIn Koehn, P., Barrault, L., Bojar, O., Bougares, F., Chatterjee, R., Costa-jussà, M. R., Federmann, C., Fishel, M., Fraser, A., Freitag, M. et al., editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Tang, Y., Tran, C., Li, X., Chen, P.-J., Goyal, N., Chaudhary, V., Gu, J. (2021). Multilingual Translation from Denoising Pre-TrainingIn Zong, C., Xia, F., Li, W., editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3450–3466, Online. Association for Computational Linguistics.

Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y. et al. (2024). A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M. et al. (2020). HuggingFace's Transformers: State-of-the-art Natural Language Processing. arXiv:1910.03771 [cs].

Xu, N., Li, Y., Xu, C., Li, Y., Li, B., Xiao, T. (2019). Analysis of Back-Translation Methods for Low-Resource Neural Machine TranslationIn Tang, J., Kan, M.-Y., Zhao, D., Li, S., editors, *Natural Language Processing and Chinese Computing*, pages 466–475, Cham. Springer International Publishing.

Zhang, Z., Bo, X., Ma, C., Li, R., Chen, X., Dai, Q., Zhu, J., Dong, Z. (2024). A Survey on the Memory Mechanism of Large Language Model based Agents. arXiv:2404.13501 [cs].

# Findings of WMT 2024 Shared Task on Low-Resource Indic Languages Translation

**Partha Pakray**
NIT Silchar

**Santanu Pal**
Wipro AI Lab45

**Advaitha Vetagiri**
NIT Silchar

**Reddi Mohana Krishna**   **Arnab Kumar Maji**   **Sandeep Kumar Dash**   **Lenin Laitonjam**
NIT Silchar         North-Eastern Hill University      NIT Mizoram          NIT Mizoram

**Lyngdoh Sarah**
North-Eastern Hill University

**Riyanka Manna**
Amrita Vishwa Vidyapeetham Amaravati

## Abstract

This paper presents the results of the low-resource Indic language translation task, organized in conjunction with the Ninth Conference on Machine Translation (WMT) 2024. In this edition, participants were challenged to develop machine translation models for four distinct language pairs: English–Assamese, English-Mizo, English-Khasi, and English-Manipuri. The task utilized the enriched IndicNE-Corp1.0 dataset, which includes an extensive collection of parallel and monolingual corpora for northeastern Indic languages. The evaluation was conducted through a comprehensive suite of automatic metrics—BLEU, TER, RIBES, METEOR, and ChrF—supplemented by meticulous human assessment to measure the translation systems' performance and accuracy. This initiative aims to drive advancements in low-resource machine translation and make a substantial contribution to the growing body of knowledge in this dynamic field.

## 1 Introduction

The low-resource Indic language translation field has witnessed significant advancements, particularly marked by the success of last year's Indic MT Shared Task. This initiative, organized alongside the Eighth Conference on Machine Translation (WMT) 2023[1] (Pal et al., 2023), demonstrated the potential and necessity of focusing on low-resourced languages. Building on the momentum and achievements of last year's task, we are pleased to continue our efforts with the Indic MT Shared Task for the Ninth Conference on Machine Translation (WMT) 2024[2].

Low-resource Indic languages represent a vast and diverse array of languages spoken across India. Despite their deep cultural and linguistic heritage,

these languages face significant challenges due to limited resources and institutional support. The obstacles are multifaceted, including smaller speaker populations, minimal governmental backing, insufficient documentation, and restricted access to modern technological tools.

India is celebrated for its linguistic diversity, with many languages spoken throughout the subcontinent. The Eighth Schedule of the Indian Constitution officially recognizes 22 languages, granting them substantial governmental support and resources. However, numerous other languages, particularly those spoken by indigenous and minority communities, often remain marginalized and under-supported. These low-resource languages encounter additional barriers, such as the absence of standardized scripts, limited lexical resources, and a dearth of linguistic research. These factors, combined with the lack of formal educational resources and declining inter-generational transmission, threaten their preservation and vitality. As a result, many of these languages risk becoming endangered, underscoring the urgent need for targeted efforts to document, revitalize, and sustain them in the face of ongoing challenges.

Given these challenges, our initiative is dedicated to documenting, revitalizing, and supporting low-resource Indic languages through innovative technological solutions. The previous year's Indic MT Shared Task concentrated on four language pairs: English–Assamese, English–Mizo, English–Khasi, and English–Manipuri — utilizing the enriched IndicNE-Corp1.0 dataset (Pal et al., 2023). The success of this task highlighted the critical need for sustained efforts in this domain. Our ongoing objective is to foster advancements in machine translation and natural language processing tailored to these languages.

The evaluation of this task employs a comprehensive set of metrics, incorporating both automatic measures—such as BLEU (Papineni et al., 2002),

---

[1] https://www2.statmt.org/wmt23/indic-mt-task.html
[2] https://www2.statmt.org/wmt24/indic-mt-task.html

TER (Snover et al., 2006), RIBES (Isozaki et al., 2010), METEOR (Banerjee and Lavie, 2005), and ChrF (Popović, 2015)—and rigorous human assessments. This dual approach ensures a thorough evaluation of the translation systems' performance, accuracy, and cultural fidelity.

Through this ongoing initiative, we aim to make a significant contribution to the preservation of linguistic diversity and cultural heritage, thereby supporting the rights and identities of minority language communities in India. By leveraging cutting-edge technologies, we strive to create a lasting impact and propel the field of low-resource language translation forward, ensuring these languages not only survive but thrive in the digital age.

## 2 Languages

### 2.1 Khasi Language and Its Dialects

Khasi belongs to the Austro-Asiatic family of languages spoken in the central and eastern regions of Meghalaya. Before 1813, the Khasi lacked its own script. During the period from 1813 to 1814, the Bengali script was employed to translate the Bible into Khasi, owing to the widespread literacy in Bengali at that time. By 1816, some translated versions of the Gospel of Matthew had been printed and distributed among Khasi speakers who were literate in Bengali. However, it was not until 1841, with the arrival of a Welsh missionary, that the Roman script was introduced, and translations were subsequently made into the standard dialect, specifically the Sohra variety.

Khasi exhibits significant dialectal diversity. Grierson (1904) identified four dialects of Khasi: Standard Khasi, Pnar or Synteng, Lyngngam, and War. Acharya (1971) reaffirmed Grierson's classification and noted the existence of additional sub-dialects, such as Bhoi, spoken in the northern open lands of Meghalaya. Bareh (1977) offers a more comprehensive list of Khasi dialects, primarily based on their geographical distribution:

- Amwi in the southern Jaiñtia hills,

- Shella in the southern Khasi hills,

- Warding in the south of the Khasi hills,

- Myriaw, Nongkhlaw, Nongspung, Maram, and Mawiang in the mid-western area of the Khasi hills,

- Cherra in the mid-southern hills,

- Mylliem, Laitlyngkot, Nongkrem, and Lyniong-Khasi in the central parts,

- Jowai in the central Jaintia hills,

- Bhoi in the north-east Khasi hills,

- Manar, Nongwah, and Jirang in the north Khasi hills,

- Khatarblang (Mawpran) in the mid-southern region, and

- Nongstoin and Langrin in the west Khasi region.

Bareh further adds that several sub-dialects exhibit variations within each group, particularly in phonology. Among these, Amwi is considered the most typical dialect. Compared to other dialects, Amwi appears to be the most rudimentary and is generally not intelligible to speakers of neighbouring dialects such as Jowai or Khad ar Blang. Amwi is said to be more agglutinative in form, potentially preserving its Mon-Khmer heritage. While its grammar resembles Jowai's, notable differences exist in morphology and phonology. Despite these distinctions, the Amwi speakers are familiar with their neighbouring dialects and can adopt them for communication.

Bareh (1977) categorizes the aforementioned dialects into three major branches:

1. **Eastern dialects**:

   - Jowai (Central Highlands),
   - Amwi and the War dialects (in the south), and
   - Bhoi Synteng in the north.

2. **Central dialects**:

   - Nongphlang or Nonglum, Cherra, and related dialects such as Nongkrem, Mylliem, Nongkhlaw, Nongspung, Rambrai, Mawsynram, Maram, Laitlyngkot, Mawphlang, etc.,
   - Bhoi East (in the north), consisting of Mawrong, Bhoi Lymbong, etc., and
   - Bhoi West (in the north), consisting of Manar, etc.,
   - War Shala (in the south), and
   - Warding (in the south).

3. **Western dialects**:

- Nongstoin
- Lyngam
- Langrin

Addressing the abovementioned dialects, Bareh notes that numerous sub-dialects exhibit phonological variations within each group. Daladier (2007:341), cited in Sidwell (2009), comments on the Mon-Khmer language group, which includes Khasi, noting that it comprises three main branches. Although now standardized and formalized through written use, Khasi retains conservative unwritten dialects, particularly in the War region. Other notable dialects include Pnar and War, with War further subdivided into four sub-dialect groups: Nongtalang, Amvi, Tremblang, and Shella. The sub-classification of Pnar dialects remains largely unexplored. Additionally, Pnar-War and War-Khasi dialects are spoken in several Jaintia villages. The War dialects of Khasi are divided into two groups: War-Khasi and War-Jaintia, spoken in the southeast corners of the Khasi and Jaintia Hills districts, respectively. Grierson (1904) also discusses the War dialects.

For the shared task, we have utilized the Sohra (Cherra) dialect of Khasi as the standard form for translation purposes. This dialect, recognized for its historical significance and broad usage in educational and religious contexts, has been established as the standardized variant of Khasi following its formalization through the introduction of the Roman script in 1841. By employing the Sohra dialect, we ensure consistency and accessibility for participants, reflecting the widely accepted linguistic norm within the Khasi-speaking community.

## 2.2 Introduction: About the Manipuri Language

Manipuri, also known as Meiteilon, is a Sino-Tibetan language predominantly spoken in the northeastern Indian state of Manipur. It is recognized as one of the 22 scheduled languages of India and serves as the lingua franca among various ethnic communities in the region, fostering communication and cultural exchange.

The language boasts a rich literary heritage, with a history of written texts dating back to ancient times. Manipuri uses the Meitei script, also known as Meitei Mayek, alongside the Bengali script for writing purposes. Despite its cultural significance, Manipuri faces linguistic preservation and modernization challenges, particularly in the digital era. There is a pressing need for computational resources and tools to support the language, which is vital for its continued use and growth.

In recent years, there has been growing interest in developing natural language processing (NLP)(Allen, 2003) tools and resources for under-resourced languages like Manipuri. However, several challenges persist in this area for the Manipuri language (Gyanendro Singh et al., 2016). One of the primary issues is the limited availability of annotated corpora and linguistic resources, which are essential for training robust machine learning models. This scarcity hinders the development of accurate NLP applications such as machine translation (Pal et al., 2023), sentiment analysis (Singh and Singh, 2017), and speech recognition (Gyanendro Singh et al., 2016).

Another significant challenge is the complexity of the Manipuri script and its morphological structure. The language exhibits rich inflectional morphology, making it difficult to apply standard NLP techniques that are typically designed for resource-rich languages like English. Moreover, the lack of standardization in digital representation further complicates computational processing, as existing tools often struggle with script conversion and text normalization.

Current research efforts are focused on addressing these challenges by creating linguistic resources, developing language-specific algorithms, and adapting existing NLP frameworks to better accommodate the unique characteristics of Manipuri. However, much work remains to be done to bridge the gap between Manipuri and other well-resourced languages in the digital domain.

## 2.3 Introduction: About the Assamese Language

Assamese, an Indo-Aryan language, is predominantly spoken in the northeastern Indian state of Assam. It serves as the official language of Assam and plays a crucial role as a lingua franca among various ethnic groups in the region, facilitating communication and cultural exchange. Assamese is also one of the 22 scheduled languages of India, underscoring its significance in the country's linguistic landscape.

The Assamese language has a rich literary tradition, with its roots extending back to the early medieval period. The script used for Assamese is derived from the ancient Brahmi script, and over

time, it evolved into its current form. Despite its historical and cultural importance, Assamese faces challenges in the modern era, particularly in the field of language technology. The development of computational tools and resources for Assamese is critical for its preservation and growth, especially in an increasingly digital world.

## 2.4 Introduction: About the Mizo Language

Mizo, a member of the Tibeto-Burman language family, is predominantly spoken in the northeastern Indian state of Mizoram. It serves as the primary language of communication among the Mizo people and is also spoken by various ethnic groups in neighbouring states and regions, including Manipur, Tripura, Assam, and even parts of Myanmar and Bangladesh. Mizo is recognized for its tonal nature and distinct phonological features, which make it a unique language in the Tibeto-Burman group.

The language has a rich oral tradition, encompassing folktales, songs, and cultural narratives that reflect the heritage of the Mizo people. The development of the written form of Mizo began in the late 19th century with the introduction of the Roman script by Christian missionaries, which facilitated the transcription of the language and the creation of written literature. Today, Mizo has a well-established literary tradition, with a substantial body of work ranging from poetry to modern prose. Despite its cultural significance, Mizo faces linguistic preservation and development challenges, particularly in the context of modern technological advancements and digital communication.

## 3 Low-Resource Indic Language Translation 2024 Shared Task

### 3.1 Overview and Task Description

Building upon the resounding success of the "Shared Task: Low-Resource Indic Language Translation" at WMT 2023, which witnessed enthusiastic participation from around the globe, we are excited to announce the continuation of this initiative at the Ninth Conference on Machine Translation (WMT 2024). The advances in machine translation (MT) have significantly enhanced the performance of translation systems, especially with the adoption of techniques such as multilingual translation and transfer learning. Despite these advancements, extending coverage to diverse low-resource languages remains a formidable challenge

due to the scarcity of parallel data needed to train robust MT systems.

The WMT 2024 Indic Machine Translation Shared Task addresses this challenge by focusing on low-resource Indic languages from diverse language families. This year, the task emphasizes the following language pairs: English-Assamese, English-Mizo, English-Khasi, and English-Manipuri. Additionally, there was an intended focus on English-Nyishi; however, this category was cancelled due to issues with training data. Similarly, other planned language pairs under the category with very limited training data, such as English-Bodo, English-Mising, and English-Kokborok, were also cancelled for this year.

### 3.2 Categories

This year's task features two main categories based on the availability of training data:

#### 3.2.1 Category 1: Moderate Training Data Available

- English ⇔ Assamese (en-as)

- English ⇔ Mizo (en-lus)

- English ⇔ Khasi (en-kha)

- English ⇔ Manipuri (en-mni)

### 3.3 Goal

The central objective of this shared task is to develop machine translation systems that produce high-quality translations despite the constraints posed by limited data availability. Participants are encouraged to explore several innovative strategies, including:

- **Monolingual Data Utilization**: Effectively leveraging monolingual data to enhance translation quality.

- **Multilingual Approaches**: Investigating the benefits of cross-lingual transfer for low-resource language pairs.

- **Transfer Learning**: Adapting models trained on resource-rich language pairs to target low-resource languages.

- **Innovative Techniques**: Experimenting with novel methods specifically tailored to low-resource settings.

### 3.4 Data

#### 3.4.1 Training

The datasets used for this task include parallel and monolingual corpora for Assamese, Khasi, Mizo, and Manipuri, drawn from the IndicNE-Corp1.0 dataset. While the dataset for Nyishi was planned, it remains unavailable this year due to data quality issues.

#### 3.4.2 Testing

For the testing section, we have created 1000 language pair sentences for each of the following language pairs:

- English ⇔ Assamese (en-as)

- English ⇔ Mizo (en-lus)

- English ⇔ Khasi (en-kha)

- English ⇔ Manipuri (en-mni)

The first 500 sentences are provided in English to be translated into the specific target language, and the last 500 sentences are provided in the target language to be translated into English.

### 3.5 Evaluation

The evaluation will be conducted using both automatic and human evaluation methods to ensure a comprehensive assessment of the translation systems. Automatic evaluation metrics include BLEU, TER, RIBES, METEOR, and ChrF. In addition, native speakers will perform human evaluations to assess the quality of the translation more rigorously.

## 4 Dataset

### 4.1 Training

The dataset for the WMT 2024 Shared Task on Low-Resource Indic Language Translation is primarily based on the IndicNE-Corp1.0 dataset [3]. This corpus was built by aggregating datasets from previous research, including significant contributions from (Laskar et al., 2020) (Laskar et al., 2022), (Khenglawt et al., 2022), and (Laitonjam and Ranbir Singh, 2021). The compiled datasets encompass both parallel and monolingual corpora across four languages: Assamese, Mizo, Khasi, and Manipuri.

---

[3] https://data.statmt.org/wmt23/indic-mt/

In earlier studies, we focused on developing parallel and monolingual corpora for English ⇔ Assamese (en-asm) (Laskar et al., 2020, 2022), English ⇔ Mizo (en-lus) (Khenglawt et al., 2022), English ⇔ Khasi (en-kha) (Laskar et al., 2021), and English ⇔ Manipuri (en-mni) (Laitonjam and Ranbir Singh, 2021). The data was sourced from a variety of online platforms including the Bible, multilingual dictionaries (such as Xobdo and Glosbe), multilingual question papers, PMIndia (Haddow and Kirefu, 2020), web pages, blogs, and online newspapers.

Table 1 shows the detailed statistics of the parallel datasets used for training and validation for each language pair.

| Type | Sentences | Tokens (eng) | Tokens (target) |
|------|-----------|--------------|-----------------|
| Assamese | 50,000 | 969,623 | 825,063 |
| Mizo | 50,000 | 981,468 | 1,062,414 |
| Khasi | 24,000 | 729,930 | 875,545 |
| Manipuri | 21,687 | 390,730 | 330,319 |

Table 1: Parallel data statistics for train and validation.

In addition to the parallel corpora, we also made monolingual data available for each language, which is presented in Table 2.

| Language | Size (MB) | Sentences | Tokens |
|----------|-----------|-----------|--------|
| Assamese | 805 | 2,624,715 | 49,232,154 |
| Mizo | 145 | 1,909,823 | 27,936,225 |
| Khasi | 104 | 182,737 | 22,140,361 |
| Manipuri | 716 | 2,144,897 | 36,514,693 |

Table 2: Monolingual data statistics for Assamese, Mizo, Khasi, and Manipuri languages.

### 4.2 Testing

The testing dataset for the 2024 shared task was meticulously curated to present a substantial challenge beyond previous years' datasets. It comprised 1000 samples for each language pair, spanning four distinct and diverse domains: News, Travel, Sports, Entertainment, and Business. This domain-specific distribution aimed to comprehensively evaluate models' performance across varied and complex linguistic contexts, reflecting real-world translation demands. A collaborative approach was employed to create these testing samples, involving four specialized teams, each dedicated to one domain. These teams were provided 1000 English sentences, which they translated into their assigned target languages. The translation teams were instructed to maintain high fidelity to the source mate-

| Language Pair | Domain | Source Sentences | Target Sentences | Task |
|---|---|---|---|---|
| en-as | Sports and Travel | 500 | 500 | English to Assamese |
| en-lus | Sports and Travel | 500 | 500 | English to Mizo |
| en-kha | Sports and Travel | 500 | 500 | English to Khasi |
| en-mni | Sports and Travel | 500 | 500 | English to Manipuri |
| en-as | Entertainment and Business | 500 | 500 | Assamese to English |
| en-lus | Entertainment and Business | 500 | 500 | Mizo to English |
| en-kha | Entertainment and Business | 500 | 500 | Khasi to English |
| en-mni | Entertainment and Business | 500 | 500 | Manipuri to English |

Table 3: Domain-specific distribution of the test dataset for each language pair.

rial while ensuring the translations were idiomatic and contextually appropriate for each domain.

The test set release process was intentionally staged to introduce additional complexity and rigour. In the first phase, 500 English sentences were released, requiring participants to translate these into the target languages. This forward translation task required participants to demonstrate their models' proficiency in capturing nuances and domain-specific terminology in the target languages. In the second phase, 500 sentences in the target languages were provided, requiring translation back into English. This reverse translation task assessed the models' ability to accurately render the meaning, tone, and subtleties of the original sentences in English, thus testing bidirectional translation capability. The combined forward and reverse tasks aimed to evaluate the accuracy, fluency, and idiomatic correctness of the translations. The careful selection of diverse domains and the structured release of the test set was intended to challenge the generalization capabilities of the participating models. The goal was to ensure that only the most robust models, capable of handling a wide range of real-world scenarios, would excel.

This approach ensures a rigorous and multifaceted evaluation, capturing the subtleties of each language pair's translation performance across different domains.

## 5 Participants and System Descriptions

In this shared task, total of 12 teams registered and contributed, as indicated in table 8, the released dataset have been distributed among participants. In table 7, we have compiled the system outputs submitted by participants, encompassing both primary and contrastive submission types.

**DLUT-NLP** (Ju et al., 2024): The participant for low-resource translation tasks involving English-Assamese, English-Mizo, English-Khasi,

| Language Pair | Submissions |
|---|---|
| English - Assamese | 11 (primary), 6 (contrastive) |
| English-Mizo | 10 (primary), 5 (contrastive) |
| English-Khasi | 10 (primary), 6 (contrastive) |
| English-Manipuri | 10 (primary), 6 (contrastive) |

Table 4: Number of participants in the low-resource Indic language translations

and English-Manipuri language pairs. It utilized a transformer-based model, with monolingual data for pre-training and parallel data for fine-tuning. Enhancements included back-translation, oversampling, and model averaging, along with knn-mt technology during inference, supported by a datastore created from parallel data.

**A3-108** (Yadav et al., 2024): The team tackled low-resource machine translation by implementing control mechanisms in transformer-based NMT models. They encoded the target sentence length as a control token in the source sentence for eight language pairs: English-Assamese, Manipuri, Khasi, and Mizo. Four variations of this encoding were tested against baseline models. Two systems were submitted for each language pair: a primary system using control tokens based on the target-to-source token length ratio, and a contrastive baseline system without control tokens. All models were trained on the provided dataset.

**SRIB-NMT** (Patil et al., 2024): The team participated in the WMT-24 challenge for translating English to four low-resource Indic languages. They used transformer models for both their primary and contrastive systems. The primary system involved pre-training language models on large amounts of text data before fine-tuning them for translation. The contrastive system improved upon this by further fine-tuning a pre-trained translation model using a technique called LoRA, resulting in better translation quality.

**YES-MT** (Bhaskar and Krishnamurthy, 2024):

The team participated, focusing on four language pairs: English to Assamese, Khasi, Manipuri, and Mizo. Their primary systems used Transformer models trained from scratch. In contrast, contrastive systems applied transfer learning with fine-tuning techniques like LoRA and Supervised Fine-Tuning (SFT) on pre-trained models such as IndicTrans2 and LLaMA 3. Their experiments explored the effectiveness of these approaches, including quantization, in enhancing translation quality for low-resource languages.

**HW-TSC** (Wei et al., 2024)**:** The team participated in the WMT-24 challenge for translating English to four low-resource Indic languages. They used transformer models for both their primary and contrastive systems. The primary system involved pre-training language models on large amounts of text data before fine-tuning them for translation. The contrastive system improved upon this by further fine-tuning a pre-trained translation model using a technique called LoRA, resulting in better translation quality.

**CycleL** (Sören Dréano, 2024)**:** The team developed a novel self-supervised Neural Machine Translation (NMT) model called CycleGN. Unlike traditional NMT models, CycleGN doesn't require parallel data. It utilizes Cycle Consistency Loss (CCL) and Masked Language Modeling (MLM) for training. The model was tested on low-resource language pairs Spanish-Aragonese and Spanish-Asturian using PILAR datasets as part of the WMT24 Shared Task. Despite computational challenges and early training termination, the results demonstrated the potential of self-supervised learning for low-resource translation scenarios.

**NLIP_Lab-IIITH** (Sahoo et al., 2024)**:** The participated team aiming to improve Manipuri and Khasi translations. They utilized mBART and IndicTrans2 models as baselines, incorporating data augmentation techniques like backtranslation and data filtering with fine-tuned LaBSE. Despite limited data, iterative fine-tuning on enhanced datasets led to significant improvements in translation quality, as measured by BLEU, chrF, and TER metrics.

**MTNLP-IIITH** (P M et al., 2024)**:** The team tackled the WMT24 Low-Resource Indic NMT challenge for Manipuri and Khasi, employing mBART and IndicTrans2 models. To overcome data scarcity, they implemented backtranslation and LaBSE-based data filtering. Despite computational constraints, iterative fine-tuning on the pro-

cessed data yielded substantial enhancements in translation quality as assessed by BLEU, chrF, and TER metrics.

**SPRING-IITM** (Sayed et al., 2024)**:** The team developed a robust translation model for four low-resource Indic languages: Khasi, Mizo, Manipuri, and Assamese. They expanded their training corpus using back translation on monolingual datasets and fine-tuned the pre-trained NLLB 3.3B model for Assamese, Mizo, and Manipuri, achieving superior performance over the baseline. For Khasi, they introduced special tokens and trained the model on a custom Khasi corpus, demonstrating significant improvements in translation quality for all four languages.

**JUNLP:** The participant focused on developing a translation system for four low-resource Indic languages: Assamese, Manipuri, Mizo, and Khasi, which are widely spoken in India's North Eastern zone. They combined all language data into a single system using Transformer architecture, enabling translation from English to any of these languages within the same framework. Their approach addresses the challenges posed by the scarcity of data for these languages.

**SRPH-LIT** (Roquea et al., 2024)**:** The team from Samsung R&D Institute Philippines joined the WMT 2024 Low-Resource Indic Language Translation task, focusing on the translation of the following pairs: English ⇔ Assamese, English ⇔ Mizo, English ⇔ Khasi, and English ⇔ Manipuri. In both directions, they adopt the standard sequence-to-sequence Transformer model for translation. The following techniques are data augmentation by back-translation, noisy channel reranking, and checking a multilingual model, which is trained on all the combined language pairs.

**ADAPT-MT** (Gajakos et al., 2024)**:** The ADAPT-MT team participated in the WMT 2024 Low-Resource Indic Language Translation task, focusing on Assamese-to-English and English-to-Assamese. They leveraged Large Language Models (LLMs) as their base systems, employing strategies like fine-tuning with WMT data, few-shot prompting, and efficient data extraction techniques to enhance translation quality. Their approaches were evaluated using BLEU, ChrF, WER, and COMET metrics, showing effective improvements in translating low-resource languages.

| Team Name | Name of University/Lab/Industry/Group |
|---|---|
| AI Lab-IITI | Indian Institute of Technology Indore |
| **NLIPLab-IITH** | **Natural Language and Information Processing Lab at IIT Hyderabad, India** |
| GUIT-NLP | Gauhati University |
| ATULYA-NITS | National Institute of Technology, Silchar |
| Lokkhi | Central Institute of Technology |
| CFILT-IITB | Indian Institute of Technology Bombay |
| CNLP-NITMZ | NIT MIZORAM |
| NITS-CNLP | National Institute of Technology, Silchar |
| DCU-ADAPT | Dublin City University |
| onemt | IIIT-H |
| CL-IIITM | Indian Institute of Information Technology |
| **A3-108** | **International Institute of Information Technology - Hyderabad** |
| **SRIB-NMT** | **Samsung Research Institute** |
| **JUNLP** | **Jadavpur University** |
| GNLP | GKV |
| BVSLP | Banasthali Vidyapith |
| LangMavericks | IIT Madras |
| BITS-P | Birla Institute of Technology & Science, Pilani |
| **DLUT-NLP** | **Dalian University of Technology** |
| JC-beginners | NJIT |
| GUIT-NLP | Gauhati University |
| SHARK | Independent Researcher |
| bjfu | Beijing Forestry University |
| **Yes-MT** | **IIIT Hyderabad** |
| **MTNLP-IIITH** | **LTRC, IIIT Hyderabad, India** |
| **SRPH-LIT** | **Samsung Research Philippines** |
| MUNI-NLP | Masaryk University |
| **CycleL** | **Dublin City University** |
| JUMT | Jadavpur University |
| mbzuai-uhh | MBZUAI, Universität Hamburg |
| Nexus | Z-AGI Labs |
| **SPRING-IITM** | **Indian Institute of Technology, Madras** |
| BV-SLP | Banasthali Vidyapith |
| **HW-TSC** | **Huawei Technologies Co., Ltd.** |
| SAILors | University of New Haven |
| NLIPLab_IITH | Natural Language and Information Processing Lab |
| **ADAPT-MT** | **ADAPT Centre, Dublin City University** |

Table 5: The following table provides an overview of the teams registered for the low-resource Indic language translation task at WMT24 and the datasets provided to them. Participation varied across different language pairs, and only 12 teams in bold completed submissions of both system outputs and system descriptions.

## 6 Results and Discussion

Results for both directions of the four language pairs in WMT 2024 are detailed as follows: English-Assamese in Table 6, English-Mizo in Table 10, English-Khasi in Table 12, and English-Manipuri in Table 8. This section provides the evaluation scores for teams that submitted system outputs and corresponding papers.

Quantitative results are evaluated using established metrics: BLEU, TER, RIBES, ChrF, and METER. BLEU measures the precision of n-grams in candidate translations relative to reference translations. TER quantifies the number of edits required to align the candidate translation with the reference. RIBES evaluates the correlation between the rank orders of words in candidate and

| Team | Test Set | BLEU | TER | RIBES | METEOR | ChrF |
|------|----------|------|-----|-------|--------|------|
| DLUT-NLP | en_to_as_primary | 0.0723 | 85.17 | 0.183 | 0.2205 | 0.3786 |
| | as_to_en_primary | 0.05 | 81.7 | 0.1361 | 0.2907 | 0.3398 |
| A3-108 | en_to_as_contrastive | 0 | 100.46 | 0.0347 | 0.0587 | 0.1817 |
| | as_to_en_contrastive | 0 | 96.44 | 0.0378 | 0.0677 | 0.1803 |
| | en_to_as_primary | 0 | 99.79 | 0.0243 | 0.05134 | 0.1773 |
| | as_to_en_primary | 0 | 96.19 | 0.0322 | 0.0671 | 0.1883 |
| SRIB-NMT | en_to_as_primary | 0.0132 | 101.83 | 0.071 | 0.0744 | 0.2215 |
| | as_en_en_contrastive | 0.2959 | 34.92 | 0.3505 | 0.7409 | 0.6488 |
| YES-MT | en_to_as_contrastive | 0.2568 | 54.63 | 0.306 | 0.5029 | 0.6518 |
| | en_to_as_primary | 0 | 101.78 | 0.0105 | 0.0292 | 0.1123 |
| HW-TSC | en_to_as_primary | 0.2516 | 55.43 | 0.2963 | 0.5124 | 0.6569 |
| | as_to_en_primary | **0.3228** | 32.71 | 0.3625 | 0.7606 | 0.6593 |
| CycleL | en_to_as_primary | 0 | 123.02 | 0.0029 | 0.0061 | 0.0886 |
| | as_to_en_primary | 0 | 101.81 | 0.0075 | 0.0249 | 0.0994 |
| NLIP_Lab-IIITH | en_to_as_primary | 0.2058 | 62.65 | 0.2674 | 0.4539 | 0.6021 |
| | as_to_en_primary | 0.1685 | 55.11 | 0.242 | 0.5746 | 0.5286 |
| | en_to_as_contrastive | 0.185 | 65.79 | 0.2583 | 0.433 | 0.5891 |
| | as_to_en_contrastive | 0.1547 | 58.12 | 0.2312 | 0.5326 | 0.5053 |
| SPRING-IITM | en_to_as_contrastive | **0.2726** | 52.79 | 0.3032 | 0.513 | 0.652 |
| | as_to_en_contrastive | 0.2669 | 39.08 | 0.3308 | 0.7066 | 0.6048 |
| JUNLP | en_to_as_primary | 0 | 134.69 | 0 | 0.0059 | 0.0563 |
| SRPH-LIT | en_to_as_primary | 0 | 1195.25 | 0 | 0.0001 | 0.1852 |
| | as_to_en_primary | 0 | 104.67 | 0.0175 | 0.0513 | 0.166 |
| ADAPT-MT | en_to_as_primary | 0.1612 | 65.96 | 0.2641 | 0.3927 | 0.5673 |
| | as_to_en_primary | 0.318 | 33.56 | 0.3778 | 0.7537 | 0.6551 |
| | as_to_en_contrastive | 0.3227 | 33.63 | 0.372 | 0.7563 | 0.6573 |

Table 6: Performance of teams in the WMT24 low-resource Indic language translation task for the English-Assamese language pair, measured across multiple metrics.

| Team | Test Set | Adequacy | Fluency | Overall Rating |
|------|----------|----------|---------|----------------|
| DLUT-NLP | en_to_as_primary | 2.5 | 3 | 2.75 |
| | as_to_en_primary | 1.8 | 2.4 | 2.1 |
| A3-108 | en_to_as_contrastive | 0.6 | 1 | 0.8 |
| | as_to_en_contrastive | 0.1 | 0.2 | 0.15 |
| | en_to_as_primary | 0.1 | 0.2 | 0.15 |
| | as_to_en_primary | 0 | 0 | 0 |
| SRIB-NMT | en_to_as_primary | 0.4 | 0.6 | 0.5 |
| | as_en_en_contrastive | 3.6 | 4.1 | 3.85 |
| YES-MT | en_to_as_contrastive | 4.3 | 4.5 | 4.4 |
| | en_to_as_primary | 0 | 0 | 0 |
| HW-TSC | en_to_as_primary | 4.1 | 4 | 4.05 |
| | as_to_en_primary | 4.6 | 4.7 | 4.65 |
| CycleL | en_to_as_primary | 0 | 0 | 0 |
| | as_to_en_primary | 0 | 0 | 0 |
| NLIP_Lab-IIITH | en_to_as_primary | 4.2 | 4.1 | 4.15 |
| | as_to_en_primary | 4.1 | 4.1 | 4.1 |
| | en_to_as_contrastive | 3.4 | 4.1 | 3.75 |
| | as_to_en_contrastive | 3.4 | 3.5 | 3.45 |
| SPRING-IITM | en_to_as_contrastive | 4.6 | 4.6 | **4.6** |
| | as_to_en_contrastive | 4.3 | 4.3 | 4.3 |
| JUNLP | en_to_as_primary | 0 | 0 | 0 |
| SRPH-LIT | en_to_as_primary | 0 | 0 | 0 |
| | as_to_en_primary | 0 | 0 | 0 |
| ADAPT-MT | en_to_as_primary | 4.2 | 4.4 | 4.3 |
| | as_to_en_primary | 4.7 | 4.7 | 4.7 |
| | as_to_en_contrastive | 4.8 | 4.8 | **4.8** |

Table 7: Human evaluation of teams in the WMT24 low-resource Indic language translation task for the English-Assamese language pair, assessed based on Adequacy, Fluency, and Overall Rating.

| Team | Test Set | BLEU | TER | RIBES | METEOR | ChrF |
|---|---|---|---|---|---|---|
| DLUT-NLP | en_to_mni_primary | 0.0077 | 96.554 | 0.0697 | 0.0711 | 0.2863 |
| | mni_to_en_primary | 0.0315 | 87.21 | 0.1297 | 0.2131 | 0.3166 |
| A3-108 | en_to_mni_contrastive | 0 | 101.73 | 0.0084 | 0.0179 | 0.1401 |
| | mni_to_en_contrastive | 0.002 | 96.45 | 0.029 | 0.0615 | 0.1865 |
| | en_to_mni_primary | 0 | 101.55 | 0.0072 | 0.0166 | 0.1415 |
| | mni_to_en_primary | 0 | 96.5 | 0.0271 | 0.0635 | 0.1889 |
| SRIB-NMT | en_to_mni_primary | 0 | 104.1 | 0.0191 | 0.0307 | 0.1889 |
| | mni_to_en_contrastive | 0.1889 | 53.05 | 0.2917 | 0.5943 | 0.571 |
| Yes-MT | en_to_mni_primary | 0 | 104.03 | 0.007 | 0.0214 | 0.1102 |
| | en_to_mni_contrastive | 0.0259 | 84.47 | 0.1312 | 0.1605 | 0.4438 |
| HW-TSC | en_to_mni_primary | 0.0211 | 87.93 | 0.1077 | 0.1406 | 0.4218 |
| | mni_to_en_primary | **0.2877** | 42.16 | 0.3532 | 0.6646 | 0.6106 |
| CycleL | en_to_mni_primary | 0 | ERROR | 0.0054 | ERROR | ERROR |
| | mni_to_en_primary | 0 | ERROR | 0.00542 | ERROR | ERROR |
| NLIP_Lab-IIITH | en_to_mni_primary | 0.0258 | 88.53 | 0.1176 | 0.1391 | 0.4062 |
| | mni_to_en_primary | 0.1106 | 67.02 | 0.2303 | 0.4557 | 0.4935 |
| | en_to_mni_contrastive | **0.0279** | 87.22 | 0.1235 | 0.1235 | 0.414 |
| | mni_to_en_contrastive | 0.1159 | 67.49 | 0.2319 | 0.4416 | 0.4748 |
| MTNLP-IIITH | en_to_mni_primary | 0 | 94.77 | 0.0737 | 0.0822 | 0.3325 |
| | mni_to_en_primary | 0.0362 | 94.79 | 0.1136 | 0.1873 | 0.2777 |
| | en_to_mni_contrastive | 0.0064 | 96.46 | 0.0628 | 0.0724 | 0.3191 |
| | mni_to_en_contrastive | 0.0484 | 101.76 | 0.1087 | 0.194 | 0.2662 |
| SPRING-IITM | en_to_mni_contrastive | 0.027 | 84.6 | 0.1185 | 0.1567 | 0.4428 |
| | mni_to_en_contrastive | 0.2088 | 48.77 | 0.3031 | 0.61 | 0.5364 |
| JUNLP | en_to_mni_primary | 0 | 101.25 | 0.0044 | 0.0239 | 0.1471 |
| SRPH-LIT | en_to_mni_primary | 0 | 940.98 | 0 | 0.0001 | 0.1568 |
| | mni_to_en_primary | 0 | 103.28 | 0.0046 | 0.0396 | 0.1729 |

Table 8: Performance of teams in the WMT24 low-resource Indic language translation task for the English-Manipuri language pair, measured across multiple metrics.

reference translations. ChrF assesses the character n-gram F-score, and METER offers a learned metric for translation quality evaluation.

Furthermore, linguistic experts proficient in the target language pairs were engaged for manual evaluations. Twenty sample sentences from the primary submission type were randomly selected for each language pair. Human evaluators assessed the candidate translations based on three criteria: adequacy, fluency, and overall rating. Adequacy gauges how well the candidate translation captures the meaning of the reference. Fluency assesses whether the candidate translation constitutes a well-formed sentence in the target language, independent of its correspondence to the reference. Overall rating integrates both adequacy and fluency to comprehensively evaluate translation quality.

For example, if the reference translation is "The cat sat on the mat," a candidate translation such as "The feline rested on the carpet" is deemed adequate as it preserves the meaning of the reference. In contrast, a candidate translation like "The cat ran across the street," although fluent, is considered inadequate due to the introduction of new information not present in the reference.

The human evaluation parameters are rated on a scale of 0–5, with higher scores reflecting superior quality. The final adequacy, fluency, and overall rating scores are the average ratings assigned to individual test sentences.

## Discussion

For the English-Assamese language pair team, SPRING-IITM achieved a high BLEU score, low TER and an overall rating of 4.6 in the human evaluation. They expanded their training corpus using back translation on monolingual datasets and fine-tuned the pre-trained NLLB 3.3B model. For the Assamese-English language pair team, HW-TSC reaches a higher BLEU score, which is even more than the en-as pair, lower TER and team ADAPT-MT gains a higher overall rating of 4.8 in human evaluation.

For the English-Manipuri language pair team, NLIP_Lab_IIITH achieved a higher BLUE score in automatic evaluation and overall rating in human evaluation compared to the other teams. They utilized mBART and In-dicTrans2 models as baselines, incorporating data augmentation techniques like back translation and data filtering with fine-

| Team | Test Set | Adequacy | Fluency | Overall Rating |
|---|---|---|---|---|
| DLUT-NLP | en_to_mni_primary | 1.95 | 3.6 | 2.75 |
| | mni_to_en_primary | 1.9 | 2.2 | 2.05 |
| A3-108 | en_to_mni_contrastive | 1.5 | 2.2 | 1.85 |
| | mni_to_en_contrastive | 1.0 | 1.15 | 1.075 |
| | en_to_mni_primary | 1.15 | 3.45 | 2.3 |
| | mni_to_en_primary | 1.0 | 1.05 | 1.025 |
| SRIB-NMT | en_to_mni_primary | 1.75 | 2.4 | 2.075 |
| | mni_to_en_contrastive | 3.95 | 3.75 | 3.85 |
| YES-MT | en_to_mni_primary | 1.1 | 2.15 | 3.25 |
| | en_to_mni_contrastive | 4.4 | 4.2 | 4.3 |
| HW-TSC | en_to_mni_primary | 4.1 | 4.6 | 4.35 |
| | mni_to_en_primary | 4.8 | 4.4 | **4.6** |
| CycleL | en_to_mni_primary | 1.25 | 3.65 | 2.45 |
| | mni_to_en_primary | 1.0 | 1.0 | 1.0 |
| NLIP_Lab-IIITH | en_to_mni_primary | 2.35 | 3.95 | 3.15 |
| | mni_to_en_primary | 3.1 | 3.65 | 3.375 |
| | en_to_mni_contrastive | 3.3 | 4.2 | **3.75** |
| | mni_to_en_contrastive | 3.2 | 3.35 | 3.275 |
| MTNLP-IIITH | en_to_mni_primary | 3.1 | 3.7 | 3.4 |
| | mni_to_en_primary | 1.0 | 1.0 | 1.0 |
| | en_to_mni_contrastive | 1.6 | 2.3 | 1.95 |
| | mni_to_en_contrastive | 1.0 | 1.0 | 1.0 |
| SPRING-IITM | en_to_mni_contrastive | 3.25 | 3.75 | 3.5 |
| | mni_to_en_contrastive | 3.8 | 4.06 | 3.93 |
| JUNLP | en_to_mni_primary | 2.7 | 2.4 | 2.55 |
| SRPH-LIT | en_to_mni_primary | 1.65 | 2.3 | 1.975 |
| | mni_to_en_primary | 1.0 | 1.0 | 2.0 |

Table 9: Human evaluation results for teams in the WMT24 low-resource Indic language translation task for the English-Manipuri language pair. The results are presented for Adequacy, Fluency, and Overall Rating on a scale from 0 to 5.

| Team | Test Set | BLEU | TER | RIBES | METEOR | ChrF |
|---|---|---|---|---|---|---|
| DLUT-NLP | en_to_lus_primary | 0.0075 | 98.17 | 0.0725 | 0.1395 | 0.2426 |
| | lus_to_en_primary | 0.0233 | 86.79 | 0.0895 | 0.2622 | 0.3162 |
| A3-108 | en_to_lus_contrastive | 0 | 92.32 | 0.0406 | 0.0978 | 0.18 |
| | lus_to_en_contrastive | 0 | 97.75 | 0.0195 | 0.0544 | 0.1633 |
| | en_to_lus_primary | 0 | 92.84 | 0.0328 | 0.0906 | 0.173 |
| | lus_to_en_primary | 0 | 96.18 | 0.0181 | 0.0587 | 0.1826 |
| SRIB-NMT | en_to_lus_primary | 0 | 102.98 | 0.0361 | 0.062 | 0.1646 |
| | lus_to_en_contrastive | 0.1127 | 64.94 | 0.2026 | 0.4784 | 0.4482 |
| YES-MT | en_to_lus_primary | 0 | 97.19 | 0.0445 | 0.0802 | 0.1282 |
| | en_to_lus_contrastive | 0.0468 | 73.07 | 0.176 | 0.4087 | 0.4151 |
| HW-TSC | en_to_lus_primary | 0.0189 | 86.38 | 0.1074 | 0.1962 | 0.2873 |
| | lus_to_en_primary | 0.0492 | 76.27 | 0.1492 | 0.3646 | 0.3769 |
| CycleL | en_to_lus_primary | 0 | 101.76 | 0.008 | 0.0477 | 0.1645 |
| | lus_to_en_primary | 0 | 100.48 | 0.0064 | 0.0311 | 0.1487 |
| NLIP_Lab-IIITH | en_to_lus_primary | 0.0303 | 81.89 | 0.1479 | 0.2575 | 0.3396 |
| | lus_to_en_primary | 0.0603 | 76.34 | 0.1739 | 0.3716 | 0.3893 |
| | en_to_lus_contrastive | 0 | 98.53 | 0.0277 | 0.0807 | 0.1792 |
| | lus_to_en_contrastive | 0.0849 | 70.28 | 0.1819 | 0.4374 | 0.4188 |
| SPRING-IITM | en_to_lus_contrastive | **0.066** | 66.06 | 0.1746 | 0.495 | 0.4979 |
| | lus_to_en_contrastive | **0.1849** | 53.19 | 0.2684 | 0.588 | 0.5044 |
| JUNLP | en_to_lus_primary | 0 | 98.71 | 0.0589 | 0.0837 | 0.1502 |
| SRPH-LIT | en_to_lus_primary | 0.0025 | 94.46 | 0.0255 | 0.0834 | 0.1891 |
| | lus_to_en_primary | 0 | 108.95 | 0.014 | 0.0421 | 0.1431 |

Table 10: Performance of teams in the WMT24 low-resource Indic language translation task for the English-Mizo language pair, measured across multiple metrics.

| Team | Test Set | Adequacy | Fluency | Overall Quality |
|---|---|---|---|---|
| DLUT-NLP | en_to_lus_primary | 0.6 | 0.5667 | 0.5833 |
| | lus_to_en_primary | 2.7 | 2.7 | 2.7 |
| A3-108 | en_to_lus_contrastive | 0 | 0 | 0 |
| | lus_to_en_contrastive | 0 | 0 | 0 |
| | en_to_lus_primary | 0 | 0 | 0 |
| | lus_to_en_primary | 0 | 0 | 0 |
| SRIB-NMT | en_to_lus_primary | 0 | 0 | 0 |
| | lus_to_en_contrastive | 4.3667 | 4.6667 | 4.5167 |
| YES-MT | en_to_lus_primary | 0 | 0 | 0 |
| | en_to_lus_contrastive | 2.65 | 2.8 | 2.725 |
| HW-TSC | en_to_lus_primary | 0.2333 | 0.1333 | 0.1833 |
| | lus_to_en_primary | 4.2333 | 4.3 | 4.2667 |
| CycleL | en_to_lus_primary | 0 | 0 | 0 |
| | lus_to_en_primary | 0 | 0 | 0 |
| NLIP_Lab-IIITH | en_to_lus_primary | 3.0333 | 3.2667 | 3.15 |
| | lus_to_en_primary | 3.3333 | 3.4333 | 3.3833 |
| | en_to_lus_contrastive | 0 | 0 | 0 |
| | lus_to_en_contrastive | 4.6 | 4.7 | 4.65 |
| SPRING-IITM | en_to_lus_contrastive | 4.5333 | 4.5667 | **4.55** |
| | lus_to_en_contrastive | 4.7667 | 4.8333 | **4.8** |
| JUNLP | en_to_lus_primary | 0 | 0 | 0 |
| SRPH-LIT | en_to_lus_primary | 0 | 0 | 0 |
| | lus_to_en_primary | 0 | 0 | 0 |

Table 11: Updated human evaluation results for teams in the WMT24 low-resource Indic language translation task for the English-Mizo language pair, based on Adequacy, Fluency, and Overall Quality scores.

| Team | Test Set | BLEU | TER | RIBES | METEOR | ChrF |
|---|---|---|---|---|---|---|
| DLUT-NLP | en_to_kha_primary | 0.0665 | 78.17 | 0.1583 | 0.2939 | 0.3512 |
| | kha_to_en_primary | 0.0253 | 81.7 | 0.1223 | 0.2834 | 0.2953 |
| A3-108 | en_to_kha_contrastive | 0.0108 | 92.92 | 0.087 | 0.1209 | 0.1905 |
| | kha_to_en_contrastive | 0 | 105.76 | 0.0094 | 0.0403 | 0.1358 |
| | en_to_kha_primary | 0.011 | 87.69 | 0.0873 | 0.1589 | 0.2296 |
| | kha_to_en_primary | 0 | 107.7 | 0.0071 | 0.0359 | 0.1348 |
| SRIB-NMT | en_to_kha_primary | 0.0054 | 103.72 | 0.0821 | 0.0969 | 0.1778 |
| | kha_to_en_contrastive | 0.042 | 80.29 | 0.1205 | 0.3283 | 0.318 |
| Yes-MT | en_to_kha_primary | 0.0029 | 159.36 | 0.0489 | 0.0511 | 0.1139 |
| | en_to_kha_contrastive | 0.0696 | 80.74 | 0.2167 | 0.2797 | 0.3541 |
| HW-TSC | en_to_kha_primary | 0.0454 | 87.75 | 0.1509 | 0.2134 | 0.2747 |
| | kh_en_primary | 0.0315 | 79.83 | 0.1275 | 0.3044 | 0.3137 |
| CycleL | en_to_kha_primary | 0.0038 | 91.86 | 0.0696 | 0.1399 | 0.2245 |
| | kha_to_en_primary | 0 | 132.21 | 0.0062 | 0.0264 | 0.0973 |
| NLIP_Lab-IIITH | en_to_kha_primary | 0.0475 | 87.16 | 0.1406 | 0.2205 | 0.2894 |
| | kha_to_en_primary | 0.0108 | 92.83 | 0.0742 | 0.1612 | 0.2488 |
| | en_to_kha_contrastive | 0.0521 | 88.9 | 0.1515 | 0.2173 | 0.288 |
| | kha_to_en_contrastive | 0.0312 | 81.23 | 0.1263 | 0.3007 | 0.312 |
| MTNLP-IIITH | en_to_kha_primary | 0.0492 | 84.79 | 0.1595 | 0.2589 | 0.3316 |
| | kha_to_en_primary | 0.0049 | ERROR | 0.25108 | ERROR | ERROR |
| | en_to_kha_contrastive | 0.0359 | 103.49 | 0.1106 | 0.1649 | 0.2333 |
| | kha_to_en_contrastive | 0.006 | 106.6 | 0.0487 | 0.102 | 0.1731 |
| SPRING-IITM | en_to_kha_contrastive | **0.1212** | 63.31 | 0.1864 | 0.4453 | 0.4455 |
| | kha_to_en_contrastive | **0.1047** | 61.43 | 0.2172 | 0.5042 | 0.4271 |
| JUNLP | en_to_kha_primary | 0 | 138.36 | 0.0079 | 0.0094 | 0.0344 |
| SRPH-LIT | en_to_kha_primary | 0.0044 | 126.94 | 0.0533 | 0.0879 | 0.1425 |
| | kha_to_en_primary | 0 | 109.81 | 0.0106 | 0.0407 | 0.1336 |

Table 12: Performance of teams in the WMT24 low-resource Indic language translation task for the English-Khasi language pair, measured across multiple metrics.

| Team | Test Set | Adequacy | Fluency | Overall Quality |
|---|---|---|---|---|
| DLUT-NLP | en_to_kha_primary | 2.43 | 3.2 | 2.815 |
| | kha_to_en_primary | 2.83 | 3.53 | 3.18 |
| A3-108 | en_to_kha_contrastive | 0.33 | 0.6 | 0.465 |
| | kha_to_en_contrastive | 0.33 | 0.36 | 0.345 |
| | en_to_kha_primary | 0.33 | 0.7 | 0.515 |
| | kha_to_en_primary | 1 | 1 | 1 |
| SRIB-NMT | en_to_kha_primary | 0.33 | 0.46 | 0.395 |
| | kha_to_en_contrastive | 3.36 | 3.6 | 3.48 |
| Yes-MT | en_to_kha_primary | 0.33 | 0.33 | |
| | en_to_kha_contrastive | 2.3 | 2.5 | 2.4 |
| HW-TSC | en_to_kha_primary | 0.43 | 0.5 | 0.465 |
| | kha_to_en_primary | 4 | 4.23 | 4.115 |
| CycleL | en_to_kha_primary | 0.33 | 0.4 | 0.365 |
| | kha_to_en_primary | 0.33 | 0.4 | 0.365 |
| NLIP_Lab-IIITH | en_to_kha_primary | 1.66 | 1.76 | 1.71 |
| | kha_to_en_primary | 2.66 | 2.93 | 2.795 |
| | en_to_kha_contrastive | 2.26 | 2.3 | 2.28 |
| | kha_to_en_contrastive | 3.47 | 3.23 | 3.35 |
| MTNLP-IIITH | en_to_kha_primary | 1.93 | 1.93 | 1.931 |
| | kha_to_en_primary | 2.83 | 2.83 | 2.83 |
| | en_to_kha_contrastive | 0.76 | 0.76 | 0.76 |
| | kha_to_en_contrastive | 1.76 | 1.83 | 1.795 |
| SPRING-IITM | en_to_kha_contrastive | 4.56 | 4.93 | **4.745** |
| | kha_to_en_contrastive | 4.93 | 4.96 | **4.945** |
| JUNLP | en_to_kha_primary | 1 | 1 | 1 |
| SRPH-LIT | en_to_kha_primary | 0 | 0 | 0 |
| | kha_to_en_primary | 0 | 0 | 0 |

Table 13: Human evaluation results for teams in the WMT24 low-resource Indic language translation task for the English-Khasi language pair, based on Adequacy, Fluency, and Overall Quality scores.

tuned LaBSE.Team HW-TSC achieved higher BLEU score as well as overall rating in human evaluation for the Manipuri-Englishwhich is significantly higher when compared to the en-mni language pair. They employed a contrastive system which improved upon this by further fine-tuning a pre-trained translation model using a technique called LoRA, resulting in better translation quality.

For the English-Mizo language pair team, SPRING-IITM outperforms all the teams in both directions of the language pairs with higher BLEU and overall ratings in human evaluation. The team developed a robust translation model for four low-resource Indic languages: Khasi, Mizo, Manipuri, and Assamese. They expanded their training corpus using back translation on monolingual datasets and fine-tuned the pre-trained NLLB 3.3B model.

Team SPRING-IITM surpassed all the teams in for the both directions of the language pairs for English-Khasi in automatic and human evaluation. They expanded their training corpus using back translation on monolingual datasets and fine-tuned the pre-trained NLLB 3.3B model. For Khasi, they introduced special tokens and trained the model on a custom Khasi corpus.

## Conclusion

The outcomes of the participating teams in the WMT 2024 translation task for four language pairs have been meticulously evaluated using both automated and human metrics. This year's shared task on low-resource Indic language translation utilised the IndicNE-Corp1.0 dataset from WMT 2023, while a newly developed test set was introduced, characterized by a higher difficulty level than the previous year. This enhanced test set aims to better assess the translation capabilities of the models across the participating languages.

The dataset features four under-resourced languages—Assamese, Mizo, Khasi, and Manipuri—from the northeastern region of India. Future initiatives will focus on expanding the dataset by adding more northeastern Indic languages and increasing the corpus size.

We will be incorporating additional languages in the next iteration of the shared task, including English $\Leftrightarrow$ Nyishi (en-nshi), English $\Leftrightarrow$ Bodo (en-bodo), English $\Leftrightarrow$ Mising (en-mrp), and English $\Leftrightarrow$ Kokborok (en-trp). This expansion aims to enhance the scope of linguistic diversity, allowing

participants to engage with a broader range of low-resource languages.

## Acknowledgements

## References

James F. Allen. 2003. *Natural language processing*, page 1218–1222. John Wiley and Sons Ltd., GBR.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Yash Bhaskar and Parameswari Krishnamurthy. 2024. Yes-mt's submission to the low-resource indic language translation shared task in wmt 2024. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Neha Gajakos, Prashanth Nayak, Rejwanul Haque, and Andy Way. 2024. Adapt-mt submissions to wmt 2024 low-resource indic language translation task. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Loitongbam Gyanendro Singh, Lenin Laitonjam, and Sanasam Ranbir Singh. 2016. Automatic syllabification for Manipuri language. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 349–357, Osaka, Japan. The COLING 2016 Organizing Committee.

Barry Haddow and Faheem Kirefu. 2020. Pmindia – a collection of parallel corpora of languages of india.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 944–952, Cambridge, MA. Association for Computational Linguistics.

Chenfei Ju, Junpeng Liu, Kaiyu Huang, and Degen Huang. 2024. Dlut-nlp machine translation systems for wmt24 low-resource indic language translation. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Vanlalmuansangi Khenglawt, Sahinur Rahman Laskar, Santanu Pal, Partha Pakray, and Ajoy Kumar Khan. 2022. Language resource building and English-to-mizo neural machine translation encountering tonal words. In *Proceedings of the WILDRE-6 Workshop within the 13th Language Resources and Evaluation Conference*, pages 48–54, Marseille, France. European Language Resources Association.

Lenin Laitonjam and Sanasam Ranbir Singh. 2021. Manipuri-English machine translation using comparable corpus. In *Proceedings of the 4th Workshop on Technologies for MT of Low Resource Languages (LoResMT2021)*, pages 78–88, Virtual. Association for Machine Translation in the Americas.

Sahinur Rahman Laskar, Abdullah Faiz Ur Rahman Khilji Darsh Kaushik, Partha Pakray, and Sivaji Bandyopadhyay. 2021. EnKhCorp1.0: An English–Khasi corpus. In *Proceedings of the 4th Workshop on Technologies for MT of Low Resource Languages (LoResMT2021)*, pages 89–95, Virtual. Association for Machine Translation in the Americas.

Sahinur Rahman Laskar, Abdullah Faiz Ur Rahman Khilji, Partha Pakray, and Sivaji Bandyopadhyay. 2020. EnAsCorp1.0: English-Assamese corpus. In *Proceedings of the 3rd Workshop on Technologies for MT of Low Resource Languages*, pages 62–68, Suzhou, China. Association for Computational Linguistics.

Sahinur Rahman Laskar, Riyanka Manna, Partha Pakray, and Sivaji Bandyopadhyay. 2022. A domain specific parallel corpus and enhanced english-assamese neural machine translation. *Computación y Sistemas*, 26(4):1669–1687.

Abhinav P M, Ketaki Shetye, and Parameswari Krishnamurthy. 2024. Mtnlp-iiith: Machine translation for low-resource indic languages. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the WMT 2023 shared task on low-resource Indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694, Singapore. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the*

*40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Pranamya Patil, Raghavendra HR, Aditya Raghuwanshi, and Kushal Verma. 2024. Srib-nmt's submission to the indic mt shared task in wmt 2024. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matthew Theodore Roquea, Carlos Rafael Catalana, Dan John Velascoa, Manuel Antonio Rufinoa, and Jan Christian Blaise Cruza. 2024. Samsung r&d institute philippines @ wmt 2024 indic mt task. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Pramit Sahoo, Maharaj Brahma, and Maunendra Sankar Desarkar. 2024. Nlip-lab-iith low-resource mt system for wmt24 indic mt shared task. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Hamees Sayed, Advait Joglekar, and Srinivasan Umesh. 2024. Spring lab iitm's submission to low resource indic language translation shared task. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Loitongbam Gyanendro Singh and Sanasam Ranbir Singh. 2017. Word polarity detection using syllable features for manipuri language. In *2017 International Conference on Asian Language Processing (IALP)*, pages 206–209.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Noel Murphy Sören Dréano, Derek Molloy. 2024. Exploration of the cyclegn framework for low-resource languages. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Bin Wei, Jiawei Zhen, Zongyao Li, Zhanglin Wu, Jiaxin Guo, Daimeng Wei, Zhiqiang Rao, Shaojun Li, Yuanchang Luo, Hengchao Shang, Jinlong Yang, Yuhao Xie, and Hao Yang. 2024. Machine translation advancements of low-resource indian languages by transfer learning. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

Saumitra Yadav, Ananya Mukherjee, and Manish Shrivastava. 2024. A3-108 controlling token generation in low resource machine translation systems. In *Proceedings of the Low-Resource Indic Language Translation Shared Task, 9th Conference on Machine Translation (WMT24), EMNLP*, November 15-16, 2024, Miami, Florida, USA.

# Findings of WMT 2024's MultiIndic22MT Shared Task for Machine Translation of 22 Indian Languages

**Raj Dabre**[1]     **Anoop Kunchukuttan**[2]

[1]National Institute of Information and Communications Technology, Japan

[2]Microsoft, India

raj.dabre@nict.go.jp     ankunchu@microsoft.com

## Abstract

This paper presents the findings of the WMT 2024's MultiIndic22MT Shared Task, focusing on Machine Translation (MT) of 22 Indian Languages. In this task, we challenged participants with building MT systems which could translate between any or all of 22 Indian languages in the 8th schedule of the Indian constitution and English. For evaluation, we focused on automatic metrics, namely, chrF, chrF++ and BLEU.

## 1 Introduction

India is a linguistically diverse region, with 1,369 distinct natively spoken languages which were identified in the census conducted in 2011. Among these native languages, 22 have been listed in the 8[th] Schedule of the Constitution of India. Furthermore, about 97% of the population of India speaks one of these 22 languages as their first language in their daily lives. It is important to note that English is widely spoken and serves as the default medium of formal communication in many areas, particularly in business, education, government, and judiciary. However, the percentage of the population speaking English is approximately 10% and in the interest of smooth and clear communication, the importance in India of language translation for effective communication, social inclusion, equitable access, and national integrity cannot be over-emphasized.

Having established that Indian language MT is important, the only way to improve it is via active involvement from MT researchers and MT system developers to push the boundaries of translation quality. To this end, we offered the first of its kind shared task focusing on MT for all 22 scheduled Indian languages. Over half a decade ago, in the Workshop on Machine Translation 2018 (WAT 2018) (Nakazawa et al., 2018), the organizers introduced the IndicMT task for the first time

spanning covering 7 Indic languages. Over the years they gradually added languages in WAT from 2018 to 2023 (Nakazawa et al., 2019, 2020, 2021, 2022, 2023), with WAT 2023 boasting 19 Indian languages. Over the years, with the increasing number of languages and datasets for Indian languages, these tasks have garnered growing attention, however the challenge still remains since Indian languages are still resource poor in comparison with European languages.

This year the multilingual Indian languages MT task, referred to as MultiIndic22MT, is hosted under the Ninth Conference on Machine Translation (WMT24) and for the first time ever, the task spans all 22 scheduled languages of India belonging to 4 language families and written in 12 scripts. The languages exhibit both genetic and contact relatedness (Kunchukuttan et al., 2018). Some of these languages are extremely low-resource. This diversity makes this language group ideal for studies in multilingual learning, language relatedness and low-resource MT. Our primary goal behind having this shared task was to attract both researchers and developers to identify effective practices for pushing the quality of Indian language Machine Translation, especially for the lower resourced languages. Our secondary goal was also to identify some interesting but yet unexplored practices, even if they do not lead to state-of-the-art MT performance.

## 2 MultiIndic22MT Shared Task

The task covered English and 22 Indic Languages, as follows:

1. Assamese

2. Bengali

3. Bodo

4. Dogri

5. Konkani

6. Gujarati

7. Hindi

8. Kannada

9. Kashmiri (Arabic script)

10. Maithili

11. Malayalam

12. Marathi

13. Manipuri (Meitei script)

14. Nepali

15. Oriya

16. Punjabi

17. Sanskrit

18. Santali

19. Sindhi (Devanagari script)

20. Tamil

21. Telugu

22. Urdu.

We evaluated user submissions on 44 translation directions (English-Indic and Indic-English). We also evaluate the performance of 5 Indic-Indic pairs: Bengali-Hindi, Tamil-Telugu, Hindi-Malayalam and Sindhi-Punjabi. We encouraged the use of multilingualism and transfer-learning by leveraging monolingual data, back-translation and (potentially) LLMs, to develop high quality systems. Although the intention is to have users develop multilingual systems and submit translations for all directions, we also welcomed submissions for specific language pairs. The link to the shared task page is here[1].

# 3 Datasets and Pre-trained Models

For this shared task, we prepared a fairly extensive list of resources for the participants to train their MT systems. We also describe the evaluation sets.

## 3.1 Datasets

We allowed participants to use existing mined as well as back-translated parallel data along with monolingual data.

**Parallel Data:** As a source of parallel corpora, we recommended using the Bharat Parallel Corpus Collection (BPCC) dataset[2] (Gala et al., 2023) which spans all 22 languages in the shared task. BPCC is a comprehensive and publicly available parallel corpus that includes both existing and new data for all 22 scheduled Indic languages. It comprises two parts: BPCC-Mined and BPCC-Human, totaling approximately 230 million bitext pairs. BPCC-Mined contains about 228 million pairs, with nearly 126 million pairs newly added as a part of this work. On the other hand, BPCC-Human consists of 2.2 million gold standard English-Indic pairs, with an additional 644K bitext pairs from English Wikipedia sentences (forming the BPCC-H-Wiki subset) and 139K sentences covering everyday use cases (forming the BPCC-H-Daily subset). It is worth highlighting that BPCC provides the first available datasets for 7 languages and significantly increases the available data for all languages covered. Note that one may pivot via English to obtain Indic-Indic parallel corpora.

**Parallel Back-translated Data:** Additionally, BPCC also contains back-translation data generated by intermediate checkpoints of IndicTrans2 (Gala et al., 2023) models for training purposes. This data consists of English original sentences translated to 22 Indic languages for a total of 401.9M back-translated sentences and Indian language original sentences translated to English for a total of 400.9M back-translated sentences. The mined, human curated and back-translated corpora represent an extensive training dataset which we expect will be sufficient for training MT systems of reasonable quality.

**Monolingual Data:** We also recommended the use of monolingual data from Varta[3] (Aralikatte et al., 2023), IndicCorp v2[4] (Doddapaneni et al., 2023) and Sangraha[5] (Khan et al., 2024) corpora. Sangraha subsumes IndicCorp v2 but does not explicitly include Varta. Sangraha covers 22 languages, containing a total of 251B tokens, of which con-

tains verified[6] (64B), unverified[7] (24B), and synthetic[8] (162B) tokens. On the other hand, Varta spans only 9B tokens and belongs to the NEWS domain, whereas Sangraha spans multiple domains. Our evaluation sets, which we will describe later, are multi-domain (including news) and hence we expected Sangaraha to be a better source but could not neglect Varta due to its domain specificity and high quality.

## 3.2 Pre-trained Models

In addition to datasets, following recently followed trends in shared tasks, we encouraged participants to leverage one or all of the following publicly available models for fine-tuning or synthetic data generation:

**IndicTrans2 (Gala et al., 2023):** This consists of the 3 IndicTrans2 models, one-to-many, many-to-one, and many-to-many, for English to Indic, Indic to English and Indic to Indic translation. These are the current state-of-the-art open-source MT systems, and we encouraged participants to build on top of these models to improve performance, especially for the lower resourced languages like Santali, Sindhi, Bodo, Dogri, Konkani, Kashmiri, Maithili and Manipuri.

**mT5 (Xue et al., 2021):** mT5 is a well known pretrained model which supports half of the Indian languages in this shared task. However, it is only pre-trained and not fine-tuned for MT and is more suitable for focused domain specific fine-tuning investigations.

**IndicBART (Dabre et al., 2022):** IndicBART is a small pre-trained model for 11 Indic languages and English which, when fine-tuned, is known to outperform mBART (Liu et al., 2020) and give comparable performance as a mT5, despite both models being twice its size.

**VartaT5 (Aralikatte et al., 2023):** This is a T5 model specific for Indic languages and is analogous to IndicBART.

**BLOOM (Workshop et al., 2023):** BLOOM is a family of decoder only pre-trained models supporting 44 languages, some of them being a subset of the Indian languages we focus on in this shared task. Model sizes range from 500 million parame-

ters to 176 billion parameters. However, BLOOM is known to be an under-trained model, and thus we expected participants to focus more on using Gemma.

**Gemma (Team et al., 2024):** Is another family of decoder only models with 2 and 7 billion parameters. Gemma is theoretically capable of tokenizing all 22 Indian languages of this task but its primary support is more in favor of the higher resource languages like Hindi, Marathi, Bengali, etc. We expected that participants would explore some prompting approaches on top of Gemma to determine its viability for Indian language translation.

## 4 Submission Criteria

We expected two types of submissions: Constrained and Unconstrained. Constrained submissions were those which used the data and models stipulated by the organizers explicitly. Unconstrained submissions were those where any other data or models were used without confirmation from the organizers. Furthermore, we encouraged primary and contrastive submissions, where participants could submit one Primary (ranked) and one Contrastive (unranked, optional).

## 5 Evaluation Sets and Metrics

**Evaluation Sets:** We provide participants with a validation set and 3 test sets. The validation set is an extension of FLORES-200 for the 22 Indian languages[9], as described in Gala et al. (2023) and consists of 997 23-way sentences. As for the test sets, 2 out of 3 are publicly available and one is a hidden test set. The publicly available sets are In22-Conv[10] and In22-Gen[11] spanning 1,503 and 1,024 23-way parallel sentences, for the conversational and general styles, respectively. The hidden test set was originally described in Chitale et al. (2024) and is an Indic language original test set where Indic sentences were translated into English by linguists. This is different from all other test sets which are English original and were translated into Indic languages. This hidden test set was released to the participants 2 weeks before the deadline and unlike In22-Conv and In22-Gen, the references

---

[6]The URLs of webpages from which the corpora were crawled were manually verified by linguists.

[7]The urls of webpages from which the corpora were crawled were unverifiable.

[8]These were obtained by translating English documents into Indian languages.

[9]https://indictrans2-public.objectstore.e2enetworks.net/flores-22_dev.zip

[10]https://huggingface.co/datasets/ai4bharat/IN22-Conv

[11]https://huggingface.co/datasets/ai4bharat/IN22-Gen

were kept hidden. This test set covered only 13 of the 22 Indic languages namely, Assamese, Bengali, Bodo, Gujarati, Hindi, Kashmiri, Malayalam, Nepali, Santali, Sanskrit, Sindhi, Telugu and Urdu. While we asked participants to work on translation to and from English for In22-Conv and In22-Gen, for the hidden test set, only translation from Indic to English direction was possible in order to keep the test set hidden[12].

**Evaluation Metrics:** We asked participants to submit their translations to us which we would then evaluate using BLEU (Papineni et al., 2002), chrF (Popović, 2015) and chrF++ (Popović, 2017) using sacreBLEU[13] (Post, 2018). We follow the appropriate tokenization of Indic languages as done by Gala et al. (2023) before computing scores.

# 6 Participants and Submissions

Although 32 teams had registered initially, only 4 teams ended up submitting systems and 3 submitted system description papers (1 withdrew). The teams and their submitted systems are as follows:

## 6.1 BV-SLP Team

The BV-SLP team (Joshi et al., 2024), short for the Banasthali Vidyapith Speech and Language Processing Lab, focused on Sindhi to English translation and only submitted translations for the hidden test set. Their approach focuses on special handling of named entities. They first extract named entities from the source Sindhi sentence and translate it first using a knowledge base of Sindhi-English named entity pairs. This intermediate output is then translated using a NMT system, which is trained to retain the translated named entities and only translate the Sindhi part. To develop the NMT system itself, they converted the existing Sindhi-English parallel corpus into a form where the Sindhi sentences had their named entities replaced with their English translations. This pre-translation approach is well known to work well for handling named entities. They used two approaches for translation itself, one (Primary) where Sindhi is directly translated into English and one (Contrastive) where Sindhi is first translated into Hindi and then into English.

---

## 6.2 NITS-CNLP Team

The NITS-CNLP team (Singh et al., 2024), short for the National Institute of Technology Silchar's Centre for Natural Language Processing, focused on English to Manipuri translation and submitted a primary and a contrastive system. Their approach was rather straightforward, where they used the English-Manipuri data from BPCC (Gala et al., 2023) and trained a transformer model. They submitted results for the In22-Conv and In22-Gen test sets. They also performed some manual evaluations.

## 6.3 NLIP-Lab Team

The NLIP-Lab (Brahma et al., 2024), short for the Natural Language and Information Processing Lab, was the only team that went all out and submitted translations for all translation directions and test sets. The NLIP-Lab team use an approach based on pre-training models using codemixed data which was synthetically created. Specifically, they take BPCC parallel data and replace words in English sentences with semantically similar words of the target Indic language sentences. They then pre-train a model with both the original and code-mixed data. They further refine their pre-trained model with original and code-mixed data obtained only from the high quality BPCC-seed datasets. Finally, they fine-tune their models only on the seed datasets without the code-mixed counterparts. They hypothesized that this leads to fairly strong MT systems.

# 7 Results and Findings

Overall, the NLIP-Lab team got 1st rank for all language pairs, directions and test sets, including In22-Conv, In22-Gen and the hidden test set for Indic to English translation.

## 7.1 Sindhi to English Translation

NLIP-Lab had a contender in the form of BV-SLP team for Sindhi to English translation but where the primary system of BV-SLP got BLEU, chrF and chrF++ scores of 19.4, 44.6 and 43.0, respectively. NLIP-Lab translations scored BLEU, chrF and chrF++ scores of 21.2, 47.1 and 45.5, respectively. This showed that NLIP-Lab's RASP pre-training and fine-tuning approach was definitely better than the named entity handling approach. The likely explanation was that NLIP-Lab used a

lot more parallel data and trained a larger model than their competitior.

## 7.2 English to Manipuri Translation

Once again, NLIP-Lab's contender for the English to Manipuri task was the NITS-CNLP lab. This was for the In22-Conv and In22-Gen test sets. NITS-CNLP got BLEU, chrF and chrF++ scores of 6.4, 28.6 and 26.6 for In22-Conv and 8.1, 32.1 and 29.4 for In22-Gen. However, NLIP-Lab got better scores of 15.2, 43.6 and 41.1 for In22-Conv and 18.2, 48.0 and 45.0 for In22-Gen. This shows that NLIP-Lab's systems are substantially better. However, this is to be expected given that NITS-CNLP did not train massively multilingual models and the latter did.

## 7.3 Did NLIP-Lab Beat IndicTrans2?

Unfortunately, NLIP-Lab's systems did not beat IndicTrans2. For the Indic to English directions, IndicTrans2 was almost 10 BLEU better on In22-Gen and almost 4 BLEU better on In22-Conv. For the English to Indic directions, however, the gap narrowed down to about 2 BLEU. This implies that despite IndicTrans2 being trained on significantly larger data (mostly backtranslated) and in multiple stages, its performance can still be approached by systems not leveraging massive amounts of data. This highlights then need for investigating better approaches for translating into Indic languages. As a side note, these same observations hold for Indic to Indic translation.

## 8 Conclusion

In this report we present the findings of the MultiIndic22MT shared task for machine translation involving 22 Indian languages. Despite the initial enthusiasm shown by participants during task registration, only 3 out of 32 teams submitted their translations and system description papers. Of these 3, only NLIP-Lab submitted translations for all directions and got first rank for all their submissions. Approaches explored varied from named entity replacement, pivot language translation (using Hindi as a pivot), code-mixed pretraining and training from scratch. Overall, code-mixed pre-training stood tall and led to the best systems. However, none of the systems could still beat IndicTrans2, indicating that there is much effort needed for pushing the state of the art for translation involving Indian languages. Given the advent of LLMs and the

focus on decoder-only architectures which are well suited for document level MT, we expect that the next batch of innovations will be focused on the same. However, most LLMs dont support Indic languages that well and thus participants may have to resort to using approaches like transliteration to bridge the gap or even reduce it between the type of scripts that LLMs have seen and those that they have not (J et al., 2024; Dabre et al., 2020, 2022; Gala et al., 2023). We hope that more people will participate in another iteration of this task with interesting approaches.

## References

Rahul Aralikatte, Ziling Cheng, Sumanth Doddapaneni, and Jackie Chi Kit Cheung. 2023. Varta: A large-scale headline-generation dataset for Indic languages. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3468–3492, Toronto, Canada. Association for Computational Linguistics.

Maharaj Brahma, Pramit Sahoo Maunendra, and Sankar Desarkar. 2024. Nlip_lab-iith multilingual mt system forwat24 mt shared task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Pranjal Chitale, Jay Gala, and Raj Dabre. 2024. An empirical study of in-context learning in LLMs for machine translation. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 7384–7406, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Raj Dabre, Chenhui Chu, and Anoop Kunchukuttan. 2020. A survey of multilingual neural machine translation. *ACM Comput. Surv.*, 53(5).

Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh Khapra, and Pratyush Kumar. 2022. IndicBART: A pre-trained model for indic natural language generation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1849–1863, Dublin, Ireland. Association for Computational Linguistics.

Sumanth Doddapaneni, Rahul Aralikatte, Gowtham Ramesh, Shreya Goyal, Mitesh M. Khapra, Anoop Kunchukuttan, and Pratyush Kumar. 2023. Towards leaving no Indic language behind: Building monolingual corpora, benchmark and models for Indic languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12402–12426, Toronto, Canada. Association for Computational Linguistics.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M

Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Jaavid J, Raj Dabre, Aswanth M, Jay Gala, Thanmay Jayakumar, Ratish Puduppully, and Anoop Kunchukuttan. 2024. RomanSetu: Efficiently unlocking multilingual capabilities of large language models via Romanization. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15593–15615, Bangkok, Thailand. Association for Computational Linguistics.

Nisheeth Joshi, Pragya Katyayan, Palak Arora, and Bharti Nathani. 2024. System description of bv-slp for sindhi-english machine translation in multiindic22mt 2024 shared task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Mohammed Khan, Priyam Mehta, Ananth Sankar, Umashankar Kumaravelan, Sumanth Doddapaneni, Suriyaprasaad B, Varun G, Sparsh Jain, Anoop Kunchukuttan, Pratyush Kumar, Raj Dabre, and Mitesh Khapra. 2024. IndicLLMSuite: A blueprint for creating pre-training and fine-tuning datasets for Indian languages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15831–15879, Bangkok, Thailand. Association for Computational Linguistics.

Anoop Kunchukuttan, Mitesh Khapra, Gurneet Singh, and Pushpak Bhattacharyya. 2018. Leveraging orthographic similarity for multilingual neural transliteration. *Transactions of the Association for Computational Linguistics*, 6:303–316.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pretraining for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Toshiaki Nakazawa, Nobushige Doi, Shohei Higashiyama, Chenchen Ding, Raj Dabre, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Yusuke Oda, Shantipriya Parida, Ondřej Bojar, and Sadao Kurohashi. 2019. Overview of the 6th workshop on Asian translation. In *Proceedings of the 6th Workshop on Asian Translation*, pages 1–35, Hong Kong, China. Association for Computational Linguistics.

Toshiaki Nakazawa, Kazutaka Kinugawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida, Makoto Morishita, Ondřej Bojar, Akiko Eriguchi, Yusuke Oda, Chenhui Chu, and Sadao Kurohashi. 2023. Overview of the 10th workshop on Asian translation. In *Proceedings of the 10th Workshop on Asian Translation*, pages 1–28, Macau SAR, China. Asia-Pacific Association for Machine Translation.

Toshiaki Nakazawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida, Anoop Kunchukuttan, Makoto Morishita, Ondřej Bojar, Chenhui Chu, Akiko Eriguchi, Kaori Abe, Yusuke Oda, and Sadao Kurohashi. 2022. Overview of the 9th workshop on Asian translation. In *Proceedings of the 9th Workshop on Asian Translation*, pages 1–36, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Shohei Higashiyama, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, Ondřej Bojar, Chenhui Chu, Akiko Eriguchi, Kaori Abe, Yusuke Oda, and Sadao Kurohashi. 2021. Overview of the 8th workshop on Asian translation. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 1–45, Online. Association for Computational Linguistics.

Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Shohei Higashiyama, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, Ondřej Bojar, and Sadao Kurohashi. 2020. Overview of the 7th workshop on Asian translation. In *Proceedings of the 7th Workshop on Asian Translation*, pages 1–44, Suzhou, China. Association for Computational Linguistics.

Toshiaki Nakazawa, Katsuhito Sudoh, Shohei Higashiyama, Chenchen Ding, Raj Dabre, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, and Sadao Kurohashi. 2018. Overview of the 5th workshop on Asian translation. In *Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation: 5th Workshop on Asian Translation: 5th Workshop on Asian Translation*, Hong Kong. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ningthoujam Justwant Singh, Kshetrimayum Boynao Singh, Ningthoujam Avichandra Singh, and Thoudam Doren Sing. 2024. Wmt24 system description for the multiindic22mt shared task on manipuri language. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open models based on gemini research and technology.

BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Frohberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, Somaieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névéol, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan

Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Daniel McDuff, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Periñán, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrimann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Marianna Nezhurina, Mario Sänger, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aroonsiri, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. 2023. Bloom: A 176b-parameter open-access multilingual language model.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

# Findings of WMT2024 English-to-Low Resource
# Multimodal Translation Task

**Shantipriya Parida**[1]**, Ondřej Bojar**[2]**, Idris Abdulmumin**[3]**,**
**Shamsuddeen Hassan Muhammad**[4]**, Ibrahim Said Ahmad**[5]

[1]Silo AI, Finland; [2]Charles University, MFF, ÚFAL, Czech Republic;
[3]Data Science for Social Impact Research Group, University of Pretoria, South Africa;
[4]Imperial College London, UK; [5]Institute for Experimental AI, Northeastern University, USA
correspondence: shantipriya.parida@silo.ai

## Abstract

This paper presents the results of the English-to-Low Resource Multimodal Translation shared tasks from the Ninth Conference on Machine Translation (WMT2024). This year, 7 teams submitted their translation results for the automatic and human evaluation.

## 1 Introduction

The Ninth Conference on Machine Translation (WMT24), held in conjunction with EMNLP 2024, hosted a number of shared tasks covering various aspects of machine translation (MT). This conference builds on 17 previous editions of WMT as a workshop or a conference. This year, Workshop on Asian Translation (WAT), the most recognized shared task campaign on Asian languages, merged with WMT, adding many new shared tasks to the venue.

Multi-modal translation, which involves incorporating non-text sources alongside text input for machine translation, has gained attention in the past years (Specia et al., 2016; Elliott et al., 2016). However, research in this area has focused on European languages such as English, German, French, Czech, and mainly used two datasets: Flickr30k (Young et al., 2014) and MS-COCO (Lin et al., 2014), where the text caption corresponds to the content of the associated image.

We organized the WMT2024 English-to-LowRes Multimodal Shared Task for Low-Resource Asian and African languages. One important difference is that in our setting, the text caption is attached to a rectangular region of the picture and not the picture as a whole. This approach provides an interesting opportunity to consider not only the broader image but also the localized visual context surrounding the described region, which may provide additional cues for more accurate translation.

## 2 Task and Datasets

In this task, participants were provided with corpora from the Visual Genome dataset in four target language: Hindi, Bengali, Malayalam, and Hausa. The specific datasets are: Hindi Visual Genome 1.1 (HVG, Parida et al., 2019)[1] for Hindi; Bengali Visual Genome (BVG, Sen et al., 2022)[2] for Bengali; Malayalam Visual Genome (MVG, Parida and Bojar, 2021)[3] for Malayalam; and Hausa Visual Genome (HaVG, Abdulmumin et al., 2022)[4] for Hausa. The datasets are split into train, test, dev and challenge test in a parallel fashion. The number of sentences in each split is provided in Table 1. Each split contains items consisting of an image, a highlighted rectangular region within the image ($x, y, width, height$), the original English caption for this region, and the reference translation in the respective target language. These components are illustrated in Figure 1. Depending on the task track, some of these components serve as the source, while others act as references or competing candidate solutions. The specific tracks for this task are listed below.

### 2.1 Text-Only Translation

Labeled "TEXT" in WAT official tables, participants translate short English captions into the target language without using visual information. Additional textual resources are allowed but must be documented in the system description paper.

### 2.2 Captioning

Labeled with the target language code, e.g., "HI," "BN," "ML," "HA", participants generate captions

---

[1]https://lindat.mff.cuni.cz/repository/xmlui/handle/11234/1-3267
[2]https://lindat.mff.cuni.cz/repository/xmlui/handle/11234/1-3722
[3]https://lindat.mff.cuni.cz/repository/xmlui/handle/11234/1-3533
[4]https://lindat.mff.cuni.cz/repository/xmlui/handle/11234/1-4749

| data split | size |
|---|---|
| train | 28,930 |
| dev | 998 |
| test | 1,595 |
| challenge test | 1,400 |

Table 1: Shared task dataset splits



Image_id: 2340461
x: 111 y: 54 width: 342 height: 216

English Text: Two elephants standing in the water
Hindi Text: पानी में खड़े दो हाथी
Bengali Text: পানিতে দাঁড়িয়ে দুটি হাতি
Malayalam Text: വെള്ളത്തിൽ നിൽക്കുന്ന രണ്ട് ആനകൾ
Hausa Text: giwaye biyu da ke tsaye a cikin ruwa

Figure 1: Example of a Data Point (Image Id, Region Detail, Source, and Target Languages

in the target language for the highlighted rectangular region in the input image.

## 2.3 Multi-Modal Translation

Labeled "MM", given an image, a rectangular region within it, and an English caption for that region, participants translate the caption into the target language. Both textual and visual information are available for this task.

## 3 Evaluation Methods

### 3.1 Automatic Evaluation

We evaluated translation results by two metrics: BLEU (Papineni et al., 2002), and RIBES (Isozaki et al., 2010). BLEU scores were calculated using SacreBLEU (Post, 2018). RIBES scores were calculated using `RIBES.py` version 1.02.4.[5] All scores for each task were calculated automatically using the corresponding reference translations by the evaluation system through which the participants make their submissions.

---

[5] http://www.kecl.ntt.co.jp/icl/lirg/ribes/index.html

**Automatic Evaluation System** The automatic evaluation system receives translation results by participants and automatically gives evaluation scores to the uploaded results. As shown in Figure 2, the system requires participants to provide the following information for each submission:

- Human Evaluation: whether or not they submit the results for human evaluation;
- Publish the results of the evaluation: whether or not they permit to publish automatic evaluation scores on the WAT2024 web page;
- Task: the task to which the results belong;
- Used Other Resources: whether or not they used additional resources; and
- Method: the type of the method including SMT, RBMT, SMT and RBMT, EBMT, NMT and Other.

Evaluation scores of translation results that participants permit to be published are disclosed via the WAT2024 evaluation web page. Participants can also submit the results for human evaluation using the same web interface. This automatic evaluation system will remain available even after WMT-WAT2024.

### 3.2 Human Evaluation

In WMT2024, we conducted **JPO adequacy evaluation**.

**JPO adequacy evaluation** The evaluation was carried out by translation experts based on the JPO adequacy evaluation criterion, which was originally defined by Japan Patent Office to assess the quality of translated patent documents.

**Sentence selection and evaluation** For the JPO adequacy evaluation, the 200 test sentences were randomly selected from the test sentences. For each test sentence, input source sentence, translation by participants' system, and reference translation were shown to the annotators. To guarantee the quality of the evaluation, each sentence was evaluated by two annotators. Note that the selected sentences are basically the same as those used in the previous workshop.

**Evaluation Criterion** Table 2 shows the JPO adequacy criterion from 5 to 1. The evaluation is performed subjectively. "Important information" represents the technical factors and their relationships. The degree of importance of each element is also considered in evaluating. The percentages

## SUBMISSION

**Submission:**

| | |
|---|---|
| Human Evaluation: | ☐ human evaluation |
| Publish the results of the evaluation: | ☑ publish |
| Team Name: | ORGANIZER |
| Task: | HINDENMMEVTEXT24en-bn ∨ |
| Submission File: | Choose file   No file chosen |
| Used Other Resources: | ☐ used other resources such as parallel corpora, monolingual corpora and parallel dictionaries in addition to official corpora |
| Method: | SMT ∨ |
| System Description (public): | | 100 characters or less |
| System Description (private): | | 100 characters or less |

Submit

Figure 2: The interface for translation results submission

| Sccore | Description |
|---|---|
| 5 | All important information is transmitted correctly. (100%) |
| 4 | Almost all important information is transmitted correctly. (80%–) |
| 3 | More than half of important information is transmitted correctly. (50%–) |
| 2 | Some of important information is transmitted correctly. (20%–) |
| 1 | Almost all important information is NOT transmitted correctly. (–20%) |

Table 2: The JPO adequacy criterion

in each grade are rough indications for the transmission degree of the source sentence meanings. The detailed criterion is described in the JPO document (in Japanese).[6]

## 4   Baseline Systems

Human evaluations were conducted as pairwise comparisons between the translation results for a specific baseline system and translation results for each participant's system. That is, the specific baseline system served as the standard for human

evaluation.

At WMT2024, we adopted some of neural machine translation (NMT) as baseline systems. The NMT baseline systems consisted of publicly available software, and the procedures for building the systems and for translating using the systems were published on the WAT web page.

**Tokenization**   The shared task datasets come untokenized and we did not use or recommend any specific external tokenizer. The standard OpenNMT-py sub-word segmentation was used for pre/post-processing for the baseline system and each participant used what they wanted.

**NMT Methods**   We used the NMT models for all tasks. For the English→Hindi, English→Malayalam, and English→Bengali Multimodal tasks we used the Transformer model (Vaswani et al., 2018) as implemented in OpenNMT-py (Klein et al., 2017) and used the "base" model with default parameters for the multimodal task baseline. We have generated the vocabulary of 32k sub-word types jointly for both the source and target languages. The vocabulary is shared between the encoder and decoder.

---

[6] http://www.jpo.go.jp/shiryou/toushin/chousa/tokkyohonyaku_hyouka.htm

# 5 Participating Teams and Results

We describe the teams' profiles and submissions as described in their respective description papers. Table 3 shows the team IDs, their respective organizations, and countries.

## 5.1 Systems' Descriptions

**DCU_NMT** participated in the English-to-Hindi track only, developing both text-only and multimodal neural machine translation (NMT) systems. They trained the text-only models from scratch on constrained data and further enhanced them with back-translated data. For the multimodal approach, they used a context-aware transformer to integrate visual features by first encoding the image captions with a BERT model and then concatenating them with the textual input. They reported that while the back-translated text-only model achieved the best performance overall, the multimodal systems, despite lacking back-translated data, outperformed the text-only baseline, indicating the potential of visual context. However, their findings revealed that the impact of visual features was inconsistent, showing less effectiveness on the challenge set, highlighting the need for further exploration into effective multimodal integration.

**ODIAGEN** participated in and reported results for all the tracks, including both text-only and multimodal translation. For text-only translation, they trained the Mistral-7B (Jiang et al., 2023) model to handle English to multiple low-resource languages: Hindi, Bengali, Malayalam, and Hausa. In the multimodal English-to-Hindi task, they employed the PaliGemma-3B (Beyer et al., 2024) model, integrating both image and text inputs. However, their findings revealed that the multimodal systems were suboptimal due to improper normalization of location coordinates, which hindered the models ability to map these coordinates accurately to the provided images. While the PaliGemma-3B model demonstrated strong performance in text translation tasks, it struggled to leverage visual context effectively, underscoring the importance of refining multimodal techniques for better accuracy.

**Arewa_NLP** participated in the English-Hausa text-only translation task, fine-tuning the OPUS-MT-en-ha transformer model. While the system performed well on standard test set, it struggled with the more complex content in the Challenge Test, suggesting a need for further training.

**v036** participated in the English-to-Indic tracks only with the help of visual context. They utilized InternVL2 (Chen et al., 2023) to extract features from the marked image region, which was then passed into a Rapid Automatic Keyword Extraction (RAKE) algorithm to generate keywords for use as hash-tags to provide context to the source text. They then used an LLM (Llama 405B) to generate chain-of-thoughts prompts, consisting the original source and target sentences, extracted keywords as hash-tags and some reasoning why that translation was generated, that serve as training data. Finally, they fine-tuned Llama 8b Instruct model, one for each language, on the generated prompts. They reported that although their predictions were mostly correct, the model failed to generate similar translations as the ground truth, indicating the need for human evaluation as the best method to assess the quality of the translations.

**Brotherhood** participated in all the tracks, leveraging LLMs such as GPT-4o and Claude 3.5 Sonnet to enhance cross-lingual image captioning without traditional training or fine-tuning (Betala and Chokshi, 2024). They used instruction-tuned prompting to generate contextual conversations around cropped images, incorporating the original English captions as context, and translated these conversations into target languages. They employed weighted prompting strategy to balance the original captions with the translated conversations for more descriptive outputs. They reported that their training-free approach minimizes error propagation from flawed datasets while offering flexibility in balancing source fidelity with descriptiveness, demonstrating promise for improving low-resource language datasets. However, they identified challenges such as dependence on LLM APIs, hallucination risks, computational demands, and the limitations of traditional metrics like BLEU for evaluating enriched descriptions, highlighting the need for more comprehensive evaluation methods.

**UNLP** participated in the English-to-Hindi, Malayalam, Bengali, and Hausa tracks. They used visual context to improve translation accuracy, employing a gated fusion mechanism to integrate visual information with textual data, combining the outputs of visual and textual encoders to create context-aware translations. For each language,

| Team ID | Organization | Country |
|---------|-------------|---------|
| DCU_NMT | Dublin City University | Ireland |
| ODIAGEN | Odia Generative AI | India |
| Arewa_NLP | FUTB, BUK, and Arewa Data Science Academy | Nigeria |
| v036 | SCB DataX, Walmart Global Tech | Thailand, India |
| Brotherhood | Indian Institute of Technology Madras | India |
| UNLP | University of Galway, and Lua Health, Galway | Ireland |
| 00-7 | Krutrim AI | India |

Table 3: List of participants who submitted translations for the WMT2024 English-to-LowRes Multimodal Translation Task

they fine-tuned their multimodal model on this combined input, ensuring a nuanced understanding of both linguistic and visual cues. The team reported that while their multimodal model consistently outperformed text-only baselines across BLEU, ChrF2, and TER metrics, some discrepancies with the ground truth translations highlight the importance of incorporating human evaluation for a more reliable assessment of translation quality.

**00-7** competed in three tracks—Image Captioning, Text-only, and Multimodal Translation—for Indic languages, developing a multimodal model that integrates a multilingual LLM with a vision module for improved translation. Their method employs a ViT image encoder to extract visual token embeddings, which are projected into the LLM space through an adapter layer, generating translations autoregressively. They achieved state-of-the-art results for Hindi on the Challenge set, while remaining competitive for other languages. Despite the models success, they observed limited impact of the vision modality on translation quality.

### 5.2 Results

**Automatic evaluation results** Tables 4 to 8 present the automatic evaluation results of the submitted systems, indicating that the systems performed competitively against each other. Despite these promising results, participants expressed a need for human evaluations, as shown in subsequent tables. This reflects a common concern among participants who suspect that their systems may outperform the scores they received, underscoring the importance of qualitative assessments in conjunction with automatic metrics.

**Human evaluation results** Tables 10 and 11 present the adequacy scores after human evalua-

| Lang. | System | ID | Type | RSRC | BLEU | RIBES |
|-------|--------|----|------|------|------|-------|
| en-hi | 00-7 | 7190 | NMT | Yes | 53.40 | 0.842400 |
| en-hi | v036 | 7406 | NMT | No | 43.20 | 0.812507 |
| en-hi | Brotherhood | 7378 | NMT | Yes | 37.90 | 0.795538 |
| en-hi | DCU_NMT | 7372 | NMT | No | 30.30 | 0.710342 |
| en-ml | 00-7 | 7195 | NMT | Yes | 39.80 | 0.739973 |
| en-ml | v036 | 7395 | NMT | No | 33.30 | 0.606598 |
| en-ml | Brotherhood | 7377 | NMT | Yes | 13.60 | 0.428194 |
| en-bn | 00-7 | 7192 | NMT | Yes | 45.30 | 0.796451 |
| en-bn | v036 | 7414 | NMT | No | 33.90 | 0.736029 |
| en-bn | Brotherhood | 7375 | NMT | Yes | 21.70 | 0.644341 |
| en-ha | Brotherhood | 7376 | NMT | Yes | 21.10 | 0.636818 |

Table 4: MMCHMM24 submissions

| Lang. | System | ID | Type | RSRC | BLEU | RIBES |
|-------|--------|----|------|------|------|-------|
| en-hi | 00-7 | 7313 | NMT | No | 54.10 | 0.858322 |
| en-hi | ODIAGEN | 7358 | Other | No | 44.10 | 0.815457 |
| en-hi | DCU_NMT | 7349 | NMT | No | 35.90 | 0.762839 |
| en-ml | 00-7 | 7327 | NMT | Yes | 34.00 | 0.651880 |
| en-ml | ODIAGEN | 7343 | Other | No | 18.10 | 0.505942 |
| en-bn | 00-7 | 7321 | NMT | Yes | 44.20 | 0.789032 |
| en-bn | ODIAGEN | 7336 | Other | No | 35.60 | 0.735341 |
| en-ha | ODIAGEN | 7366 | Other | No | 24.40 | 0.663630 |

Table 5: MMCHTEXT24 submissions

| Lang. | System | ID | Type | RSRC | BLEU | RIBES |
|-------|--------|----|------|------|------|-------|
| en-hi | v036 | 7411 | NMT | No | 44.60 | 0.833853 |
| en-hi | 00-7 | 7325 | NMT | No | 43.70 | 0.813357 |
| en-hi | DCU_NMT | 7351 | NMT | No | 40.60 | 0.806358 |
| en-hi | UNLP | 7392 | NMT | No | 40.30 | 0.800532 |
| en-hi | Brotherhood | 7379 | NMT | Yes | 29.70 | 0.725450 |
| en-ml | 00-7 | 7194 | NMT | Yes | 51.40 | 0.780907 |
| en-ml | v036 | 7396 | NMT | No | 42.70 | 0.700828 |
| en-ml | UNLP | 7393 | NMT | No | 32.20 | 0.626281 |
| en-ml | Brotherhood | 7382 | NMT | Yes | 15.10 | 0.410674 |
| en-bn | 00-7 | 7191 | NMT | Yes | 46.40 | 0.775597 |
| en-bn | v036 | 7418 | NMT | No | 44.10 | 0.737924 |
| en-bn | UNLP | 7391 | NMT | No | 42.10 | 0.766589 |
| en-bn | Brotherhood | 7381 | NMT | Yes | 22.10 | 0.575370 |
| en-ha | UNLP | 7394 | NMT | No | 41.80 | 0.723997 |
| en-ha | Brotherhood | 7380 | NMT | Yes | 17.70 | 0.580239 |

Table 6: MMEVMM24 submissions

tion. The scores reinforce the need for human evaluations to actually determine the quality of multi-

| Lang. | System | ID | Type | RSRC | BLEU | RIBES |
|-------|--------|----|------|------|------|-------|
| en-hi | 00-7 | 7322 | NMT | Yes | 43.30 | 0.812578 |
| en-hi | DCU_NMT | 7348 | NMT | Yes | 42.70 | 0.817949 |
| en-hi | ODIAGEN | 7335 | Other | No | 41.60 | 0.821154 |
| en-ml | 00-7 | 7326 | NMT | Yes | 37.80 | 0.633752 |
| en-ml | ODIAGEN | 7365 | Other | No | 33.10 | 0.668374 |
| en-bn | 00-7 | 7320 | NMT | No | 45.10 | 0.766452 |
| en-bn | ODIAGEN | 7363 | Other | No | 43.70 | 0.789757 |
| en-ha | ODIAGEN | 7344 | Other | No | 49.80 | 0.812898 |
| en-ha | Arewa_NLP | 7314 | SMT | No | 40.70 | 0.755910 |

Table 7: MMEVTEXT24 submissions

| Lang. | System | ID | Type | RSRC | BLEU | RIBES |
|-------|--------|----|------|------|------|-------|
| en-hi | 00-7 | 7385 | NMT | Yes | 2.80 | 0.183643 |
| en-ml | 00-7 | 7389 | NMT | Yes | 0.90 | 0.064375 |
| en-bn | 00-7 | 7386 | NMT | No | 1.80 | 0.105044 |

Table 8: MMEVHI24 submissions

| Lang. | System | ID | Type | RSRC | BLEU | RIBES |
|-------|--------|----|------|------|------|-------|
| en-hi | 00-7 | 7346 | NMT | No | 1.30 | 0.125551 |
| en-ml | 00-7 | 7390 | NMT | Yes | 0.30 | 0.039097 |
| en-bn | 00-7 | 7387 | NMT | Yes | 0.40 | 0.041301 |

Table 9: MMCHHI24 submissions

modal generations. The number of sentences that were marked 4 and 5 (almost all or all information transmitted) in system 7375 Brotherhood in Table 10 indicates a higher performance than what the automatic metrics suggest for the same system in Table 4.

| Lang. | System | ID | JPO adequacy scores | | | | | |
|-------|--------|----|---|---|---|---|---|---|
| | | | # | 1 | 2 | 3 | 4 | 5 |
| en-bn | v036 | 7414 | 1 | 2 | 6 | 29 | 84 | 79 |
| | | | 2 | 7 | 23 | 47 | 85 | 38 |
| en-bn | Brotherhood | 7375 | 1 | 0 | 1 | 16 | 71 | 112 |
| | | | 2 | 1 | 10 | 11 | 46 | 132 |
| en-ha | Brotherhood | 7376 | 1 | 11 | 21 | 40 | 48 | 80 |
| | | | 2 | 16 | 29 | 50 | 68 | 37 |

Table 10: MMCHMM24 Human Evaluations on random 200 Test Sentences

## 6 Conclusion and Future Directions

This paper presents an overview of the English-to-Low Resource Multimodal Translation shared tasks at WMT2024. The task attracted strong participation from numerous teams. Out of these, 7 teams submitted system description papers detailing their approaches and results. In the future, we aim to expand the range of low-resource

| Lang. | System | ID | JPO adequacy scores | | | | | |
|-------|--------|----|---|---|---|---|---|---|
| | | | # | 1 | 2 | 3 | 4 | 5 |
| en-bn | ODIAGEN | 7336 | 1 | 15 | 18 | 55 | 66 | 46 |
| | | | 2 | 46 | 43 | 48 | 40 | 23 |
| en-ha | ODIAGEN | 7366 | 1 | 18 | 29 | 62 | 61 | 30 |
| | | | 2 | 26 | 58 | 66 | 36 | 14 |

Table 11: MMCHTEXT24 Human Evaluations on random 200 Test Sentences

languages, with a particular focus on multimodal translation, and encourage greater participation from more teams.

## Acknowledgements

## Ethical Considerations

The authors do not see ethical or privacy concerns that would prevent the use of the data used in the study. The datasets do not contain personal data. Personal data of annotators needed when the datasets were prepared and when the outputs were evaluated were processed in compliance with the GDPR and national law.

## References

Idris Abdulmumin, Satya Ranjan Dash, Musa Abdullahi Dawud, Shantipriya Parida, Shamsuddeen Muhammad, Ibrahim Sa'id Ahmad, Subhadarshi Panda, Ondřej Bojar, Bashir Shehu Galadanci, and Bello Shehu Bello. 2022. Hausa Visual Genome: A Dataset for Multi-Modal English to Hausa Machine Translation. In *Proceedings of the Language Resources and Evaluation Conference*, pages 6471–6479, Marseille, France. European Language Resources Association.

Siddharth Betala and Ishan Chokshi. 2024. Brotherhood at wmt 2024: Leveraging llm-generated contextual conversations for cross-lingual image captioning. *arXiv preprint arXiv:2409.15052*.

Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. 2024. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*.

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*.

Desmond Elliott, Douwe Kiela, and Angeliki Lazaridou. 2016. Multimodal learning and reasoning. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, Berlin, Germany. Association for Computational Linguistics.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, EMNLP '10, pages 944–952, Stroudsburg, PA, USA. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. Opennmt: Opensource toolkit for neural machine translation. In *Proceedings of ACL 2017, System Demonstrations*, pages 67–72. Association for Computational Linguistics.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision – ECCV 2014*, pages 740–755, Cham. Springer International Publishing.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *ACL*, pages 311–318.

Shantipriya Parida and Ondřej Bojar. 2021. Malayalam visual genome 1.0. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

Shantipriya Parida, Ondřej Bojar, and Satya Ranjan Dash. 2019. Hindi Visual Genome: A Dataset for Multimodal English-to-Hindi Machine Translation. *Computación y Sistemas*, 23(4):1499–1505.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Arghyadeep Sen, Shantipriya Parida, Ketan Kotwal, Subhadarshi Panda, Ondřej Bojar, and Satya Ranjan Dash. 2022. Bengali visual genome 1.0. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

Lucia Specia, Stella Frank, Khalil Sima'an, and Desmond Elliott. 2016. A shared task on multimodal machine translation and crosslingual image description. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 543–553, Berlin, Germany. Association for Computational Linguistics.

Ashish Vaswani, Samy Bengio, Eugene Brevdo, Francois Chollet, Aidan Gomez, Stephan Gouws, Llion Jones, Łukasz Kaiser, Nal Kalchbrenner, Niki Parmar, Ryan Sepassi, Noam Shazeer, and Jakob Uszkoreit. 2018. Tensor2tensor for neural machine translation. In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Papers)*, pages 193–199. Association for Machine Translation in the Americas.

Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78.

# Findings of the WMT 2024 Shared Task
# Translation into Low-Resource Languages of Spain:
# Blending Rule-Based and Neural Systems

**Felipe Sánchez-Martínez,**[†] **Juan Antonio Pérez-Ortiz,**[*†]
**Aarón Galiano-Jiménez,**[†] **Antoni Oliver**[‡]

[†]Universitat d'Alacant  {fsanchez,japerez,aaron.galiano}@ua.es
[*]Valencian Graduate School and Research Network of Artificial Intelligence, ValgrAI
[‡]Universitat Oberta de Catalunya  aoliverg@uoc.edu

## Abstract

This paper presents the results of the Ninth Conference on Machine Translation (WMT24) Shared Task "Translation into Low-Resource Languages of Spain". The task focused on the development of machine translation systems for three language pairs: Spanish–Aragonese, Spanish–Aranese, and Spanish–Asturian. 17 teams participated in the shared task with a total of 87 submissions. The baseline system for all language pairs was Apertium, a rule-based machine translation system that still performs competitively well, even in an era dominated by more advanced nonsymbolic approaches. We report and discuss the results of the submitted systems, highlighting the strengths of both neural and rule-based approaches.

## 1 Introduction

In Spain, a diverse linguistic landscape exists, including, beyond the widely recognized Spanish, other languages such as Basque, Catalan, and Galician. Although Spanish is obviously at the forefront in terms of the volume of resources available for training data-driven machine translation (MT) systems, the capabilities and richness of the other languages should not be underestimated. Basque, Catalan, and Galician, which might have been considered limited in resources in the past, actually possess a significant amount of data that facilitate their integration into modern MT technologies. In fact, these three languages have been recently included among the list of up to 100 languages in well-known multilingual systems such as mBERT[1] (Devlin et al., 2019), XLM-R (Conneau et al., 2020), mBART (Liu et al., 2020), mT5 (Xue et al., 2021) or NLLB-200 (Costa-jussà et al., 2024). However, Spain is home to additional languages with much fewer resources, especially in the form

of bilingual data. This task focuses on three of them, namely, Aragonese, Aranese, and Asturian, all of them Romance languages. In particular, participants were asked to submit MT systems from Spanish into any of these three languages.

An interesting fact about our three low-resource languages is that they have open rule-based MT systems available for the Apertium framework. Apertium (Forcada et al., 2011) is a free/open-source rule-based architecture for MT that consists of a pipeline of modules performing morphological analysis, part-of-speech tagging, lexical transfer, lexical selection, chunk-level or recursive structural transfer, and morphological generation.

Another important aspect is that the target languages of the shared task have undergone various orthographic conventions and standards, and the datasets, as well as the MT systems, available may not necessarily adhere to the current conventions adopted by the language academies and used in the test sets.

**Submission platform.** We utilized the open-source OCELoT platform[2] to collect translation submissions. The platform offers anonymized public leaderboards and has been employed in several previous WMT tasks. Submission privileges were restricted to registered and verified teams with accurate contact information, and each team was limited to a maximum of seven submissions per test set.

**Main goals.** The primary objectives of this shared task can be summarized as follows:

- To push the boundaries of MT system development when the amount of resources is extremely scarce.

- To explore the transferability among low-resource Romance languages when translating from Spanish.

---

[1]https://huggingface.co/google-bert/bert-base-multilingual-cased

[2]https://github.com/AppraiseDev/OCELoT

- To find the best way to use pre-trained models of any kind for the translation between Spanish and low-resource Romance languages.

- To create publicly available corpora for MT development and evaluation.

**Main findings.** The main conclusions of the shared task and insights gained are outlined next:

- The best systems result in automatic evaluation scores that are statistically significantly higher than the baseline system, namely, the Apertium rule-based system.

- However, the absolute differences in BLEU and chrF2 scores are not very large (up to 2 BLEU and 1 chrF2 points) in the case of Aragonese and Aranese, which suggests that the rule-based system may still play a role in the translation of these languages, given the fact that they are considerably less resource-hungry than the neural counterparts.

- In the case of Asturian, while the best systems generally maintain a similar range of differences with Apertium, there is one stand-out system that extends the gap significantly, achieving up to 5 BLEU and more than 3 chrF2 points higher. It is worth noting that this winning system primarily leverages a commercial large language model (LLM) through few-shot learning, sampling a new output if the LLM generates a translation that is unexpectedly short or long compared to the source. This underscores the increasing potential of cutting-edge LLMs and the implications for smaller, specialized systems, which may soon be outpaced by new models, even for low-resource languages like Asturian.

The structure of the paper is as follows. Sec. 2 provides an overview of the three target languages. Sec. 3 then outlines the different submission categories based on the resources used, whereas Sec. 4 describes the training data and resources provided to participants, as well as the development and test data used. Sec. 5 briefly describe the systems submitted within each category. The automatic evaluation results are then reported and discussed in Sec. 6. Finally, Sec. 7 concludes the paper with summarizing remarks.

## 2 Languages

Aragonese (Glottocode[3] arag1245), a Romance language mainly spoken in the Pyrenees valleys of Aragón, is primarily used in rural communities and among older generations; intergenerational transmission is severely at risk. It has around 25 000 speakers (Reyes et al., 2017, Table 5).[4] Although recognized as cultural heritage, it does not hold official status, which hampers its broader use and preservation. Despite these challenges, efforts to revitalize the language continue, supported by educational initiatives and cultural programs.

Aranese (Glottocode aran1260) is a variety of the Occitan language spoken in the Val d'Aran, Catalonia, where it holds official status alongside Catalan and Spanish. It is spoken by approximately 4 500 people (Generalitat de Catalunya, 2019, page 4), though its use has been declining due to the dominance of Spanish and Catalan in the region. Despite its small number of speakers, Aranese remains protected by local laws, and efforts to promote its use in education and public life are ongoing.

Asturian (Glottocode astu1246), another Romance language, is spoken by around 250 000 people (Llera Ramo, 2018, Figure 6) in Asturias, though it lacks official status. Like the other languages, Asturian is recognized and protected as cultural heritage, and there are efforts to increase its presence in schools and public life. Many speakers have a passive understanding of the language, and there is a strong cultural identity linked to it.

## 3 Submission Categories

Participants could submit their work in one of three categories,[5] depending on the corpora used, the models employed, and the reproducibility of the results: *constrained*, *open*, and *closed*.

**Constrained submissions.** These submissions are limited to using only the resources (corpora, dictionaries, Apertium-based systems or data, and orthographic conventions) listed in Section 4.1. Participants may also use publicly available pre-trained language or translation models, as long as their size does not exceed 1 billion parameters (1B),

---

[3] https://glottolog.org

[4] Here and in the following figures, we provide the data for the number of people who can speak the language, including those with even a basic level of proficiency.

[5] The organizers reserved the right to assign submissions to the appropriate category if there was any uncertainty.

as specified in their model cards.[6] This size restriction also applies to neural systems used for auxiliary purposes, such as generating synthetic data. The developed systems could be either bilingual or multilingual, and do not necessarily needed to cover all the target languages.

**Open submissions.** Submissions in this category can utilize any resources (corpora, pre-trained models, etc.) in any language, with no size restrictions, as long as the resources are publicly available under open-source licenses to ensure reproducibility. MT systems or large language models available online also fall into this category, provided that the resulting outputs are made available to the public.

**Closed submissions.** Closed submissions face no restrictions on the availability of resources (corpora, pre-trained models, etc.) used for training.

## 4 Data and Resources

This section describes the training corpora and resources provided to the participants for the *constrained* submissions (Sec. 4.1), as well as the development and test corpora used for all submission categories (Sec. 4.2).

### 4.1 Training Corpora and Resources

**Training data.** The shared task included a *constrained* submission category that restricted the resources participants could use to develop their systems, as outlined in Sec. 3. In addition to the FLORES+ dev set (see Sec. 4.2), which could be used for training or validation, participants in this category were provided with the following resources:

- Any resource from OPUS, particularly the largely uncurated resources available for Spanish–Aragonese,[7] Spanish–Occitan,[8] and Spanish–Asturian.[9] Using data from OPUS included monolingual data on the source or target sides or any other bilingual corpus.

- Data from the PILAR dataset (Galiano-Jiménez et al., 2024b), a collection of low-resource language corpora from the Iberian

Peninsula. PILAR contains monolingual and parallel resources for research and development in Romance languages, with data for Aragonese (monolingual web crawled and literary texts), Aranese (bilingual Spanish–Aranese legal provisions from the Diari Oficial de la Generalitat de Catalunya, web crawled texts, and classic literary works), and Asturian (literary and popular science writings).

Systems submitted to the other categories (open and closed) could use the resources listed above, but they were not restricted to them.

**Language identification.** Participants also had access to tools such as Idiomata Cognitor (Galiano-Jiménez et al., 2024a), a highly accurate language identifier for the target languages and other Romance languages.[10]

**Apertium data.** For participants interested in integrating linguistic data into their systems or generating synthetic data, links were provided to Apertium's resources for Aragonese,[11] Spanish–Aragonese,[12] Aragonese–Catalan,[13] Spanish–Asturian,[14] Asturian,[15] Occitan–Spanish,[16] and Occitan–Catalan.[17]

**Other MT systems.** In addition to Apertium-based MT systems, participants were informed of other available MT systems, which could also follow different orthographic conventions to those used in the test sets: the *traduze*[18] system for Aragonese–Spanish; the *Softcatalà*[19] neural Aranese–Catalan system; and the *eslema*[20] MT system for Asturian–Spanish.

**Dictionaries.** Dictionaries, whether monolingual or bilingual, could serve as valuable complementary resources for participants. The following dictionaries were suggested as potential sources:

---

[6]For example, NLLB-200-600M, among others, meets this requirement: https://huggingface.co/facebook/nllb-200-distilled-600M.
[7]https://opus.nlpl.eu/results/es&an/corpus-result-table
[8]https://opus.nlpl.eu/results/es&oc/corpus-result-table
[9]https://opus.nlpl.eu/results/es&ast/corpus-result-table

[10]https://github.com/transducens/idiomata_cognitor
[11]https://github.com/apertium/apertium-arg
[12]https://github.com/apertium/apertium-spa-arg
[13]https://github.com/apertium/apertium-arg-cat
[14]https://github.com/apertium/apertium-spa-ast
[15]https://github.com/apertium/apertium-ast
[16]https://github.com/apertium/apertium-oci-spa
[17]https://github.com/apertium/apertium-oci-cat
[18]https://traduze.aragon.es/
[19]https://github.com/Softcatala/nmt-softcatala
[20]https://eslema.it.uniovi.es/comun/traductor.php

*Diccionari der aranés*[21] by Institut d'Estudis Aranesi;[22] and the *Diccionariu de la Llingua Asturiana*, available online[23] with a limit of 500 query results.

**Orthographic standards.** Participants were informed that the target languages have exhibited various orthographic conventions over time. The evaluation and test sets adhere to contemporary standards, as supported by their respective language academies. The following documents reflect these standards: *Normes ortogràfiques*[24] by the Academia de la Llingua Asturiana; *Ortografía de l'aragonés*[25] by the Academia Aragonesa de la Lengua; and *Gramatica der occitan aranés*[26] published by the Institut d'Estudis Aranesi.

### 4.2 Development and Test Data

Ad-hoc versions of the FLORES+ datasets were purposefully created for the three languages in the shared task. FLORES+ is a multilingual translation benchmark that began with a limited set of languages (Guzmán et al., 2019), was later expanded to 101 languages (Goyal et al., 2022), and most recently to 200 languages (Costa-jussà et al., 2024). In late 2023, the Open Language Data Initiative[27] (OLDI) took over leadership in extending the dataset to new languages and renamed it FLORES+. Specifically, OLDI proposed a shared task[28] to extend FLORES+ to more languages for the Ninth Conference on Machine Translation (WMT24). The Aragonese, Aranese and Asturian versions of FLORES+ used in this shared task were submitted to the OLDI's task as well.

The sentences in FLORES+ are translations of English sentences sampled equally from Wikinews (an international news source), Wikijunior (a collection of age-appropriate non-fiction books), and Wikivoyage (a travel guide). The dataset consists

of a development set (dev) of 997 sentences and a development test set (devtest) of 1012 sentences.

Participants in this shared task were initially provided with the FLORES+ dev set in March 2024 and encouraged to use it for system development, as it closely mirrors the test set in terms of orthographic, grammatical, and domain aspects. Participants had a deadline of July 12, 2024, to submit translations of the Spanish side of the devtest set. Only after that deadline, was the devtest set for Aragonese, Aranese, and Asturian publicly released.

The following provides a brief overview of the FLORES+ datasets for each language, whereas a more detailed explanation of the creation process is available in the paper by Pérez-Ortiz et al. (2024).

For the Aragonese and Aranese datasets, a first draft of the dev and devtest sets were initially generated using the Spanish–Aragonese and Catalan–Aranese Apertium (Forcada et al., 2011) rule-based system. These machine translations were post-edited by language experts and then reviewed by native speakers, including members of the Academia Aragonesa de la Lengua[29] and the Institut d'Estudis Aranesi.[30] The post-editing step is justified by three factors: the lack of resources to hire qualified translators for a from-scratch translation, the common practice of post-editing for these languages, and the high degree of similarity between Spanish and these languages, which makes Apertium translations reliable and less prone to unnatural *translationese*. In the case of Asturian, professional translations originally included in FLORES-101 were reviewed twice by native speakers, including members of the Academia de la Llingua Asturiana.[31]

Pérez-Ortiz et al. (2024, Table 2) report the extent of changes made to the dev and devtest sets after both the initial and final revisions by the language academies. The data reveal significant modifications to the output of Apertium, with a TER score[32] of approximately 26% for Aragonese, 64% for Aranese, and 7% for Asturian, after the two rounds of revision.

---

[21] https://www.diccionari.cat/cerca/diccionari-der-aranes

[22] A PDF version can be downloaded from http://www.institutestudisaranesi.cat/wp-content/uploads/2021/04/DICCIONARI-DER-ARANÉS.pdf.

[23] https://diccionariu.alladixital.org/

[24] https://alladixital.org/wp-content/uploads/2024/01/Normes-Ortografiques-8a-edicion-FINAL-3.pdf

[25] https://academiaaragonesadelalengua.org/sites/default/files/ficheros-pdf/ortografia-aragones.pdf

[26] http://www.institutestudisaranesi.cat/wp-content/uploads/2021/04/gramatica-aranes.pdf

[27] https://oldi.org

[28] https://www2.statmt.org/wmt24/open-data.html

[29] https://academiaaragonesadelalengua.org

[30] http://www.institutestudisaranesi.cat

[31] https://www.academiadelallingua.com

[32] The translation error rate (TER) metric (Snover et al., 2006) is employed here to quantify the number of edits needed to transform the sentences from the initial versions into their corresponding counterparts in the final corpus.

## 5 Teams Participating in the Shared Task

We received a total of 87 submissions from 17 different teams. Table 1 lists the teams that participated in the shared task, along with the language pairs they worked on and the reference, if available, to their system description paper.

Along with their translations of the test set, participants submitted an extended abstract describing their systems and the resources used. Based on that information, we provide a brief overview of the systems developed by each of the participants.

**CUNI-GA.** The CUNI-GA team's contribution (Hrabal et al., 2024) for the three language pairs in the shared task involved the QLoRA fine-tuning of two open-source large language models (LLMs): Aya-23-8B and Command-R 35B. They used a small back-translated dataset, specifically the literary section of the PILAR corpus (Galiano-Jiménez et al., 2024b), which was back-translated using Apertium. Both LLMs were fine-tuned with a single joint model covering all the languages.

**CycleL.** The Dublin City University presented two systems (Spanish–Aragonese and Spanish–Asturian) to the constrained task (Dréano et al., 2024). They employed CycleGN, a fully self-supervised NMT framework that does not rely on parallel data. For this shared task, they exclusively used the PILAR corpus, applying sentence permutations to ensure the dataset remained non-parallel.

**Helsinki-NLP.** This team submitted models (de Gibert et al., 2024) exclusively to the unconstrained open track. Alongside the data provided in the task, such as PILAR, they utilized additional monolingual resources like Wikipedia dumps and dictionary definitions. To enhance their training data, they generated synthetic parallel data through back-translation using an OPUS-MT model. Their data filtering process incorporated language identification using the Idiomata Cognitor tool, as well as the OpusCleaner (Bogoychev et al., 2023) and OpusFilter (Aulamo et al., 2021) tools, to clean and refine their datasets.

For their models, Helsinki-NLP considered various initial systems, including OPUS-MT models (Tiedemann et al., 2024) and different sizes of NLLB-200 models, ranging from 600M to 3.3B parameters. They ultimately chose a multilingual OPUS-MT model based on the transformer-big architecture and produced an ensemble model

after fine-tuning. Their other submissions used sequence-level distillation to train smaller student models that integrated rule-based translation. This was done by translating parallel sentences using both their neural best system and Apertium, selecting the output with the best chrF score relative to the reference, and training smaller transformer-based models on the distilled data. The sizes of distilled models ranged from the transformer *base* architecture to even smaller models obtained via the OpusDistillery tool.[33] Their different models showed statistically significant differences in most cases, except for Asturian, with the distilled models providing competitive translation performance.

**HW-TSC.** The Huawei Translation Service Center participated in the *constrained category* by submitting three systems, one for each of the target languages. Their submissions (Luo et al., 2024) were based on a transformer-big architecture with an expanded number of encoder layers (25). They started by training multilingual systems on sampled training data to obtain both one-to-many and many-to-one pre-trained models, which were then further trained on the original bilingual data to create translation models between Spanish and Aragonese, Aranese, and Asturian in both directions. Additionally, they utilized synthetic corpora generated via Apertium (forward translation) and through back-translation using the aforementioned multilingual models. LaBSE (Feng et al., 2022) denoising was applied to filter out noisy parallel sentences from both the provided training data and the generated synthetic data. Finally, transductive ensemble learning was employed to aggregate multiple models for inference.

**ILENIA-MT.** For the constrained submission (Sant et al., 2024), the team leveraged synthetic corpus generation through Apertium, primarily using data from OPUS and PILAR. Synthetic data was generated by translating from Spanish to Aragonese and Aranese (pivoting through Catalan in this case) using Apertium, while for Asturian, the team directly used NLLB-200-600M. Additional monolingual data was sourced from orthography dictionaries as supplementary resources. A comprehensive data filtering process was applied, involving the removal of noisy sentences using LABSE-based

---

| Submission Name | Language Pairs | System Description |
|---|---|---|
| Apertium (baseline) | Aragonese, Aranese, Asturian | (Forcada et al., 2011) |
| CUNI-GA | Aragonese, Aranese, Asturian | (Hrabal et al., 2024) |
| CycleL | Aragonese, Asturian | (Dréano et al., 2024) |
| Helsinki-NLP | Aragonese, Aranese, Asturian | (de Gibert et al., 2024) |
| HW-TSC | Aragonese, Aranese, Asturian | (Luo et al., 2024) |
| ILENIA-MT | Aragonese, Aranese, Asturian | (Sant et al., 2024) |
| imaxin | Asturian | (González, 2024) |
| LCT-LAP | Aragonese, Aranese, Asturian | (Bär et al., 2024) |
| Mora translate | Asturian | (Menan et al., 2024) |
| SJTU-MT | Aragonese, Aranese, Asturian | (Hu et al., 2024) |
| SRPH-LIT | Aragonese, Aranese, Asturian | (Velasco et al., 2024) |
| Stevens Inst. of Tech. | Aragonese | (no associated paper) |
| TAN-IBE | Aragonese, Aranese, Asturian | (Oliver, 2024) |
| TIM-UNIGE | Aragonese, Aranese | (Mutal and Ormaechea, 2024) |
| TRIBBLE | Aragonese, Aranese, Asturian | (Kuzmin et al., 2024) |
| UAlacant | Aragonese, Aranese, Asturian | (Galiano-Jiménez et al., 2024) |
| Vicomtech | Aragonese, Aranese, Asturian | (Ponce et al., 2024) |
| Z-AGI Labs | Aragonese, Aranese, Asturian | (no associated paper) |

Table 1: Participants in the WMT24 Shared Task "Translation into Low-Resource Languages of Spain". Apertium has its own row, but it is not an actual participant; it rather serves as the baseline system.

embeddings (Feng et al., 2022), sentence length filtering, and language recognition with the Idiomata Cognitor tool. The NLLB-200-600M model was then fine-tuned with all the resulting parallel and synthetic data. To handle unsupported languages in NLLB, new language tags were added for Aragonese and Aranese, initialized with embeddings from Spanish and Occitan, respectively.

For the open submission (Sant et al., 2024), ILENIA-MT used Apertium to generate a large amount of synthetic data, translating 30 million sentences sourced from Spanish monolingual corpora. A transformer model was then trained from scratch. This approach resulted in slightly lower scores than the constrained submission.

**imaxin software.** This team presented an improved version of the Apertium system for the Spanish–Asturian language pair (González, 2024). The team has enhanced Apertium both in terms of syntax, by developing new constraint grammar and transfer rules, and in the lexical domain, by expanding the dictionaries.

**LCT-LAP.** The University of the Basque Country submitted three systems to the *constrained category* (Bär et al., 2024). These systems were obtained by fine-tuning OPUS-MT pre-trained mod-

els for two high-resource Romance languages: Spanish–Galician was used as the starting point for Spanish–Asturian, and Spanish–Catalan was used for Spanish–Aragonese and Spanish–Aranese. The fine-tuning was conducted on OPUS corpora, with noisy parallel sentences removed from the provided training data, and on synthetic corpora generated with Apertium by translating monolingual corpora in PILAR. Before utilizing the OPUS corpora, Idiomata Cognitor was employed to remove parallel sentences not in the desired language, and Apertium was then employed to translate one side of the parallel corpus, followed by BLEU scoring to filter out low-quality parallel sentences.

**Mora translate.** This team participated in the Spanish–Asturian language pair with a constrained submission (Menan et al., 2024). Their main contribution is a dual-stage data filtering system that combines statistical methods for both bilingual and monolingual data, along with a filtering method based on Jensen-Shannon divergence (Lin, 1991). They used the filtered CCMatrix and Wikimedia corpora, and utilized the PILAR corpus for Asturian and the Spanish portion of the English–Spanish Wikimedia corpus as monolingual data. Training was conducted in two phases: (1) training the entire model using the filtered Spanish–Asturian CCMatrix, and (2) fine-tuning the

best model by unfreezing only the decoder. For fine-tuning, several datasets were combined, including monolingual data translated with NLLB-200-600M.

**SJTU-MT.** The systems submitted by this team (Hu et al., 2024) are based on strategies that differ significantly from traditional methods. The submissions covered all three target languages with notable variations in approach for each language pair.

For Aragonese and Aranese, the team generated a pseudo-parallel corpus using Apertium. They sampled one million Spanish sentences from the NLLB Spanish corpus in OPUS, then translated these into Aragonese and Aranese using Apertium to create a synthetic parallel corpus. Both models used a small LLM, Qwen2-0.5B,[34] which was first fine-tuned on this synthetic corpus.

For Aragonese, the model underwent an additional step of few-shot fine-tuning. This involved assembling five-shot examples using sentences from the FLORES+ dev set, providing these as context before training the model to translate sentences. At inference time, when a new sentence is inputted for translation, they use the BM25 ranking function (Robertson and Zaragoza, 2009) to identify the five most relevant examples from the FLORES+ dev set to replicate the same few-shot format introduced during training.[35]

For Aranese, after supervised fine-tuning of Qwen2 on the pseudo-corpus, an additional step involved applying the recently proposed *contrastive preference optimization* (CPO) algorithm (Xu et al., 2024). This method, which moves beyond standard training that replicates a reference translation, employs reinforcement learning loss functions to push models towards preferred translations while steering away from suboptimal ones. In order to apply CPO to their model, the team used Apertium translations as the least preferred and the target references from FLORES+ as the most preferred. Despite discouraging Apertium-like translations at times during training, this process improved the system's performance for Aranese.[36]

For Aragonese and Aranese, if the generated translations were significantly shorter or longer than the input, they were replaced with Apertium translations.

For Asturian, the team employed a completely different approach and participated in the open track. They used the large language model Claude 3.5 Sonnet,[37] utilizing a simple prompting strategy: when translating a new FLORES+ devtest sentence, they retrieved the 20 most similar examples from the FLORES+ dev set using the BM25 ranking function, providing these as suggestions to the model. If the translations produced were significantly shorter or longer than the input, rather than relying on Apertium as before, they simply resampled the model's output until achieving a translation within an acceptable length range.

**SRPH-LIT.** Samsung R&D Institute Philippines submitted three translation systems to the *constrained category* (Velasco et al., 2024), each addressing one of the three language pairs with a standard sequence-to-sequence transformer architecture. For each language pair, three systems were trained and combined using a noisy-channel re-ranking strategy to enhance output selection during decoding. The training data included filtered OPUS corpora —using ratio-based and LaBSE-based embedding methods— as well as synthetic data generated through back-translation with Apertium. Due to limited direct translation support in Apertium, translations for Aragonese–Spanish followed the path Aragonese–Catalan, Catalan–Interlingua, Interlingua–Spanish, while Aranese–Spanish followed the path Aranese–Catalan, Catalan–Spanish.

**Stevens** The Stevens Institute of Technology participated with a *constrained* model for Spanish–Aragonese. They leveraged NLLB-200 via a multi-stage fine-tuning process applied to both Aragonese–Spanish and Spanish–Aragonese translations. To supplement the limited available parallel data, they developed a back-translation system, generating synthetic parallel data and refining it by selecting the top 20% based on L2 cosine similarity. This iterative process enhanced both the back-translation model and the final forward translation system. The final system was an ensemble, created by averaging the weights of the two best-performing models on the development corpus.

---

[34] https://github.com/QwenLM/Qwen

[35] Notably, this approach leveraged the dev set not only during training but also during inference, as opposed to many other systems.

[36] This confirms that the use of post-editing of Apertium translations as an initial step in obtaining the FLORES+ data for our target languages did not overly bias the translations toward an Apertium-like style.

[37] https://www.anthropic.com/news/claude-3-5-sonnet

**TAN-IBE.** The TAN-IBE team (Oliver, 2024) presented systems for all language pairs in the shared task. To address the lack of resources and the low quality of existing corpora, the team: (a) cleaned the existing corpora; (b) created new corpora from Wikipedia; (c) experimented with back-translation and synthetic corpora; and (d) explored multilingual systems. All training was conducted using a transformer Marian-NMT model. For the Spanish–Asturian pair, they submitted an open system using a cleaned version of the NLLB corpus and the newly created Wikipedia corpus. For Spanish–Aragonese and Spanish–Aranese, they submitted constrained systems, using the cleaned existing corpora, the new Wikipedia corpus, and synthetic corpora generated with Apertium.

**TIM-UNIGE.** This team started by generating synthetic data for both forward and back-translation (Mutal and Ormaechea, 2024). They employed a two-phase synthetic data generation strategy using the BLOOMZ-560M model (Muennighoff et al., 2023) to fit within the constraints of the task. In the first phase, they fine-tuned BLOOMZ[38] on monolingual data from the target and related languages, using the task of predicting the next token. The FLORES+ dataset served as the validation set. This approach allowed the model to generate additional monolingual text, which they obtained by sampling various prefix lengths from the FLORES+ sentences and completing them through the model.

This synthetic monolingual text was then passed through Apertium to create additional parallel data, which was used to train their models. The training involved either building a transformer from scratch, using a pretrained Helsinki model (72M parameters), or the NLLB-200-600M model. A curriculum learning strategy was employed during training, where multiple phases gradually incorporated smaller subsets of higher-quality data, reduced the learning rate, and shortened the training steps. The final step involved fine-tuning exclusively on the FLORES+ dataset.

**TRIBBLE.** Universitat Pompeu Fabra and the Polish Academy of Science jointly submitted a model in the *constrained category* for translating into the three languages addressed in the shared task (Kuzmin et al., 2024). Their system, built on

the NLLB-200-600M model, was trained on corpora sampled from OPUS and PILAR, along with synthetic data generated using Apertium. They further refined the data by utilizing Idiomata Cognitor and fastText to filter out sentence pairs in undesired languages. To reduce noise in the parallel corpus, they translated the Spanish side into the target language using Apertium and computed similarity scores based on the Levenshtein distance, discarding low-similarity bilingual sentence pairs.

**UAlacant.** The models submitted by Universitat d'Alacant (Galiano-Jiménez et al., 2024) use both parallel and monolingual corpora, supplemented by synthetic corpora, in the three translation directions of this task. The systems use corpora from OPUS and PILAR, together with synthetic data generated by Apertium. All submissions are classified as open due to the use of the NLLB-200-1.3B model, which exceeds the 1B parameter limit. For each translation direction, they submitted three models by fine-tuning NLLB-200.

The first approach combines parallel corpora for translation with monolingual data used in a denoising task, helping the system to learn from target language corpora even in the absence of sufficient bilingual data. A second approach introduces synthetic data, including both back-translation, where monolingual target language texts are translated into Spanish, and synthetic corpora generated by translating Spanish texts into the target languages. The third approach trains on multiple language pairs simultaneously, including Spanish, Aragonese, Asturian, Aranese and related Romance languages, such as Catalan, Galician and Valencian, using their linguistic similarities to enhance knowledge transfer and improve performance across languages. This results in a multilingual system capable of translating between all the languages.

**Vicomtech.** This team submitted systems for both the *constrained* and *open categories* (Ponce et al., 2024). For the *constrained category*, they exploited synthetic data generation using Apertium, like many other participants, combining it with the available parallel data. A filtering process that included language identification using Idiomata Cognitor, cross-lingual embeddings, and sentence length ratios was considered to clean the training data. Their neural translation models were built using transformer-base architectures (6 layers in the encoder and 6 in the decoder), alongside the NLLB-

---

[38]BLOOMZ has the characteristic of having been exposed to the FLORES datasets during its training.

| Rank | id | Team | BLEU |
|------|-----|------|------|
| **Open submission** | | | |
| 1 | 636 | ILENIA-MT | 62.7 |
| | 637 | ILENIA-MT | 62.6 |
| 2 | — | Apertium | 61.1 |
| | 633 | Vicomtech | 61.0 |
| 3 | 554 | UAlacant | 60.2 |
| 4 | 495 | UAlacant | 59.8 |
| 5 | 577 | Helsinki-NLP | 52.7 |
| 6 | 549 | Helsinki-NLP | 51.5 |
| 7 | 563 | Helsinki-NLP | 50.6 |
| 8 | 523 | Helsinki-NLP | 49.1 |
| 9 | 548 | UAlacant | 37.8 |
| 10 | 647 | CUNI-GA | 36.1 |
| **Constrained submission** | | | |
| 1 | 663 | SJTU-MT | 63.2 |
| | 607 | HW-TSC | 63.0 |
| 2 | 529 | ILENIA-MT | 62.3 |
| | 558 | ILENIA-MT | 62.2 |
| | 634 | ILENIA-MT | 62.2 |
| | 642 | TIM-UNIGE | 61.9 |
| 3 | 526 | ILENIA-MT | 61.6 |
| | — | Apertium | 61.1 |
| | 504 | Vicomtech | 61.1 |
| | 644 | TIM-UNIGE | 61.1 |
| 4 | 586 | TIM-UNIGE | 60.7 |
| | 539 | TIM-UNIGE | 60.5 |
| 5 | 649 | Stevens | 59.8 |
| 6 | 613 | Stevens | 57.5 |
| | 584 | TAN-IBE | 57.3 |
| 7 | 622 | TRIBBLE | 49.2 |
| 8 | 530 | LCT-LAP | 38.9 |
| | 662 | Stevens | 38.7 |
| | 513 | Stevens | 37.5 |
| 9 | 507 | SRPH-LIT | 28.2 |
| 10 | 534 | Z-AGI Labs | 24.3 |
| 11 | 533 | Z-AGI Labs | 22.1 |
| 12 | 591 | CycleL | 0.2 |

Table 2: BLEU scores computed over the FLORES+ devtest set for Spanish–Aragonese.

| Rank | id | Team | BLEU |
|------|-----|------|------|
| **Open submission** | | | |
| 1 | — | Apertium | 28.8 |
| | 627 | Vicomtech | 28.8 |
| 2 | 555 | UAlacant | 28.5 |
| 3 | 587 | ILENIA-MT | 27.3 |
| | 656 | CUNI-GA | 27.1 |
| | 552 | UAlacant | 27.0 |
| 4 | 578 | Helsinki-NLP | 24.3 |
| 5 | 562 | Helsinki-NLP | 22.4 |
| 6 | 550 | Helsinki-NLP | 22.1 |
| 7 | 524 | Helsinki-NLP | 21.6 |
| **Constrained submission** | | | |
| 1 | 621 | SJTU-MT | 30.4 |
| | 641 | TIM-UNIGE | 30.2 |
| | 527 | ILENIA-MT | 30.1 |
| | 619 | TIM-UNIGE | 30.1 |
| | 617 | TIM-UNIGE | 30.0 |
| | 640 | TIM-UNIGE | 29.9 |
| | 625 | Vicomtech | 29.8 |
| 2 | 575 | TIM-UNIGE | 28.9 |
| | — | Apertium | 28.8 |
| 3 | 494 | TIM-UNIGE | 28.2 |
| 4 | 610 | TAN-IBE | 26.9 |
| 5 | 608 | HW-TSC | 26.3 |
| 6 | 623 | TRIBBLE | 23.9 |
| 7 | 531 | LCT-LAP | 21.8 |
| 8 | 581 | SRPH-LIT | 7.7 |
| 9 | 536 | Z-AGI Labs | 3.8 |
| 10 | 535 | Z-AGI Labs | 3.7 |

Table 3: BLEU scores computed over the FLORES+ devtest set for Spanish–Aranese.

LLMs might have, even if this knowledge is likely limited due to exposure to only small amounts of text. Their *open category* models combined data from Apertium, other NMT systems, and LLM-generated data, resulting in slightly better scores for Asturian over the constrained models.

**Z-AGI Labs.** This team participated in all language pairs of the shared task. They fine-tuned the NLLB and Helsinki-NLP/OpusMT models using the OPUS dataset provided on the shared task website.

## 6 Results and Discussion

We measured the translation quality of the different systems submitted to the shared task when translating the FLORES+ devtest dataset by means of

200-600M model, finding the non-pretrained neural models to work slightly better.

For the *open category*, a key highlight of their approach, as emphasized by the authors, was the use of the Llama3-8B LLM (Dubey, 2024) to generate synthetic data in the reverse direction, i.e., from the target low-resource languages into Spanish. This approach allowed their MT systems to exploit whatever knowledge of the target languages the

BLEU (Papineni et al., 2002)[39] and chrF2 (Popović, 2015).[40]. We did not use neural-based metrics, such as COMET (Rei et al., 2020), as they are not available for the target languages. Neither did we conduct a manual evaluation because of the lack of resources to hire qualified translators.

| Rank | id | Team | BLEU |
|------|-----|------|------|
| **Open submission** | | | |
| 1 | 576 | SJTU-MT | 23.2 |
| | 551 | Helsinki-NLP | 18.2 |
| | 564 | Helsinki-NLP | 18.0 |
| 2 | 579 | Helsinki-NLP | 18.0 |
| | 568 | TAN-IBE | 18.0 |
| | 629 | Vicomtech | 18.0 |
| 3 | 525 | Helsinki-NLP | 17.9 |
| 4 | 553 | UAlacant | 17.4 |
| | — | Apertium | 17.0 |
| 5 | 556 | UAlacant | 16.9 |
| | 497 | UAlacant | 16.8 |
| 6 | 632 | ILENIA-MT | 16.7 |
| 7 | 648 | CUNI-GA | 15.2 |
| **Constrained submission** | | | |
| 1 | 606 | HW-TSC | 19.8 |
| 2 | 528 | ILENIA-MT | 18.4 |
| | 624 | TRIBBLE | 17.9 |
| | 547 | Mora translate | 17.6 |
| | 630 | Vicomtech | 17.6 |
| | 532 | LCT-LAP | 17.5 |
| | 546 | SRPH-LIT | 17.5 |
| 3 | 522 | Mora translate | 17.4 |
| | 590 | Mora translate | 17.4 |
| | 519 | Mora translate | 17.4 |
| | 512 | Mora translate | 17.4 |
| | 543 | Mora translate | 17.3 |
| 4 | — | Apertium | 17.0 |
| 5 | 538 | Z-AGI Labs | 7.6 |
| 6 | 537 | Z-AGI Labs | 6.4 |
| 7 | 597 | CycleL | 0.1 |
| **Closed submission** | | | |
| 1 | 580 | imaxin | 17.6 |
| 2 | — | Apertium | 17.0 |

Table 4: BLEU scores computed over the FLORES+ devtest set for Spanish–Asturian.

As already mentioned in the introduction, the three language pairs have an Apertium MT system

available; we therefore include Apertium among the systems evalauted in this section. The specific versions of Apertium used for each language are: Spanish–Aragonese 0.6.0,[41] Spanish–Aranese 1.0.8,[42] Spanish–Asturian 1.1.1.[43]

Tables 2, 3 and 4 show the BLEU scores attained by each system for Aragonese, Aranese and Asturian, respectively. Similarly, tables 5, 6 and 7 show the results obtained with chrF2. Each table reports, in addition to the BLEU or chrF2 scores, a ranking of the systems from best (#1) to worse. This ranking was derived using a statistical significance test conducted through paired approximate randomization (Riezler and Maxwell, 2005) with SacreBLEU.[44] Systems within the same rank do not exhibit statistically significant differences.

The ranking process involved an iterative approach. We began by selecting the best system for each metric as the control translation. The translations provided by other systems were then compared to this control to determine if the differences were statistically significant. Systems whose output did not differ significantly from the control were associated with it and removed from the pool of translations. The next best system then became the new control translation, and the process was repeated. This iterative process continued until no system remained in the pool.

In the **Spanish–Aragonese** translation task (Tables 2 and 5), the *open* submission results show the ILENIA-MT team achieving the highest BLEU score of 62.7. Notably, ILENIA-MT's performance is consistent, as their second submission scores nearly identical at 62.6. The Apertium baseline and the submission by Vicomtech closely follow, with BLEU scores of 61.1 and 61.0, respectively. UAlacant's entries, which rank 3rd and 4th with scores of 60.2 and 59.8, demonstrate strong competitiveness as well, outperforming the submissions from Helsinki-NLP, which rank lower.

In the *constrained* submission category, SJTU-MT and HW-TSC lead with BLEU scores of 63.2 and 63.0, respectively, surpassing the top scores from the *open* submissions —a difference that is statistically significant—. SJTU-MT's approach

---

| Rank | id | Team | chrF2 |
|---|---|---|---|
| **Open submission** | | | |
| 1 | 636 | ILENIA-MT | 80.0 |
| | 637 | ILENIA-MT | 80.0 |
| 2 | — | Apertium | 79.3 |
| | 633 | Vicomtech | 79.3 |
| 3 | 554 | UAlacant | 78.9 |
| 4 | 495 | UAlacant | 78.8 |
| 5 | 577 | Helsinki-NLP | 75.9 |
| 6 | 549 | Helsinki-NLP | 75.6 |
| 7 | 563 | Helsinki-NLP | 75.4 |
| 8 | 523 | Helsinki-NLP | 74.6 |
| 9 | 548 | UAlacant | 67.5 |
| 10 | 647 | CUNI-GA | 67.8 |
| **Constrained submission** | | | |
| 1 | 607 | HW-TSC | 80.3 |
| | 663 | SJTU-MT | 80.1 |
| | 529 | ILENIA-MT | 79.9 |
| | 558 | ILENIA-MT | 79.9 |
| | 634 | ILENIA-MT | 79.9 |
| 2 | 526 | ILENIA-MT | 79.5 |
| | 642 | TIM-UNIGE | 79.5 |
| | — | Apertium | 79.3 |
| | 504 | Vicomtech | 79.3 |
| 3 | 586 | TIM-UNIGE | 79.0 |
| | 539 | TIM-UNIGE | 79.0 |
| | 644 | TIM-UNIGE | 79.0 |
| | 649 | Stevens | 78.7 |
| 4 | 584 | TAN-IBE | 78.1 |
| 5 | 613 | Stevens | 77.2 |
| 6 | 622 | TRIBBLE | 73.6 |
| 7 | 530 | LCT-LAP | 68.6 |
| 8 | 513 | Stevens | 67.4 |
| 9 | 662 | Stevens | 62.0 |
| | 534 | Z-AGI Labs | 61.8 |
| 10 | 533 | Z-AGI Labs | 60.6 |
| 11 | 507 | SRPH-LIT | 58.4 |
| 12 | 591 | CycleL | 13.7 |

Table 5: chrF2 computed over the FLORES+ devtest set for Spanish–Aragonese.

| Rank | id | Team | chrF2 |
|---|---|---|---|
| **Open submission** | | | |
| 1 | — | Apertium | 49.4 |
| | 627 | Vicomtech | 49.4 |
| 2 | 555 | UAlacant | 49.3 |
| 3 | 587 | ILENIA-MT | 48.8 |
| | 656 | CUNI-GA | 48.5 |
| | 552 | UAlacant | 48.3 |
| 4 | 578 | Helsinki-NLP | 46.6 |
| 5 | 562 | Helsinki-NLP | 45.7 |
| 6 | 550 | Helsinki-NLP | 45.1 |
| 7 | 524 | Helsinki-NLP | 45.0 |
| **Constrained submission** | | | |
| 1 | 527 | ILENIA-MT | 50.1 |
| | 621 | SJTU-MT | 49.9 |
| | 619 | TIM-UNIGE | 49.8 |
| | 617 | TIM-UNIGE | 49.7 |
| | 625 | Vicomtech | 49.8 |
| 2 | 641 | TIM-UNIGE | 49.6 |
| | — | Apertium | 49.4 |
| | 640 | TIM-UNIGE | 49.3 |
| | 575 | TIM-UNIGE | 49.2 |
| 3 | 494 | TIM-UNIGE | 48.8 |
| | 610 | TAN-IBE | 48.8 |
| 4 | 608 | HW-TSC | 47.9 |
| 5 | 623 | TRIBBLE | 46.1 |
| 6 | 531 | LCT-LAP | 45.5 |
| 7 | 581 | SRPH-LIT | 34.8 |
| 8 | 536 | Z-AGI Labs | 32.8 |
| 9 | 535 | Z-AGI Labs | 31.9 |

Table 6: chrF2 computed over the FLORES+ devtest set for Spanish–Aranese.

stands out as one of the most innovative in the task, employing strategies that diverge significantly from traditional methods, whereas HW-TSC used the largest number of layers in the encoder (25) of all the submissions. ILENIA-MT continues to perform strongly, with their top entry scoring 62.3, closely followed by TIM-UNIGE, Apertium and Vicomtech.

For the **Spanish–Aranese** pair (Tables 3 and 6),

the *open* submission category presents a narrower range of BLEU scores compared to Aragonese. Apertium and Vicomtech share the top position with a BLEU score of 28.8, closely followed by UAlacant with 28.5.

In the *constrained submission* category, SJTU-MT once again leads, achieving a BLEU score of 30.4. This time, several other teams — TIM-UNIGE, ILENIA-MT, and Vicomtech— join SJTU-MT at the top, all with scores outperforming the *open* submissions by a statistically significant margin.

The lower BLEU scores for the **Spanish–Asturian** language pair (Tables 4 and 7), are likely due to the way the FLORES+ dev and devtest datasets were constructed, with translations originating from English rather than Spanish. In

| Rank | id | Team | chrF2 |
|---|---|---|---|
| **Open submission** | | | |
| 1 | 576 | SJTU-MT | 55.2 |
| | 551 | Helsinki-NLP | 51.6 |
| | 564 | Helsinki-NLP | 51.6 |
| 2 | 579 | Helsinki-NLP | 51.5 |
| | 568 | TAN-IBE | 51.6 |
| | 629 | Vicomtech | 51.6 |
| 3 | 525 | Helsinki-NLP | 51.4 |
| | 556 | UAlacant | 50.9 |
| | 497 | UAlacant | 50.9 |
| 4 | — | Apertium | 50.8 |
| | 553 | UAlacant | 50.7 |
| 5 | 632 | ILENIA-MT | 50.5 |
| 6 | 648 | CUNI-GA | 48.9 |
| **Constrained submission** | | | |
| 1 | 606 | HW-TSC | 52.2 |
| | 528 | ILENIA-MT | 52.1 |
| | 547 | Mora translate | 51.4 |
| | 630 | Vicomtech | 51.2 |
| | 519 | Mora translate | 51.2 |
| 2 | 590 | Mora translate | 51.2 |
| | 522 | Mora translate | 51.0 |
| | 512 | Mora translate | 51.0 |
| | 543 | Mora translate | 51.0 |
| | — | Apertium | 50.8 |
| | 532 | LCT-LAP | 50.7 |
| 3 | 624 | TRIBBLE | 50.5 |
| | 546 | SRPH-LIT | 50.0 |
| 4 | 538 | Z-AGI Labs | 44.4 |
| 5 | 537 | Z-AGI Labs | 42.7 |
| 6 | 597 | CycleL | 15.9 |
| **Closed submission** | | | |
| 1 | 580 | imaxin | 51.2 |
| 2 | — | Apertium | 50.8 |

Table 7: chrF2 computed over the FLORES+ devtest set for Spanish–Asturian.

the *open* submission category, SJTU-MT leads with a score of 23.2, significantly outperforming the second-best system, Helsinki-NLP, by 5 BLEU points, with the latter's scores clustering around 18.0.

The *constrained* submission results show HW-TSC leading with a BLEU score of 19.8, followed by ILENIA-MT at 18.4. Despite the constrained environment, HW-TSC's results indicate that their extensive use of encoder layers and synthetic data generation proved beneficial.

## 7 Conclusions

This paper has presented the outcomes of the Ninth Conference on Machine Translation (WMT24) Shared Task on Translation into Low-Resource Languages of Spain. The challenge centred on building MT systems for three Romance language pairs: Spanish–Aragonese, Spanish–Aranese, and Spanish–Asturian. In total, 17 teams took part in this shared task.

Across all three language pairs, there is some variability in performance both between and within the categories (open, constrained, and closed). Top-performing teams such as SJTU-MT, ILENIA-MT, and HW-TSC consistently achieved high rankings across multiple pairs. The results also underscore the challenges posed by low-resource languages, where factors such as data availability and the choice of methods —e.g., synthetic data generation, fine-tuning strategies, or transformer model size— significantly affect performance. Notably, most of the best-performing systems utilized the Apertium rule-based system to generate synthetic data, highlighting the ongoing relevance of these approaches in complementing neural methods.

## Acknowledgements

## References

Mikko Aulamo, Sami Virpioja, Yves Scherrer, and Jörg Tiedemann. 2021. Boosting neural machine translation from Finnish to Northern Sámi with rule-based backtranslation. In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 351–356, Reykjavik, Iceland (Online). Linköping University Electronic Press, Sweden.

Martin Bär, Elisa Forcada Rodríguez, and María García-Abadillo Velasco. 2024. Robustness of fine-tuned

LLMs for machine translation with varying noise levels: Insights for Asturian, Aragonese and Aranese. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Nikolay Bogoychev, Jelmer van der Linde, Graeme Nail, Barry Haddow, Jaume Zaragoza-Bernabeu, Gema Ramírez-Sánchez, Lukas Weymann, Tudor Nicolae Mateiu, Jindřich Helcl, and Mikko Aulamo. 2023. Opuscleaner and opustrainer, open source toolkits for training machine translation and large language models.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Jeff Wang, and NLLB Team. 2024. Scaling neural machine translation to 200 languages. *Nature*, 630(8018):841–846.

Ona de Gibert, Mikko Aulamo, Yves Scherrer, and Jörg Tiedemann. 2024. Hybrid distillation from RBMT and NMT: Helsinki-NLP's submission to the Shared Task on Translation into Low-Resource Languages of Spain. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Sören Dréano, Derek Molloy, and Noel Murphy. 2024. Exploration of the CycleGN framework for low-resource languages. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Abhimanyu et al. Dubey. 2024. The Llama 3 herd of models.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Mikel L. Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M. Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine Translation*, 25(2):127–144.

Aarón Galiano-Jiménez, Víctor M Sánchez-Cartagena, Juan Antonio Pérez-Ortiz, and Felipe Sánchez-Martínez. 2024. Universitat d'Alacant's submission to the WMT 2024 Shared Task on Translating into Low-Resource Languages of Spain. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024a. Idiomata cognitor.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024b. Pan-Iberian Language Archival Resource.

Generalitat de Catalunya. 2019. *Els usos lingüístics de la població de l'Aran: Principals resultats de l'Enquesta d'usos lingüístics de la població. 2018*. Generalitat de Catalunya, Barcelona.

Sofía García González. 2024. Enhanced Apertium system: Translation into low-resource languages of Spain Spanish–Asturian. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc'Aurelio Ranzato. 2019. The FLORES evaluation datasets for low-resource machine translation: Nepali–English and Sinhala–English. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6098–6111, Hong Kong, China. Association for Computational Linguistics.

Miroslav Hrabal, Josef Jon, Martin Popel, Nam H Luu, Danil Semin, and Ondřej Bojar. 2024. CUNI at WMT24 general translation task: LLMs, (Q)LoRA,

CPO and model merging. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Tianxiang Hu, Haoxiang Sun, Ruize Gao, Jialong Tang, Pei Zhang, Baosong Yang, and Rui Wang. 2024. SJTU system description for the WMT24 Low-Resource Languages of Spain task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Igor Kuzmin, Piotr Przybyła, Euan McGill, and Horacio Saggion. 2024. TRIBBLE - TRanslating IBerian languages Based on Limited E-resources: System description. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

J. Lin. 1991. Divergence measures based on the shannon entropy. *IEEE Transactions on Information Theory*, 37(1):145–151.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pretraining for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Francisco J. Llera Ramo. 2018. *III Estudio Sociolingüístico de Asturias 2017: Avance de resultados*. Academia de la Llingua Asturiana, Uvieu. Estaya Sociollingüística, colección 7.

Yuanchang Luo, Zhanglin Wu, Daimeng Wei, Hengchao Shang, Zongyao Li, Jiaxin Guo, and Zhiqiang et al. Rao. 2024. Multilingual transfer and domain adaptation for low-resource languages of Spain. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Velayuthan Menan, Dilith Jayakody, Nisansa de Silva, Aloka Fernando, and Surangika Ranatunga. 2024. Back to the stats: Rescuing low resource neural machine translation with statistical methods. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.

Jonathan Mutal and Lucía Ormaechea. 2024. TIM-UNIGE translation into low-resource languages of Spain for WMT24. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Antoni Oliver. 2024. TAN-IBE participation in the Shared Task: Translation into Low-Resource Languages of Spain. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Víctor M. Sánchez-Cartagena, Miquel Esplà-Gomis, Aarón Galiano-Jiménez, Antoni Oliver, Claudi Aventín-Boya, Cristina Valdés, Alejandro Pardos, and Juan Pablo Martínez. 2024. FLORES+ datasets for Aragonese, Aranese, Asturian and Valencian. In *Proceedings of the Ninth Conference on Machine Translation*, pages 00–00, Miami. Association for Computational Linguistics.

David Ponce, Harritxu Gete, and Thierry Etchegoyhen. 2024. Vicomtech@WMT 2024: Shared Task on Translation into Low-Resource Languages of Spain. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Anchel Reyes, Chabier Gimeno, Miguel Montañés, Natxo Sorolla, Pep Espluga, and Juan Pablo Martínez. 2017. *L'aragonés y lo catalán en l'actualidat: analisi d'o censo de población y viviendas de 2011*. Seminario Aragonés de Sociolingüística, Asociación Aragonesa de Sociolochía, Universidad de Zaragoza. Primera parte, febrero 2017.

Stefan Riezler and John T. Maxwell. 2005. On some pitfalls in automatic evaluation and significance testing for MT. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 57–64, Ann Arbor, Michigan. Association for Computational Linguistics.

Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.

Aleix Sant, Daniel Bardanca Outeiriño, José Ramom Pichel Campos, Francesca De Luca Fornaciari, Carlos Escolano, Javier García Gilabert, Pablo Gamallo Otero, Audrey Mash, Xixian Liao, and Maite Melero. 2024. Training and fine-tuning NMT models for low-resource languages using Apertium-based synthetic corpora. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Jörg Tiedemann, Mikko Aulamo, Daria Bakshandaeva, Michele Boggia, Stig-Arne Grönroos, Tommi Nieminen, Alessandro Raganato, Yves Scherrer, Raúl Vázquez, and Sami Virpioja. 2024. Democratizing neural machine translation with opus-mt. *Language Resources and Evaluation*, 58(2):713–755.

Dan John Velasco, Manuel Antonio Rufino, and Jan Christian Blaise Cruz. 2024. Samsung R&D Institute Philippines @ WMT 2024 Low-Resource Languages of Spain Shared Task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 55204–55224. PMLR.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

# Findings of the WMT 2024 Shared Task on Discourse-Level Literary Translation

**Longyue Wang, Siyou Liu, Chenyang Lyu, Wenxiang Jiao, Xing Wang, Jiahao Xu,
Zhaopeng Tu, Yan Gu, Weiyu Chen, Minghao Wu, Liting Zhou,
Philipp Koehn, Andy Way, Yulin Yuan**

vincentwang0229@gmail.com

## Abstract

Following last year's WMT, we (Tencent AI Lab and China Literature Ltd.) have continued to host literary translation shared task (Wang et al., 2023) this year, the second edition of the *Discourse-Level Literary Translation*.

First, we (Tencent AI Lab and China Literature Ltd.) release a copyrighted and document-level Chinese-English web novel corpus. Furthermore, we put forth an industry-endorsed criteria to guide human evaluation process. This year, we totally received 10 submissions from 5 academia and industry teams. We employ both automatic and human evaluations to measure the performance of the submitted systems. The official ranking of the systems is based on the overall human judgments. In addition, our extensive analysis reveals a series of interesting findings on literary and discourse-aware MT. We release data, system outputs, and leaderboard at https://www2.statmt.org/wmt24/literary-translation-task.html.

## 1 Introduction

With the swift progression of MT and the notable advancements in Large Language Models (LLM) (**??**), our curiosity is piqued regarding the efficacy of MT and LLM in the realm of literary translation. We aim to explore the extent to which these technologies can aid in addressing the intricate challenges of translating literary works. Therefore, we hold the first edition of the *Discourse-Level Literary Translation* in WMT 2023. Literary texts encompass a wide range of forms, including novels, short stories, poetry, plays, essays, and more. Among these, *web novels*, also known as online or internet novels, represent a unique and rapidly growing subset of literature. Their popularity, accessibility, and diverse genres set them apart. As they provide not only an extensive volume of text but also exhibit distinctive linguistic features, cultural phenomena, and simulations of societies, web novels can serve as valuable resources and challenging for MT research.

## Limitations

We discuss the potential limitations of this edition of the shared task as follows:

- *Language Pair*. This year, we only focus on Chinese→English direction. However, we have a long-term plan to continuously organize this task, and will extend the copyrighted dataset into Chinese-Russian and Chinese-German language pairs next year.

- *Literary Genre*. This year, we mainly used the Web Fiction Corpus which is only one type of literary text. We use Web Fiction for two reasons: (1) its literariness is less complicated than others (e.g. poetry, masterpiece); (2) such bilingual data are numerous and continuously increased. We will consider to extend more literary genres such as poetric translation in the next year.

- *Discourse Benchmark*. We have accumulated some discourse- and context-aware benchmarks (**???**). These benchmarks are pivotal for assessing the proficiency of LLMs in handling complex language structures and contextual nuances. As participation of LLM-based systems in our shared tasks increases, we anticipate integrating these benchmarks more comprehensively into our future evaluations to better measure and understand the evolution of LLM capabilities in linguistic context and discourse comprehension.

## References

Longyue Wang, Zhaopeng Tu, Yan Gu, Siyou Liu, Dian Yu, Qingsong Ma, Chenyang Lyu, Liting Zhou, Chao-Hong Liu, Yufeng Ma, Weiyu Chen, Yvette Graham, Bonnie Webber, Philipp Koehn, Andy Way, Yulin Yuan, and Shuming Shi. 2023. Findings of the WMT 2023 shared task on discourse-level literary translation: A fresh orb in the cosmos of LLMs. In *Proceedings of the Eighth Conference on Machine*

| Type | System | Sent-Level | | | | Doc-Level |
|------|--------|------|------|------|------|------|
| | | **BLEU**$^\uparrow$ | **chrF2**$^\uparrow$ | **COMET**$^\uparrow$ | **TER**$^\downarrow$ | **d-BLEU**$^\uparrow$ |
| *Baselines* | Google$^\star$ | 37.4 | 57.0 | 80.50 | 57.4 | 47.3 |
| | Llama-MT$^\star$ | n/a | n/a | n/a | n/a | 43.1 |
| | GPT-4$^\star$ | n/a | n/a | n/a | n/a | 43.7 |
| *Primary* | Cloudsheep$^\star$ | 39.5 | 57.5 | 81.22 | 55.5 | 48.5 |
| | HW-TSC | 40.5 | 58.5 | 82.61 | 56.0 | 50.2 |
| | NLP2CT-UM$^\star$ | **41.6** | **58.7** | **83.56** | **52.7** | **50.9** |
| | NTU$^\star$ | 20.9 | 41.9 | 74.53 | 73.9 | 34.6 |
| | SJTU-LoveFiction$^\star$ | 35.1 | 54.7 | 80.79 | 62.1 | 47.2 |
| *Contrastive* | HW-TSC | 40.6 | 58.6 | 82.59 | 55.9 | 50.3 |
| | NLP2CT-UM$_1^\star$ | 41.6 | 58.7 | 83.54 | 52.8 | 50.8 |
| | NLP2CT-UM$_2^\star$ | 41.5 | 58.6 | 83.38 | 52.8 | 50.7 |
| | SJTU-LoveFiction$_1^\star$ | 35.7 | 56.0 | 82.67 | 59.7 | 46.3 |
| | SJTU-LoveFiction$_2^\star$ | 38.6 | 56.5 | 82.49 | 57.1 | 49.6 |

Table 1: Evaluation results of baseline and participants' systems in terms of **automatic evaluation methods**, including 1) sentence-level metrics BLEU, chrF, COMET, TER; and 2) document-level metrics d-BLEU. Systems marked with $^\star$ are unconstrained, while others are constrained. The COMET is calculated with *unbabel-comet* using *Reference 1* while others are calculated with *sacrebleu* using two references. The best primary constrained and unconstrained systems are highlighted.

# A   Example Appendix

This is a section in the appendix.

# Findings of the WMT 2024 Shared Task on Chat Translation

**Wafaa Mohammed**[1]    **Sweta Agrawal**[2]    **M. Amin Farajian**[3]
**Vera Cabarrão**[3]    **Bryan Eikema**[1]    **Ana C. Farinha**[3]    **José G. C. de Souza**[3]

[1]University of Amsterdam, Netherlands
[2]Instituto de Telecomunicações, Lisbon, Portugal
[3]Unbabel, Lisbon, Portugal

## Abstract

This paper presents the findings from the third edition of the Chat Translation Shared Task. As with previous editions, the task involved translating bilingual customer support conversations, specifically focusing on the impact of conversation context in translation quality and evaluation. We also include two new language pairs: English↔Korean and English↔Dutch, in addition to the set of language pairs from previous editions: English↔German, English↔French, and English↔Brazilian Portuguese.

We received 22 primary submissions and 32 contrastive submissions from eight teams, with each language pair having participation from at least three teams. We evaluated the systems comprehensively using both automatic metrics and human judgments via a direct assessment framework. The official rankings for each language pair were determined based on human evaluation scores, considering performance in both translation directions—agent and customer. Our analysis shows that while the systems excelled at translating individual turns, there is room for improvement in overall conversation-level translation quality.

## 1 Introduction

Translating conversational text, in particular customer support chats, is an important and challenging application for machine translation (MT) technology. According to a 2020 survey from CSA Research, 75% of shoppers are more likely to make another purchase if customer support is offered in their native language, making it appealing for businesses to invest in multilingual support.[1] However, there are several key challenges to translating chats: customer support chats typically feature short text exchanges between agents and customers (see Table 1), leading to fragmented sentences and

---

[1] https://csa-research.com/Featured-Content/For-Global-Enterprises/Global-Growth/CRWB-Series/CRWB-B2C

omission of information (implied by the context). This makes it difficult for MT systems to produce coherent translations that maintain the intended meaning of the text (Farajian et al., 2020). Furthermore, chats often use colloquial language and are characterized by informality and grammatical inaccuracies (Gonçalves et al., 2022). Consequently, translating such content poses a dual challenge: not only must a system accurately translate between languages, but it should also effectively model the nuances and ambiguity in a dialogue.

While recent advancements in MT systems, driven by LLMs, have proven effective in various tasks, bilingual chat translation remains underexplored. The **Chat Translation Shared Task** aims to bridge this gap by promoting research and development of MT systems designed specifically for conversational translation. This year's edition places special emphasis on the role of conversation context, encouraging teams to examine how context influences translation in the inherently ambiguous and dynamic nature of chat interactions. Following the success of the previous two editions of the Chat Translation Shared Task (Farajian et al., 2020; Farinha et al., 2022), this year we organized the third edition of the task with the following improvements:

- We expanded the set of language pairs to include English↔Korean (EN-KO) and English↔Dutch (EN-NL), in addition to languages from previous editions: English↔German (EN-DE), English↔French (EN-FR), and English↔Brazilian Portuguese (EN-PT).

- We carefully curated the evaluation sets to enable the evaluation of effective context utilization on systems' performance.

- We conducted a comprehensive evaluation of all systems using: a) automatic metrics (both neural and lexical) that assess translation quality and the accuracy of modeling discourse phenomena

701

| | |
|---|---|
| 🧑 customer | Hallo, ich komme nicht in meine Sum up pos was denn no App rein |
| | Hello, I can not get into my sum up pos what then no app |
| 🎧 agent | I am sorry to hear that. |
| | Es tut mir leid, das zu erfahren. |
| 🎧 agent | Let me see what I can do for you |
| | Lassen Sie mich sehen, was ich für Sie tun kann. |
| 🎧 agent | Could you please tell me what error message you can see while logging in to your POS? |
| | Könnten Sie mir bitte sagen, welche Fehlermeldung Sie sehen können, während Sie sich bei Ihrem POS anmelden? |
| 🧑 customer | Wenn ich auf die App gehe, erscheint dieses Gerät hinzufügen. |
| | When I go to the app, it shows Add this device. |
| 🎧 agent | Could you please try to connect the App with the POS? |
| | Könnten Sie bitte versuchen, die App mit dem POS zu verbinden? |
| 🧑 customer | die App ist die PRS-ORG pos app |
| | the app is the PRS-ORG app |
| 🧑 customer | ich habe die Frage daher nicht verstanden |
| | so I did not understand the question |
| 🎧 agent | Could you please elaborate on your query? |
| | Könnten Sie bitte Ihre Anfrage näher erläutern? |

Table 1: An example of a EN-DE conversation between a *customer* (🧑) and an *agent* (🎧) from MAIA dataset.

using MUDA (Fernandes et al., 2023b), b) human direct assessments by professional linguists, and c) LLM-based fine-grained error analysis following the MQM framework.

We received a total of 22 primary submissions, 6 submissions for en↔de, 5 for en↔fr, 4 for en↔nl, 4 for en↔pt-br, and 3 for en↔ko. Six out of the eight teams used large language models (LLMs) as their base translation model, implementing various strategies such as finetuning on shared task data, augmenting training data with synthetic datasets, prompting strategies, quality-aware decoding, and several ways of leveraging conversational context to improve translation quality. With these multi-faceted solutions explored by several teams, this year's shared task yields valuable insights into the effectiveness of LLMs in translating conversational texts. We summarize the key findings from the shared task below:

- Incorporating contextual information from previous turns almost always improved translation quality. However, the optimal method for introducing context (whether through summary, graph, or raw context) still requires further investigation.

- Human evaluation showed that turn-level translation quality was consistently high across all participating systems and language pairs. Nonetheless, there is room for improvement in translating texts from later turns and at the conversation level as a whole.

- The UNBABEL-IT submission achieved the best results across most language pairs and evaluation

criteria, except on the EN-DE and EN-FR tasks according to automatic metrics.

These findings suggest that future editions of the shared task could benefit from a) designing evaluation frameworks, both automatic and human, that specifically target dialogue-specific criteria to better understand system limitations (Yeh et al., 2021; A, 2022; Deriu et al., 2021); b) expanding the datasets to include more challenging domains (e.g. patient-physician conversation or everyday dialogues) and contexts (e.g. multimodal chats) for a more thorough evaluation of MT systems.

## 2  Task Description

As in previous editions of the task, we evaluate the effectiveness of a translation layer in translating text from the customer's language to the agent's language (e.g., English) and vice versa. We provide real bilingual customer support data for five different language pairs and encourage the participants to use conversation context. They are asked to submit translations for both directions (agent and customer). We detail the shared task dataset provided to the participants and evaluation in § 2.1 and § 2.2 respectively.

### 2.1  Data: The MAIA 2.0 Corpus

The MAIA 2.0 corpus builds upon the dataset released in the previous edition (Farinha et al., 2022) and includes two additional language pairs: Dutch and Korean. Furthermore, we expanded the sizes of the existing language pairs, ensuring that each language pair contained approximately 20k segments. The dataset encompasses dialogues across diverse

| LP | train | | | | dev | | | | test | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | # seg | # conv | # length | # words | # seg | # conv | # length | # words | # seg | # conv | # length | #words |
| EN-NL | 15.5k | 595 | 26.0 | 8.6 | 2.5k | 72 | 35.4 | 9.8 | 2k | 58 | 34.7 | 10.2 |
| EN-PT | 15.0k | 435 | 34.7 | 8.0 | 2.5k | 96 | 26.6 | 8.8 | 2k | 73 | 27.9 | 8.8 |
| EN-DE | 17.8k | 493 | 36.1 | 8.5 | 2.5k | 82 | 31.3 | 9.4 | 2k | 67 | 30.5 | 9.4 |
| EN-KO | 16.1k | 423 | 38.1 | 8.5 | 1.8k | 38 | 50.9 | 10.5 | 2k | 42 | 47.2 | 9.6 |
| EN-FR | 15.0k | 264 | 56.9 | 7.7 | 3.0k | 90 | 33.4 | 10.1 | 2k | 65 | 32.2 | 10.1 |

Table 2: Dataset statistics with the number of segments (#seg), number of conversations (#conv), average conversation length (#length), and average number of words per turn (#words) in each split. Note that for KO customer parts, we considered the English reference translation to calculate the number of words.

topics, including account registration issues, payment and delivery clarifications, and after-sale services in various industries such as retail and gaming. The new dataset was automatically anonymized using Unbabel's proprietary anonymization tool, followed by a manual validation performed by expert linguists, to comply with the General Data Protection Regulation (GDPR). The corpus is released under the CC-BY-NC-4.0 license and can be freely used for research purposes only. Please note that, as the license states, no commercial uses are permitted for this corpus.

**Training and Evaluation Datasets.** We provide both training and evaluation (development and test) sets that participants can use to build their systems. Table 2 presents each data splits' statistics, including the number of segments, conversations, and average conversation length. We construct the development and test sets by selecting conversations that exhibit the highest counts of context-dependent discourse phenomena tags, as extracted using Multilingual Discourse Aware (MUDA) tagger (Fernandes et al., 2023b).

## 2.2 Evaluation

We perform a comprehensive evaluation of all submitted systems, using both automatic and human evaluation. Official rankings are determined based on the human assessment scores for both customer and agent translations. We outline the various evaluations conducted below:

### 2.2.1 Automatic Evaluation

We use COMET (Rei et al., 2022) as our primary evaluation metric for assessing translation quality of the submitted systems.[2] Additionally, we report lexical metrics: BLEU and CHRF using the SacreBLEU library (Post, 2018). We also include CONTEXTCOMETQE (Agrawal et al., 2024),

[2]Unbabel/wmt22-comet-da

a reference-free metric that uses bilingual context (previous two turns) to assess the translation quality of the current turn. As efficient discourse handling is not directly reflected in standard MT metrics (both lexical and neural), we report the F1 accuracy on the MUDA-tagged discourse phenomena. We considered 4 context-dependent discourse phenomena in our analysis:

- **Lexical cohesion:** Entities may have multiple possible translations in the target language, but the same entity should be referred to by the same word in a conversation.

- **Formality:** Korean uses honorifics to indicate formality, which are special titles or words expressing courtesy or respect for position. In other languages, speakers use second-person pronouns to refer to someone more formally or informally, depending on their relationship with the addressee. Formality should be consistent throughout a conversation.

- **Pronoun resolution:** Some highly inflected languages use gendered pronouns based on semantic or morphological rules. To assign the correct pronoun, it is therefore necessary to use the conversation's context to distinguish the grammatical gender of the pronoun's antecedent.

- **Verb forms:** Verbs must be translated consistently using the form that reflects the tone, and mood of both parties in the conversation.

### 2.2.2 Manual Evaluation

We use the DA+SQM (Direct Assessment + Scalar Quality Metric) evaluation framework, following the campaigns conducted by the WMT General Translation track over the past years, implemented via the Appraise framework (Federmann, 2018) to collect human assessments of translation quality

Figure 1: Screen capture of the Appraise interface used by professional linguists to perform human evaluation.

| LP | Threshold | # Chats | # Systems | # annotated segments |
|---|---|---|---|---|
| EN-NL | 35 | 27 | 5 | 3830 |
| EN-PT | 28 | 41 | 5 | 4700 |
| EN-DE | 31 | 36 | 7 | 6629 |
| EN-KO | 48 | 24 | 4 | 3648 |
| EN-FR | 33 | 37 | 6 | 6324 |

Table 3: Statistics of the conversations and instances sampled for the human evaluation step.

for the submitted systems. We ask professional linguists hired via the UpWork[3] platform to evaluate each turn in a conversation within the full context and provide a conversation-level quality score on a continuous scale from 0 to 100. They were instructed to pay special attention to conversation-level properties such as the consistency of style, selection of terms, formality, etc in addition to the correctness criteria. The quality scale includes seven labeled tick marks representing various quality levels based on both accuracy and grammatical correctness (Figure 1).

**Data Selection** For the human evaluation, we retain conversations with up to a given number of turns to make the evaluation manageable. The number of turns for each language pair is specified in Table 3 ("Threshold"), together with the number of conversations and instances retained.

**Measure** We generate turn-level and conversation-level system rankings for each language pair by aggregating the direct assessment scores provided by the linguists at the turn level and the conversation level respectively.

### 2.2.3 LLM-based Error Assessments

LLM-based evaluation has garnered a lot of interest from the community for conducting human-like evaluations. This shift is largely driven by the increasing complexity and scale of language models, making them capable of capturing nuanced understanding and performance of models in real-world tasks. For MT, LLM-based metrics are used to provide fine-grained error assessments over the nature, type, and severity of the errors following the MQM framework (Fernandes et al., 2023a; Lu et al., 2024; Kocmi and Federmann, 2023). Recently, Agrawal et al. (2024) show that context-aware prompting for deriving MQM assessment using LLMs can achieve better correlation with human judgments than the standard MQM prompt for chat translation evaluation, even surpassing COMET.

Hence, we complement our evaluation with an LLM-based fine-grained assessment of MT outputs derived using CONTEXTMQM (Agrawal et al., 2024). The prompt includes the past eight bilingual source sentences as context and one in-domain in-context example with MQM assessment to elicit MQM-like evaluation from GPT-4o-mini[4] for all systems submitted for the EN-DE track.[5] Like MQM, we compute the segment-level *error* score aggregating the number of minor, major, and critical errors, weighted by factors of 10, 5, and 1, respectively.

---

[3]upwork.com

[4]gpt-4o-mini-2024-07-18 accessed on 10-2-2024.

[5]Due to budget constraints, we conduct this evaluation only on EN-DE, which had the highest number (eight) of participating teams.

## 3 Participants

This section provides a brief description of each participant's systems (§ 3.1). Table 4 summarizes details about the team's institutions and the language directions they participated in. Participants were asked to submit up to three systems per language direction: one primary (explicitly marked) and up to two contrastive systems. Next, we discuss the commonalities and differences between the different submissions § 3.2.

### 3.1 Systems

#### 3.1.1 NLLB-3.3B (Baseline)

For our baseline model, we used the NLLB-3.3B multilingual machine translation model (Costajussà et al., 2022) based on an encoder-decoder Transformer architecture (Vaswani et al., 2017). NLLB-3.3B is trained to support over 200 languages, including those of interest in this shared task: English, German, French, Dutch, Brazilian Portuguese, and Korean. We opted for a sentence-level baseline that does not incorporate additional context and used a beam size of 4 for generating translation hypotheses.

#### 3.1.2 UNBABEL-IT

The joint submission of Unbabel and IT includes one primary submission and two contrastive submissions per language pair. The systems are based on Tower-7B models and are trained on the chat datasets released by the shared task. Their primary system uses contextual MBR re-ranking over a set of 50 candidates to get the best hypothesis. Additionally, the first contrastive submission is a 70B variant of the Tower model specialized to have general purpose translation capabilities and the second one uses greedy decoding with the 7B model finetuned on chat datasets.

#### 3.1.3 DEEPTEXT LAB

DEEPTEXT LAB participated in the English-Korean language pair with a single primary system. Their submission leverages Google's Gemma-2-27B model [6], using the most recent two turns and summaries of previous turns as context, all within the same document. The turn summaries are generated using the GPT-4o-mini model. Their system was trained solely using the training data provided by the shared task.

| TEAM | INSTITUTION | DIRECTIONS |
|------|-------------|------------|
| DeepText Lab | Yonsei University | EN-KO |
| HW-TSC | Huawei Translation Service Center | EN-DE |
| Multitan-GML | Université Paris Cité | EN-FR |
| SETU-ADAPT | ADAPT research centre & Dublin City University | EN-DE, EN-FR |
| SheffieldGATE | University of Sheffield | EN-DE, EN-NL, EN-PT |
| CLTeam | Vrije Universiteit Amsterdam | EN-DE, EN-NL, EN-FR, EN-PT |
| DCUGenNLP | Dublin City University | ALL |
| Unbabel-IT | Unbabel & Instituto de Telecomunicações | ALL |
| Baseline | Organizers | ALL |

Table 4: The participating teams, their affiliations, and the language directions that they participated.

#### 3.1.4 HW-TSC

Huawei Translation Service Center (HW-TSC) team submitted a primary and two contrastive systems for English↔German language pair. Their system is a 25-6 transformer encoder-decoder model with a feed-forward dimension of 4096 and 16 self-attention layers. Their primary submission uses a model from the previous edition of the shared task as a baseline, finetuned on this edition's training data, followed by a second finetuning on the validation data. Next, they use MBR reranking to select the optimal candidate with COMET as the utility function using outputs generated from a diverse set of models. Their system then undergoes a self-training step on the MBR output. The contrastive submissions include models trained with different finetuning strategies (e.g. excluding the finetuning on the dev set).

#### 3.1.5 SHEFFIELDGATE

The SHEFFIELDGATE team participated in English↔German, English↔Dutch, and English↔Brazilian Portuguese, with one primary system per language pair. Their system performs low-rank (Hu et al., 2022) instruction-tuning with the training and validation datasets provided by the shared task on the Llama-3-8B-Instruct [7] model. To incorporate contextual information and dependencies between chat messages, they introduce a context-aware sliding window approach that incorporates translations generated at each turn into the prompt.

---

[6]google/gemma-2-27b-it

[7]meta-llama/Meta-Llama-3-8B-Instruct

| Participant | Base Model | Chat Context? | In-domain Training? | Multilingual? | Synthetic Data? | Decoding |
|---|---|---|---|---|---|---|
| DeepText Lab | Gemma-2-27B | ✓(summary) | ✓ | ✗ | ✗ | NR |
| HW-TSC | Transformer 25-6 (from scratch) | ✗ | ✓ | ✗ | ✓ | MBR |
| Multitan-GML | Commercial* | ✓ | ✓ | ✗ | ✗ | NR |
| SETU-ADAPT | Llama-3-8B (EN-DE) | ✓(few-shot) | ✓ | ✗ | ✗ | NR |
| | NLLB-200-600M (EN-FR) | ✗ | ✓ | ✗ | ✓ | NR |
| SheffieldGATE | Llama-8b-Instruct | ✓ | ✓ | ✓ | ✗ | NR |
| CLTeam | TowerInstruct-7B-v0.2 | ✓(graph) | ✗ | ✓ | ✗ | NR |
| DCUGenNLP | Llama3.1-8b | NR | ✓ | ✓ | ✗ | NR |
| Unbabel-IT | TowerBase-7B | ✓ | ✓ | ✓ | ✗ | MBR |
| Baseline | NLLB-3.3B | ✗ | ✗ | ✓ | ✗ | Beam (4) |

Table 5: Summary of approaches for all primary submissions. NR: Not reported.

### 3.1.6 SETU-ADAPT

SETU-ADAPT team submitted 3 (one primary and two contrastive) systems based on different pre-trained models: NLLB[8], MBART-50[9] and Llama-3-8B[10]. Their primary system for EN-DE uses a Llama-3-8B backbone finetuned on the in-domain chat and a synthetic dataset generated by back-translating domain-specific monolingual sentences. For EN-FR, they finetune an NLLB-600M model. During inference, with the LLM-based models, they perform few-shot prompting using examples retrieved via similarity search from the training dataset. Their contrastive systems are based on the encoder-decoder models but use the same datasets for training.

### 3.1.7 MULTITAN-GML

MULTITAN-GML's primary system finetunes a "Dialog" in-domain specialized model hosted on the Model Studio Lite server [11] with 2022 Chat Task (train, valid, test) and 2024 Chat Task (valid) datasets. Their two contrastive submissions use outputs from NLLB-3.3B model and the Deep_translator API respectively. All outputs are post-edited using GPT-4o.

### 3.1.8 DCUGENNLP

DCUGENNLP team submitted a total of 15 systems (one primary and two contrastive) for all the five language pairs. Their primary system finetunes a Llama-3.1-8B model on a mix of the chat task's training data and datasets from other WMT tracks. They also include synthetically generated customer-service data generated using one of their contrastive submission. Other contrastive submis-

sions use Mistral-7B as base models with optional prompt tuning or finetuning of adapter layers.

### 3.1.9 CLTEAM

CLTEAM submitted one primary and one contrastive systems for each of the English↔German, English↔French, English↔Dutch, and English↔Brazilian Portuguese language pairs. Their system uses TowerInstruct-7B-v0.2 [12] model as the base LLM. For their primary submission, they prompt the model with both the dialogue history represented using a graph and the source sequence to be translated. To generate the graph, they prompt GPT-4o to extract entities and relationships from the dialogue data, creating triples from these elements. For the contrastive submission, they prompt the model with only the source sequence to be translated.

## 3.2 Discussion

Table 5 presents a summary of approaches used by all the submitted systems. We highlight some key aspects below:

**Model Architecture** Most teams except CLTEAM and HW-TSC finetuned general-purpose pre-trained LLMs. Where CLTEAM used an off-the-shelf translation-finetuned LLM, HW-TSC opted for a custom bilingual encoder-decoder model for their participation.

**Training Data** All teams used the provided training and development data, sourced from the current and previous versions of the task. HW-TSC went a step further by generating a synthetic parallel corpus. They did this by forward translating source-side monolingual data into target-side text and backtranslating target-side monolingual into source-side texts. SETU-ADAPT similarly used

---

[8] facebook/nllb-200-distilled-600M

[9] facebook/mbart-large-50-many-to-many-mmt

[10] unsloth/llama-3-8b-bnb-4bit

[11] modelstudio-lite

[12] Unbabel/TowerInstruct-7B-v0.2

| SYSTEM | EN-DE | | EN-FR | | EN-NL | | EN-PT | | EN-KO | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | EN | FR | EN | NL | EN | PT | EN | KO | EN |
| DeepText Lab | | | | | | | | | 93.03 | 94.11 |
| HW-TSC | **93.58** | **93.30** | | | | | | | | |
| MULTITAN-GML | | | 90.09 | 92.42 | | | | | | |
| ADAPT | 90.59 | 90.97 | 82.19 | 82.69 | | | | | | |
| SheffieldGATE | 88.67 | 90.10 | | | 88.93 | 89.71 | 90.05 | 88.12 | | |
| CLTeam | 90.90 | 91.63 | 91.37 | 91.90 | 91.31 | 91.22 | 91.77 | 90.12 | | |
| DCUGenNLP | 90.49 | 91.10 | 91.05 | 90.73 | 91.32 | 90.96 | 93.24 | 89.66 | 91.50 | 93.41 |
| Unbabel-IT | 93.22 | 92.48 | **92.96** | **92.71** | **94.36** | **93.38** | **94.76** | **92.46** | **94.96** | **95.16** |
| NLLB-3.3B | 90.56 | 89.03 | 91.06 | 89.18 | 87.86 | 88.45 | 86.33 | 86.10 | 87.26 | 88.05 |
| Δ (Best) | +3.02 | +4.27 | +1.9 | + 3.53 | +6.50 | +4.93 | +8.43 | +6.36 | +7.70 | +7.11 |

Table 6: COMET results on the official test set. Δ (Best): improvement over baseline.

| SYSTEM | EN-DE | | EN-FR | | EN-NL | | EN-PT | | EN-KO | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | EN | FR | EN | NL | EN | PT | EN | KO | EN |
| DeepText Lab | | | | | | | | | 57.67 | 77.96 |
| HW-TSC | **82.66** | **84.03** | | | | | | | | |
| MULTITAN-GML | | | 79.54 | **82.71** | | | | | | |
| ADAPT | 69.50 | 76.63 | 63.92 | 55.98 | | | | | | |
| SheffieldGATE | 64.94 | 72.04 | | | 60.01 | 68.31 | 67.67 | 66.38 | | |
| CLTeam | 69.87 | 75.39 | 74.66 | 77.41 | 63.59 | 73.00 | 71.38 | 69.45 | | |
| DCUGenNLP | 69.84 | 73.64 | 73.73 | 73.78 | 67.44 | 70.47 | 75.24 | 67.27 | 49.02 | 75.35 |
| Unbabel-IT | 77.23 | 79.87 | **80.51** | 78.57 | **80.25** | **78.60** | **82.55** | **76.01** | **62.29** | **81.57** |
| NLLB-3.3B | 70.22 | 71.79 | 76.03 | 76.37 | 59.55 | 68.62 | 58.60 | 67.13 | 34.50 | 69.87 |
| Δ (Best) | +12.44 | +12.24 | +4.48 | +6.34 | +20.70 | +9.98 | +23.95 | +8.88 | +27.79 | +11.70 |

Table 7: CHRF results on the official test set. Δ (Best): improvement over baseline.

back translation to generate more in-domain data for their EN-FR submission.

**Inference** Both UNBABEL-IT and HW-TSC leveraged a quality-aware decoding (QAD) approach (Fernandes et al., 2022) for further improving the quality of outputs during inference. While HW-TSC optimized for COMET, UNBABEL-IT used a context-aware COMET metric as a utility for selecting the best candidate. HW-TSC also used MBR outputs to further finetune the model.

**Context Usage** Different strategies were employed to incorporate conversation context into the translation process. UNBABEL-IT, SHEFFIELDGATE, and MULTITAN-GML utilized the previous turns of the conversation as context to maintain continuity and coherence in translations. DEEPTEXT LAB used both the previous two turns as well as the summary of all the previous conversation turns except the last

two, generated by GPT-4o-mini. This allowed the model to focus on the essential part of the previous content without being overwhelmed by excessive details. On the other hand, CLTEAM used a graph representation of the conversation's history as context, capturing the connectivity between various concepts thus serving as a compressed memory of the dialogue context. SETU-ADAPT used few shot examples extracted from the training data using sentence-embedding similarity.

All teams that participated for more than one language pair opted for a multilingual system except for SETU-ADAPT team who submitted two different systems for each language pair they participated in (EN-DE, EN-FR).

## 4 Overall Results

We present the results of the automatic evaluation for all participating systems for all language pairs

| SYSTEM | EN-DE | | EN-FR | | EN-NL | | EN-PT | | EN-KO | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | EN | FR | EN | NL | EN | PT | EN | KO | EN |
| DeepText Lab | | | | | | | | | 15.99 | 16.15 |
| HW-TSC | 20.79 | 23.37 | | | | | | | | |
| MULTITAN-GML | | | 0.31 | 0.21 | | | | | | |
| ADAPT | 15.63 | 17.97 | -23.31 | -22.88 | | | | | | |
| SheffieldGATE | 17.87 | 17.54 | | | 13.72 | 14.39 | 5.87 | 3.46 | | |
| CLTeam | 19.65 | 21.15 | 8.22 | 7.26 | 19.00 | 19.20 | 8.64 | 7.68 | | |
| DCUGenNLP | 17.27 | 20.38 | 5.11 | 4.80 | 16.55 | 16.09 | 8.69 | 6.70 | 15.84 | 15.74 |
| Unbabel-IT | **24.41** | **26.15** | **10.67** | **10.00** | **23.93** | **23.39** | **12.74** | **10.59** | **21.64** | **21.08** |
| NLLB-3.3B | 15.56 | 19.09 | 1.24 | 0.77 | 9.35 | 8.04 | -5.51 | -6.75 | 4.11 | 4.13 |
| Δ (Best) | +8.85 | +7.06 | +9.43 | +9.23 | +14.58 | +15.35 | +18.25 | +17.34 | +17.53 | +16.95 |

Table 8: CONTEXTCOMETQE results on the official test set. Δ (Best): improvement over baseline.

| SYSTEM | EN-DE | | EN-FR | | EN-NL | | EN-PT | | EN-KO | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | EN | FR | EN | NL | EN | PT | EN | KO | EN |
| DeepText Lab | | | | | | | | | 37.65 | 66.98 |
| HW-TSC | **68.76** | **71.27** | | | | | | | | |
| MULTITAN-GML | | | 65.43 | **71.80** | | | | | | |
| ADAPT | 51.39 | 59.90 | 33.17 | 28.56 | | | | | | |
| SheffieldGATE | 41.15 | 50.72 | | | 33.62 | 46.54 | 42.58 | 42.25 | | |
| CLTeam | 50.41 | 55.71 | 57.05 | 61.09 | 39.29 | 55.41 | 46.42 | 49.34 | | |
| DCUGenNLP | 49.97 | 57.29 | 56.32 | 56.39 | 46.38 | 52.15 | 56.36 | 45.87 | 27.66 | 62.13 |
| Unbabel-IT | 61.45 | 62.86 | **66.41** | 63.18 | **65.70** | **63.75** | **67.86** | **59.05** | **41.54** | **71.01** |
| NLLB-3.3B | 50.43 | 52.09 | 59.21 | 58.07 | 33.55 | 48.47 | 28.25 | 45.59 | 12.46 | 49.76 |
| Δ (Best) | +18.33 | +19.18 | +7.20 | +13.73 | +32.15 | +15.28 | +39.61 | +13.46 | +29.08 | +21.25 |

Table 9: BLEU results on the official test set. Δ (Best): improvement over baseline.

in § 4.1. We then discuss findings from human evaluation in § 4.2, followed by an LLM-based error assessment of submitted systems for the EN-DE task in § 4.3.

## 4.1 Automatic Evaluation

Tables 6-9 show the results of automatic evaluations on the official test set using COMET, CHRF CONTEXTCOMETQE and BLEU respectively -– most participant systems improve the translation quality according to both neural (COMET, CONTEXTCOMETQE) and lexical (CHRF, BLEU) metrics over the NLLB-3.3B model, except the SETU-ADAPT system for EN-FR. This can be explained by the fact that SETU-ADAPT finetunes an NLLB-600M model for EN-FR, which, albeit from the same family of models as our baseline (NLLB-3.3B), is significantly smaller in size.

The UNBABEL-IT submission consistently outperforms all other systems, except the EN-DE translation task, where the winning submission according to COMET, BLEU, and CHRF is HW-TSC. Similarly, MULTITAN-GML scores the best on BLEU and CHRF when translating French into English. Interestingly both systems (UNBABEL-

IT and HW-TSC) use MBR decoding with CONTEXTCOMET and COMET respectively, suggesting that inference optimization techniques like quality-aware decoding methods (Fernandes et al., 2022) can be useful in pushing the translation quality of strong MT systems. However, as we will see in §4.2, this difference is not reflected in human assessments and in automatic metrics (CONTEXTMQM and CONTEXTCOMETQE), with different methods scoring the two systems differently. This highlights the importance of carefully selecting the optimized metrics and the evaluation criteria, as over-optimizing certain metrics may lead to mixed or misleading outcomes (Fernandes et al., 2022).

UNBABEL-IT's submission also achieves the highest scores across all settings according to CONTEXTCOMETQE. However, we observe that the range of quality scores produced by the CONTEXTCOMETQE model, when aggregated at the system level, significantly deviates from the typical range of this metric.[13] While Agrawal et al. (2024) demonstrate its effectiveness as a segment-level

[13]System-level scores are higher when the context is not considered.

metric with improved correlation to human judgments, further investigation is necessary to understand how these system-level scores should be interpreted. For instance, MULTITAN-GML, which performs well on lexical metrics such as BLEU and CHRF, receives a notably lower score with CONTEXTCOMETQE.

| System | Precision | Recall | F1 |
|---|---|---|---|
| HW-TSC | 76.7 | 86.2 | 81.2 |
| SETU-ADAPT | 75.0 | 69.2 | 72.0 |
| SheffieldGATE | 73.0 | 70.8 | 71.9 |
| CLTeam | 75.7 | 81.5 | 78.5 |
| DCUGenNLP | 74.6 | 81.5 | 77.9 |
| Unbabel-IT | 75.4 | 66.2 | 70.5 |
| NLLB-3.3B | 74.3 | 84.6 | 79.1 |

Table 10: MUDA scores for EN-DE pronouns.

**Discourse Phenomena Analysis** Figure 2 shows the F1 accuracy for all systems in correctly using the discourse markers across multiple phenomena for all language pairs. The baseline system (NLLB-3.3B) has competitive accuracy with submitted systems on higher resource language pairs (EN→DE and EN→FR). For all settings except "pronouns" for German and "formality" for German and French, UNBABEL-IT achieves the highest accuracy across the board. Surprisingly, the MUDA F1 score for correctly generating German pronouns is worse for UNBABEL-IT relative to the baseline. A qualitative analysis shows that this is due to pronouns being under-generated in UNBABEL-IT's translations resulting in high precision but low recall scores as shown in Table 10.

To validate the observations and findings derived from automatic metrics, we now turn to human evaluation of the submitted systems for a more reliable assessment of translation quality.

## 4.2 Human Evaluation

We present the human evaluation results at both turn and conversation levels in Tables 11 and 12 respectively.

**Overall results.** UNBABEL-IT outperforms all systems on both turn-level and conversation-level evaluation, surpassing the HW-TSC system that achieved the highest COMET scores on EN-DE



Figure 2: MUDA F1 scores across all settings.

translation pair.[14] The translation quality according

---

[14]We note that the human evaluation for EN-DE, like other LPs, was conducted on a subset of the dataset (limited to a maximum of 30 turns per conversation).

| SYSTEM | EN-DE | | EN-FR | | EN-NL | | EN-PT | | EN-KO | |
|---|---|---|---|---|---|---|---|---|---|---|
| | XX | EN | XX | EN | XX | EN | XX | EN | XX | EN |
| DeepText Lab | | | | | | | | | 91.35 | 95.71 |
| HW-TSC | 88.47 | 90.41 | | | | | | | | |
| MULTITAN-GML | | | 81.83 | 84.62 | | | | | | |
| ADAPT | 82.55 | 88.83 | 70.22 | 65.53 | | | | | | |
| SheffieldGATE | 78.63 | 88.85 | | | 85.62 | 94.18 | 73.34 | 81.53 | | |
| CLTeam | 83.12 | 89.12 | 84.28 | 85.79 | 93.39 | 95.83 | 74.14 | 80.52 | | |
| DCUGenNLP | 84.56 | 88.60 | 85.72 | 83.26 | 91.30 | 94.61 | 80.21 | 81.55 | 89.71 | 96.15 |
| Unbabel-IT | **89.42** | **92.74** | **90.24** | **90.00** | **98.16** | **97.40** | **82.04** | **82.37** | **93.39** | **96.31** |
| NLLB-3.3B | 78.05 | 87.57 | 80.59 | 77.82 | 82.66 | 90.98 | 61.27 | 73.98 | 79.13 | 90.47 |

Table 11: Human Evaluation results aggregated at the turn level on the official test set.



Figure 3: Conversation-level DA scores.

| | EN-DE | EN-FR | EN-NL | EN-PT | EN-KO |
|---|---|---|---|---|---|
| DeepText Lab | | | | | 90.04 |
| HW-TSC | 81.19 | | | | |
| MULTITAN-GML | | 68.59 | | | |
| ADAPT | 75.75 | 59.65 | | | |
| SheffieldGATE | 75.72 | | 70.81 | 68.27 | |
| CLTeam | 78.61 | 73.32 | 84.37 | 69.85 | |
| DCUGenNLP | 77.03 | 72.27 | 76.41 | 73.78 | 89.83 |
| Unbabel-IT | **84.22** | **79.62** | **92.22** | **78.00** | **93.21** |
| NLLB-3.3B | 74.50 | 67.81 | 53.07 | 56.37 | 85.63 |

Table 12: Human Evaluation results aggregated at the conversation level on the official test set.

to direct assessment scores of all systems evaluated across all language pairs is high ($> 65$) at both conversation and turn levels. This could be because of the nature of the chat dataset which contains very short texts (the number of words per turn across language pairs is less than 8, see Table 2).

**Conversation-level results.** Figure 3 shows the distribution of scores assigned at the conversation level for all systems and language pairs. Confirming the automatic results, NLLB-3.3B scores the lowest and with the highest standard deviation for EN-KO, EN-NL and EN-PT. We also observe that EN-NL generally exhibits the largest standard deviation. After analyzing the outputs, we found that EN-NL has the highest number of segments (and conversations) receiving either a score of 0 (when hallucinating or copying source text verbatim) or 100, indicating a significant variation in translation quality for this language pair. Although there are sentences with mid-range scores, the dominance of segments with extremely high or low scores greatly influences the overall results, substantially raising the standard deviation.

**Turn-level results.** Figure 4 illustrates DA scores with the increase in the number of turns. For most systems and language pairs, translation quality deteriorates over successive turns, indicating a decline in the systems' ability to maintain consistency and accuracy in prolonged dialogues. This decline is particularly evident in the baseline sys-

Figure 4: Turn-level DA score across different language pairs through a chat.



Figure 5: Turn avg. vs conversation-level DA scores.

tem, which does not leverage contextual information from previous turns to generate translations. Interestingly, however, despite not using contextual information, HW-TSC's system maintains translation quality across successive turns. This can likely be attributed to rigorous training on in-domain data, both authentic and synthetically generated.

**Turn Vs. Conversation Quality results.** Overall, **conversation-level quality is lower than turn-level scores** suggesting that there are aspects beyond translation accuracy that might impact the overall translation quality and user experience. This is corroborated by the observation that the Spearman correlation between the average turn-level score and conversation-level DA score, though high, is 0.722. For future evaluations, it

might be worth investigating dialogue-oriented human assessment (Mendonca et al., 2023) to understand how turn-level scores impact conversation-level quality.

While direct assessments from experts provide a reliable measure of translation quality, DA scores fall short in offering insights into when and how errors occur, as well as their types and nature. Therefore, to assess the severity of errors generated by these systems, we now turn to LLM-based fine-grained error assessment of translation outputs.

### 4.3 LLM-based Evaluation

| System | % Perfect | # Minor | # Major | # Critical | Avg. Score |
|--------|-----------|---------|---------|------------|------------|
| HW-TSC | 89.12 | 100 | 88 | 59 | -0.554 |
| SETU-ADAPT | 82.61 | 158 | 139 | 99 | -0.903 |
| SheffieldGATE | 77.95 | 220 | 178 | 95 | -1.009 |
| CLTeam | 86.28 | 139 | 82 | 79 | -0.656 |
| DCUGenNLP | 83.10 | 143 | 158 | 80 | -0.849 |
| Unbabel-IT | **94.41** | **51** | **47** | **18** | **-0.228** |
| NLLB-3.3B | 80.50 | 161 | 143 | 117 | -1.002 |

Table 13: CONTEXTMQM scores for EN-DE.

Table 13 shows the results from using LLM-based error assessments via CONTEXTMQM. UNBABEL-IT leads the pack with 94.41% perfect translations. It also has the lowest number of errors in each category (minor, major, and critical), with an average error score of -0.228 (less than 1 minor error), the best among all systems. All systems, however, manage to achieve over 77% perfect

translations, meaning the overall quality across the board is strong.

Despite the positive results, there are notable differences in error distribution. For example, both the SHEFFIELDGATE and SETU-ADAPT models, while maintaining a reasonable percentage of perfect translations (82.61% and 77.95%, respectively), suffer from a significantly higher number of errors across all categories—minor, major, and critical. This suggests that when these systems do make errors, they tend to be more frequent and more serious, dragging down their overall performance compared to other systems. Interestingly, contrary to human evaluation but in line with other automatic measures, DCUGENNLP scores worse than CLTEAM submission, highlighting limitations of existing evaluation methods to discern systems with close translation quality.

## 5 Conclusions

This paper presents the findings of the Chat Translation Shared Task 2024. This year, we expanded the set of language pairs to include two additional languages (EN-KO and EN-NL). We created the evaluation sets with a focus on context usage when assessing system performance. We also employed a range of complementary evaluation methods to assess all systems, including automatic metrics that focus on translation quality, as well as fine-grained error assessments and analysis of specific discourse phenomena.

We find that the best systems finetune strong pre-trained LLMs using multilingual in-domain data and use contextual information (such as graphs, summaries or raw context) during training and inference. Additionally, using synthetic data during training improved translation quality. Furthermore, QAD strategies were effective in aligning translations with quality expectations.

As future work, a possible direction is to leverage reference-free discourse quality metrics that can give complementary insights to the translation evaluation approaches we tried this year. It might also be worth investigating human and automatic evaluation frameworks that assess specific dimensions relevant to chat (e.g. fluidity, coherence, consistency, etc).

## Limitations

Due to budget constraints, we conducted human evaluations using DA on a subset of the test set, which limited the number of turns evaluated for each language pair. For similar cost-related reasons, we ran CONTEXTMQM on a single language pair that received the highest number of submissions. Additionally, we note that our analysis of discourse-specific phenomena is constrained by the quality of taggers, which only annotate specific properties based on predefined rules and may not fully capture all levels of ambiguity present in chat datasets.

## Ethics Statement

## Acknowledgements

## References

Sujan Reddy A. 2022. Automating human evaluation of dialogue systems. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop*,

pages 229–234, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.

Sweta Agrawal, M. Amin Farajian, Patrick Fernandes, Ricardo Rei, and André F. T. Martins. 2024. Is context helpful for chat translation evaluation? *CoRR*, abs/2403.08314.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Y. Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *CoRR*, abs/2207.04672.

Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echegoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2021. Survey on evaluation methods for dialogue systems. *Artificial Intelligence Review*, 54:755–810.

M. Amin Farajian, António V. Lopes, André F. T. Martins, Sameen Maruf, and Gholamreza Haffari. 2020. Findings of the WMT 2020 shared task on chat translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 65–75, Online. Association for Computational Linguistics.

Ana C Farinha, M. Amin Farajian, Marianna Buchicchio, Patrick Fernandes, José G. C. de Souza, Helena Moniz, and André F. T. Martins. 2022. Findings of the WMT 2022 shared task on chat translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 724–743, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Christian Federmann. 2018. Appraise evaluation framework for machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 86–88, Santa Fe, New Mexico. Association for Computational Linguistics.

Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André Martins, Graham Neubig, Ankush Garg, Jonathan Clark, Markus Freitag, and Orhan Firat. 2023a. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1066–1083, Singapore. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig,

and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Patrick Fernandes, Kayo Yin, Emmy Liu, André Martins, and Graham Neubig. 2023b. When does translation require context? a data-driven, multilingual exploration. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 606–626, Toronto, Canada. Association for Computational Linguistics.

Madalena Gonçalves, Marianna Buchicchio, Craig Stewart, Helena Moniz, and Alon Lavie. 2022. Agent and user-generated content and its impact on customer support MT. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 201–210, Ghent, Belgium. European Association for Machine Translation.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. *International Conference on Learning Representations*.

Tom Kocmi and Christian Federmann. 2023. GEMBA-MQM: Detecting translation quality error spans with GPT-4. In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775, Singapore. Association for Computational Linguistics.

Qingyu Lu, Baopu Qiu, Liang Ding, Kanjian Zhang, Tom Kocmi, and Dacheng Tao. 2024. Error analysis prompting enables human-like translation evaluation in large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 8801–8816, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

John Mendonca, Patrícia Pereira, Miguel Menezes, Vera Cabarrão, Ana C Farinha, Helena Moniz, Alon Lavie, and Isabel Trancoso. 2023. Dialogue quality and emotion annotations for customer support conversations. In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 9–21, Singapore. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*,

pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.

Yi-Ting Yeh, Maxine Eskenazi, and Shikib Mehri. 2021. A comprehensive assessment of dialog evaluation metrics. In *The First Workshop on Evaluations and Assessments of Neural Conversation Systems*, pages 15–33, Online. Association for Computational Linguistics.

# Findings of the WMT 2024 Shared Task on Non-Repetitive Translation

**Kazutaka Kinugawa[1], Hideya Mino[1], Isao Goto[2], Naoto Shirai[1]**

[1]NHK Science & Technology Research Laboratories, Tokyo, Japan
[2]Ehime University, Ehime, Japan
{kinugawa.k-jg,mino.h-gq,shirai.n-hk}@nhk.or.jp
goto.isao.fn@ehime-u.ac.jp

## Abstract

The repetition of words in an English sentence can create a monotonous or awkward impression. In such cases, repetition should be avoided appropriately. To evaluate the performance of machine translation (MT) systems in avoiding such repetition and outputting more polished translations, we presented the shared task of controlling the lexical choice of MT systems. From Japanese–English parallel news articles, we collected several hundred sentence pairs in which the source sentences containing repeated words were translated in a style that avoided repetition. Participants were required to encourage the MT system to output tokens in a *non-repetitive* manner while maintaining translation quality. We conducted human and automatic evaluations of systems submitted by two teams based on an encoder-decoder Transformer and a large language model, respectively. From the experimental results and analysis, we report a series of findings on this task.

## 1 Introduction

The development of neural models has improved the performance of machine translation (MT) significantly (Sutskever et al., 2014; Bahdanau et al., 2015; Vaswani et al., 2017). MT systems are now used in a variety of real-world scenarios; however, challenges remain for such systems that assist human writers. Specifically, the MT output must not only be *adequate* and *fluent* but also follow the writing style of the target domain. For example, it is advisable for an application in the English news domain to follow rules such as the use of active rather than passive voice, the use of the affirmative rather than the negative, and the avoidance of redundant phrases (Block, 1994; Cappon, 2019; Papper, 2021). Among these writing style rules, we focus on the rule regarding the repetition of words in the English news domain. Generally, common words repeated in a sentence can create a monotonous

---



Figure 1: Motivating example from a Japanese–English parallel news article along with a consistent translation ("En(trans)") for comparison. Repeated words and their counterparts are highlighted. "入学" is intentionally removed (*reduction*), probably because it is contextually obvious. In this paper, we distinguish this type of removal from undertranslation. Additionally, "入学" and "入学" are translated differently as "join" and "enter," respectively (*substitution*).

---

or awkward impression, and in such cases, repetition should be avoided appropriately (Burstein and Wolska, 2003). Typical workarounds are (1) the removal of redundant terms, if possible (Strunk and White, 1999) or (2) the use of alternative words, such as synonyms, as substitutes.[1] In this paper, we refer to translation techniques (1) and (2) as *reduction* and *substitution*, respectively, and call the translation style using these techniques a *non-repetitive style*. Figure 1 shows an example of a non-repetitive style translation from a Japanese–English parallel news article. We observe that human writers in the English news domain often translate Japanese text with such reduction and substitution.[2] These translation techniques arise from the difference between the styles of the source and

---

[1]https://effectiviology.com/
writing-tips-from-the-elements-of-style/#Avoid_
repetition

[2]Other examples are listed in Appendix A.

target languages; that is, an article in the source language was originally produced by a writer who attached importance to conveying content without the reader misunderstanding it by using the same expressions consistently, and then it was translated into the target language by a writer who was encouraged to (or preferred to) translate it in a more diverse or concise way. The assumption in this task is that sentences translated in a simple word-by-word manner cannot be suited to the target domain. We could thus associate these translation techniques with a type of rewriting. Although this task focuses on the news domain, the monotony or awkwardness arising from the repetition of words in English is also a common problem in other domains.

Given this motivation, we presented a shared task for non-repetitive translation. To configure appropriate settings, we limited the task to one-to-one or two-to-two translations. We hypothesized that the closer the distance between repeated words, the greater the need to translate using reduction or substitution. Additionally, we targeted the repetition of common words because such words tend to be substituted according to the findings of Guillou (2013). We qualitatively categorized several patterns of non-repetitive style translations, and then collected several hundred instances in which some words were repeated in the source sentence and translated using reduction or substitution, which we used as development and test data. In the remainder of this paper, we first explain the research background of this task (§2). Next, we describe the task definition (§3), dataset we prepared (§4), evaluation methods (§5), and submitted systems (§6). Finally, we present the results and some analysis (§7).

## 2 Related Work

**Contrast with Consistent Translation**    In the context of MT research, lexically consistent translation (in this paper, we also refer to this as *repetitive style translation*) has been studied actively (Pu et al., 2017; Kuang et al., 2018; Tu et al., 2018; Lyu et al., 2021, 2022). A representative study is the hypothesis of "one translation per discourse," which was advocated by Carpuat (2009). The motivation for these studies is the assumption that translating text in a consistent style should be encouraged because this style is unambiguous and accurate for readers. Moreover, from the viewpoint of experimental evaluation, many researchers have

reported that BLEU scores improved as a result of encouraging consistent translation (Lyu et al., 2021, 2022). However, it is debatable whether all words should be translated consistently. Translation consistency can depend on several factors, such as the target domain, type of words, and translation direction (Guillou, 2013). For example, it is indisputable that technical terms in the patent domain should be translated consistently. By contrast, Guillou (2013) reported that high-frequency verbs are often translated in diverse ways in English–French translation. While improving document-level consistency based on the postprocess approach, Zhang et al. (2023) also mentioned the side effect of the loss of translation diversity. From another point of view, consistent translation has the risk of leading to a robotic wording and giving a monotonous or awkward impression to readers, as shown in Figure 1. By contrast, Cappon (2019) claimed that excessive substitution may obscure the meaning of the sentence. In monolingual writing, this phenomenon is derided as *the elegant variation.*[3] To summarize, there is a trade-off between ambiguity and monotony. This task particularly focuses on the latter aspect, which has not often been addressed in previous studies. To the best of our knowledge, no test sets exist for directly evaluating such a translation style.

**Reduction and Substitution**    Although several studies have been conducted related to non-repetitive translation, the scope of our research is different. First, several researchers have addressed the problem of controlling the output length of MT systems (Lakew et al., 2019; Schioppa et al., 2021). Typically, special tokens representing the output length at several discrete levels are inserted into source sentences. Although this approach is associated with reduction, our task requires a more meticulous omission of specific words in sentences. Regarding substitutions, MT systems are sometimes required to select infrequent words from the vocabulary. However, researchers have reported that MT systems are biased toward outputting high-frequency target words (Ott et al., 2018; Gu et al., 2020) and tend to produce lexically poorer translations than humans (Vanmassenhove et al., 2019, 2021). Gu et al. (2020) designed the objective function so that low-frequency target tokens were more likely to be output. However, they conducted

---

[3]https://en.wikipedia.org/wiki/Elegant_
variation

the experiment using regular corpora and did not present a perspective on in what scenarios low-frequency words should be output. By contrast, we set up a more specific scenario.

This task is also related to research outside of translation technique. Neural models have the traditional problem of not outputting the end-of-sequence token while generating the same tokens endlessly. To alleviate this problem, several approaches including learning-based methods (Welleck et al., 2020) and decoding-based methods (Keskar et al., 2019), have been proposed. Although the goal is different, these studies are also relevant to our task in the sense that word repetition should be avoided.

## 3 Task Definition

Our task focused on lexical choice in MT, particularly choice regarding repeated words in a source sentence. The translation direction was Japanese to English. Participants were required to control an MT system using reduction or substitution so that it did not output the same words for certain repeated words in a source sentence. Simultaneously, participants also needed to maintain translation quality as much as possible.

The challenges underlying this task included the following:

- Maintaining the balance between translation quality and controlling the output: Translation quality can be degraded when the non-repetitive style is enforced inappropriately.

- Avoiding bias toward high-frequency bilingual word pairs: Generally, for a given source word, high-frequency target words associated with it are more likely to be output. This can make it difficult to determine appropriate substitutions for some words.

- Predicting which words can be reduced or substituted: It is not easy to make an appropriate prediction because it depends on the context within the sentence.

- Mining training instances: Translations with reduction can be particularly difficult to identify in noisy corpora because of the challenge of discriminating them from undertranslations.

## 4 Dataset

We prepared the training, development, and test data for this task. They were all sourced from Japanese–English news articles published by Jiji Press LTD., a Japanese news agency. We annotated the development and test data for this task, whereas the training data comprised a regular MT corpus.

### 4.1 Development and Test Data

We provided development and test sets for this task, which we refer to as Jiji 2023 data and Jiji 2024 data, respectively. These data included 162 and 479 instances, respectively. The Jiji 2023 data were originally built for the Non-Repetitive Translation Task in WAT 2023 (Nakazawa et al., 2023). We reviewed the data and filtered out some instances this year. By contrast, the Jiji 2024 data were newly created in this year. In both datasets, all Japanese sentences contained some repeated words that were translated into English with reduction or substitution. From Japanese–English news articles, we first automatically created sentence pairs based on lexical similarities using the method of (Utiyama and Isahara, 2007) and then manually selected appropriate instances. To reduce the negative effects of imbalanced content in the source and target sentences, the Japanese sentences in the Jiji 2023 and 2024 data were manually translated from English by professional translators while preserving as much of the vocabulary of the original Japanese sentences as possible. Both the released development and test sets contained raw and tagged parallel data. In the tagged data, we marked repeated words in the source sentence and their counterparts in the target sentence with tags, which indicated that these words were evaluation targets. Examples are shown in Table 1. The respective attributes inside the tags indicate the following:

id: This indicates the IDs of repeated words. In the example, two tagged repeated words are included, that is, "機能" ("id=0") and "製品" ("id=1"). The number of instances including multiple tagged repeated words, such as this example, are limited. Additionally, the number of types of repeated words in one instance is one or two.

ref: This indicates the IDs of pairs of source words and their counterparts, such as ("製品," "models") (i.e., "id=1" and "ref=0") and ("製品," "products") (i.e., "id=1" and "ref=1").

| | |
|---|---|
| Ja | JEMAの担当者は白物家電について、「<target id=0 ref=0 type=s>機能<\target>を絞った低価格<target id=1 ref=0 type=s>製品<\target>、高価格な高<target id=0 ref=1 type=s>機能<\target><target id=1 ref=1 type=s>製品<\target>とも好調だ」と述べている。 |
| En | "Shipments have been robust for both low-priced <target id=1 ref=0 type=s>models<\target> with reduced <target id=0 ref=0 type=s>functions<\target> and expensive <target id=0 ref=1 type=s>high-spec<\target> <target id=1 ref=1 type=s>products<\target>," a JEMA official said. |

Table 1: Examples of tagged instances in the development and test data. The tags are highlighted.

| Split | # Parallel sentences |
|---|---|
| train | 200K |
| dev | 479 |
| test | 1851 |

Table 2: Statistics of the Jiji 2020 data. Note that "dev" and "test" in the table are different from the Jiji 2023 and 2024 data.

type: This indicates whether tagged source words are substituted ("s") or reduced ("r").

Note that not all words repeated in the source sentence were evaluation targets. This is because some words, such as proper nouns and technical terms, should be translated consistently, even if they were repeated in the sentence. We provided the tagged development data to help to tune the model during training. However, participants could not use the tagged test data when submitting the system results. In this task, the systems had to detect repeated words that could be reduced or substituted on their own.

## 4.2 Training Data

Regarding the training data, we provided all the data from the WAT 2020 Newswire tasks (Nakazawa et al., 2020), which were also constructed from Jiji news articles and have been continuously used in WAT since 2020 (Nakazawa et al., 2020, 2021, 2022, 2023). For simplicity, we refer to these data as Jiji 2020 data. The main files in the Jiji 2020 data are shown in Table 2. These data are a regular parallel corpus. They were not annotated specifically for this task but were in exactly the same domain as the Jiji 2023 and 2024 data. Although the development and test sets in the Jiji 2020 data, which are described as "dev" and "test" in Table 2, were not directly related to the evaluation of this task, they could be used to measure basic translation performance during training. Unfortunately, the number of parallel sentences in the Jiji 2020 data was limited. Thus, we allowed participants to use any other corpora for training.

## 5 Evaluation

We conducted both human and automatic evaluation. We based the main results of this task on the human evaluation and prepared the automatic evaluation as secondary metrics. Again, the goal of this task was to control an MT system to output translations in a non-repetitive style while maintaining translation quality.

## 5.1 Human Evaluation

We evaluated system performance using the total number of outputs that met both acceptable translation adequacy and appropriate lexical choice. Both aspects were checked by three human translators, who were assigned by the authors.

**Translation Style** Regarding the evaluation for lexical choice, the human translators checked whether the translations for the tagged source words were correctly written in a non-repetitive style. Whether untagged repeated words were translated in a repetitive or non-repetitive way did not affect this evaluation. Moreover, the technique (i.e., reduction or substitution) did not have to be consistent with that of the reference translation. In our preliminary investigations, we qualitatively studied the lexical choices of several translators, and observed cases in which one translator chose substitution, and another chose reduction. Additionally, the systems did not have to choose the same words used in the reference, provided the meaning was appropriate. The determination of substitution or repetition was essentially based on the word stem. For example, conversions between voice (e.g., "attack" and "be attacked"), tense (e.g., "study" and "studied"), and parts of speech (e.g., "problematic" and "problem") were not considered to be substitutions. Conversions to idioms (e.g., "visit" and "pay a visit") were an exception and handled as substitutions. This evaluation is not trivial. For example, it is difficult to establish uniform guidelines for determining the correctness of synonyms in substitution and whether they are appropriate reductions

The $i$-th source sentence

そのうち、21**団体**(id=1)で被害が確認され、11**団体**(id=1)が**調査**(id=2)困難とし、14**団体**(id=1)が**調査**(id=2)中としている。

The system output for the $i$-th source sentence

Of them, 21 have been confirmed to have suffered damage, 11 have found it difficult to **investigate**(id=2), and 14 are under **investigation**(id=2).

The evaluation results for the $i$-th test instance



Figure 2: Example of human evaluation for the $i$-th test instance. "団体" (id=1) is undertranslated (at least one counterpart should appear in the output in this case) (the label is thus "M"), and "調査" (id=2) is translated in a repetitive style (the label is thus "C"). For simplicity, 1-indexed IDs are used for the repeated words.

or inappropriate omissions. Thus, we adopted a majority vote by the three human translators in this evaluation process.

Next, we explain the evaluation procedure for the $i$-th test instance, which is also illustrated in Figure 2. Other test instances were also evaluated in the same manner and all results were finally aggregated. First, each translator labeled the translations for the tagged source words in the $i$-th test instance as "S" (substitution), "R" (reduction), "C" (consistent, i.e., repetitive), or "M" (mistranslation or undertranslation). Note that "S," "R," and "C" implicitly indicate that the meaning of the translation is correct. Let the label for the $j$-th evaluation target in the $i$-th test instance given by the $k$-th translator be $t_{i,j}^{(k)}$. Next, the three labels $t_{i,j}^{(1)}$, $t_{i,j}^{(2)}$, and $t_{i,j}^{(3)}$ were reduced to one by a majority vote, which we denote by $t_{i,j}$. Because the number of types of labels was more than two, three labels could all be different. Although we assumed that such a case was limited, we introduced an additional heuristic rule to determine the label as follows:

- If the label set was equal to {"C,""R,""S"}, "S" was assigned to $t_{i,j}$: Because two translators thought it was correctly translated in a non-repetitive style, the label should be "R" or "S." Next, because two translators thought

the word was not reduced, the label was determined to be "S."

- If the label set was equal to {"M,""R,""S"}, "R" was assigned to $t_{i,j}$: Because two translators thought it was correctly translated in a non-repetitive style, the label should be "R" or "S." Next, the label "M" was assigned probably because that translator thought some necessary word was omitted. Thus, the label was determined to be "S."

- If the label set was equal to {"M,""C,""S"}, "S" was assigned to $t_{i,j}$: Because two translators thought the meaning of the translation was correct, the label should be "C" or "S." Next, the label "M" was assigned probably because that translator thought some word had a slightly different nuance. Thus, the label was determined to be "S."

- If the label set was equal to {"M,""C,""R"}, "R" was assigned to $t_{i,j}$: Because two translators thought the meaning of the translation was correct, the label should be "C" or "S." Next, the label "M" was assigned probably because that translator thought some necessary word was omitted. Thus, the label was determined to be "R."

Finally, one representative label was assigned to the $i$-th test instance, which we denote by $t_i$. Representative labels were chosen from "<NON-REP>," "<REP>," and "<INCORRECT>." For test instances including only one target, the representative label $t_i$ was simply mapped from $t_{i,1}$ as follows:

- If $t_{i,1}$ was "R" or "S," "<NON-REP>" was assigned to $t_i$.

- If $t_{i,1}$ was "M," "<INCORRECT>" was assigned to $t_i$.

- If $t_{i,1}$ was "C," "<REP>" was assigned to $t_i$.

For test instances including two targets, the representative label $t_i$ was determined as follows:

- If $t_{i,1}$ was "R" or "S," and $t_{i,2}$ was "R" or "S," "<NON-REP>" was assigned to $t_i$.

- If $t_{i,1}$ was "M" or $t_{i,2}$ was "M," "<INCORRECT>" was assigned to $t_i$.

- Otherwise, "<REP>" was assigned to $t_i$.

**Translation Accuracy**    In this task, the content of the system output may be omitted incorrectly or obscured if reduction or substitution is enforced inappropriately. Thus, we measured translation adequacy for system outputs. The evaluation framework was based on Japanese Patent Office (JPO) adequacy.[4] This criterion is well established and has also been used in domains other than patents.

Specifically, the $k$-th translator assigned a five-level discrete score $s_i^{(k)} \in \{1, 2, 3, 4, 5\}$ to the $i$-th system output. Next, we averaged $s_i^{(1)}$, $s_i^{(2)}$, and $s_i^{(3)}$ to $s_i$. Additionally, to view the balance between translation style and adequacy, we reflected the style label $t_i$ in the adequacy score $s_i$. If the translation style was not "<NON-REP>," we reduced the adequacy score $s_i$ to 0. We refer to this metric as *filtered adequacy* and denote it by $s_i'$.

## 5.2   Automatic Evaluation

We also automatically predicted whether the target word was translated in a repetitive style. Note that "<NON-REP>" and "<INCORRECT>" could not be discriminated in this process. Thus, we introduced one more label "<NOT-REP>," which indicated "<NON-REP>" or "<INCORRECT>."

Figure 3: Yes/no flowchart for predicting translation styles.

As a preprocess, we built a bilingual dictionary from the Jiji 2020 data and JParaCrawl v3.0 (Morishita et al., 2022). We aggregated translations of evaluation target words in the Jiji 2024 data by running the AWESOME aligner (Dou and Neubig, 2021) on the above corpora. Let the $j$-th evaluation target word in the $i$-th source sentence be $w_{i,j}$. Based on the alignment results, we obtained a set of possible counterparts of $w_{i,j}$, which we denoted by $S_{w_{i,j}}$. We then removed low-frequency counterparts from $S_{w_{i,j}}$ to limit the maximum dictionary size $|S_{w_{i,j}}|$ to 10. We predicted a style label by applying several simple binary classifications in order of reliability confidence as follows:

(1) *Do all tokens appear once each?*: If all content words appear once each in the $i$-th system translation, this output is classified as "<NOT-REP>."

(2) *Are estimated counterparts all the same and included in the dictionary?*: First, we estimate counterparts of $w_{i,j}$ using the word aligner. If these counterparts are all the same and exist in $S_{w_{i,j}}$, this output is classified as "<REP>."

(3) *Do any tokens in the dictionary appear more than once?*: If any word in $S_{w_x}$ appears more than once, this output is classified as a repetitive style; otherwise, the output is classified

as "<NOT-REP>."

We designed the third block to mitigate misclassification caused by alignment errors in (2). The above procedures are illustrated in Figure 3. Finally, we calculated the percentage of instances labeled as "<REP>" in the test set. We refer to this metric as *repetition rate*.

To measure translation quality, we also computed BLUE scores (Papineni et al., 2002) using SacreBleu (Post, 2018).[5]

## 6 Systems

In this shared task, two teams submitted the system and description paper. In this section, we provide an overview of the submitted systems and the baseline system that we built. For comparison, resources used by each system are listed in Table 3.

### 6.1 Baseline

As a baseline, we built an MT system using fairseq (Ott et al., 2019). We adopted Transformer (big) (Vaswani et al., 2017) as the architecture, and used the Jiji 2020 and JParaCrawl v3.0 (Morishita et al., 2022) as training data. We based the method on the tagging approach (Sennrich et al., 2016; Lakew et al., 2019; Johnson et al., 2017; Schioppa et al., 2021). Specifically, we introduced style and domain tags, and combined them. First, from the Jiji 2020 data and JParaCrawl v3.0, we mined sentence pairs in which some content words were repeated in the source sentence and no content words were repeated in the target sentence. We detected content words in Japanese and English sentences using GiNZA[6] and spaCy,[7] respectively. To avoid selecting noisy instances, we excluded parallel sentences with lexical similarity scores less than 0.7 from the tagging. Specifically, we prepended the style tag "<NON-REP>" and all repeated words to the source sentences as follows:

**Src:** <NON-REP> <文書> 米国立公文書館が
文書を保管していた。

Second, we also attached the domain tag "<JIJI>" to training instances from the Jiji 2020 data. Similarly, we did not tag sentence pairs with lexical similarity scores less than 0.7. We prepended the domain tag to the source sentences as follows:

| System | Resource |
|---|---|
| Baseline | JParaCrawl v3.0, Jiji 2020 |
| SYSTRAN | all Ja–En data from OPUS, Jiji 2020 |
| Waseda Riko | Claude 3.5 Sonnet, examples from the task website[8] |

Table 3: Comparison of resources used by each system. Although the Waseda Riko system did not explicitly use the data on which Claude 3.5 Sonnet was built, they are also listed as "Claude 3.5 Sonnet" in the table.

**Src:** <JIJI> <NON-REP> <文書> 米国立公文書
館が文書を保管していた。

For inference, we prepended the style and domain tags to all the test source sentences. In this system, we adopted the same hyperparameter settings as Morishita et al. (2022).

### 6.2 SYSTRAN (Avila and Crego, 2024)

The team introduced a *repetition penalty* in the fine-tuning phase. The method was inspired by label smoothing (Szegedy et al., 2015). For training instances including word repetition in the target sentence, the ground-truth score corresponding to the repeated word was decreased from 1. The team automatically detected such instances using the spaCy tokenizer[9] and GIZA++ toolkit (Och and Ney, 2003). Specifically, the repetition penalty was combined with label smoothing, and is formulated as follows:

$$q'_t = (1 - \epsilon)(1 - \alpha_t)q_t + \frac{\epsilon}{V},$$

where $q_t$ indicates a one-hot vector used as the ground-truth label at the $t$-th time step, at which a repeated word appears, $\epsilon$ is a hyperparameter for label smoothing and $V$ is the vocabulary size. $\alpha_t$ is also a hyperparameter used to control the degree to which word repetition is discouraged. The team first trained a Transformer encoder-decoder model on parallel sentences from OPUS[10] and then fine-tuned the model on parallel sentences from the Jiji 2020 data using the above technique. To avoid feeding noisy instances into the model, the team used back-translated sentences instead of the original sentence pairs in the fine-tuning stage.

### 6.3 Waseda Riko (Wang et al., 2024)

The team built a large language model (LLM)-based pipeline. The procedure was composed of

---

| System | Adequacy (↑) | Translation Style | | | Filtered Adequacy (↑) |
|---|---|---|---|---|---|
| | | % <NON-REP> | % <REP> | % <INCORRECT> | |
| Waseda Riko | **4.6** | **89.8** | 8.1 | 2.1 | **4.1** |
| SYSTRAN | 3.9 | 32.3 | 53.8 | 13.6 | 1.3 |
| Baseline | 3.9 | 50.2 | 27.4 | 22.3 | 2.1 |

Table 4: Human evaluation results.

| System | BLEU (%) (↑) | Translation Style | | Repetition Rate (%) (↓) |
|---|---|---|---|---|
| | | # <NOT-REP> | # <REP> | |
| Waseda Riko | 24.4 | 413 | 57 | **12.1** |
| SYSTRAN | 28.9 | 214 | 256 | 54.5 |
| Baseline | **29.1** | 332 | 138 | 29.4 |

Table 5: Automatic evaluation results.

the following four steps:

(1) Preprocess: Detect repeated words from the source sentence using the MeCab tokenizer (Kudo et al., 2004) and tag these possible repeated words.

(2) Translation: Instruct the LLM to translate the tagged source sentence in a non-repetitive manner using a few-shot prompt (Brown et al., 2020).

(3) Proofreading: Instruct the LLM to review the output in the previous step and rewrite the translation as needed to enhance the result.

(4) Postprocess: Tag the counterparts in the target sentence.[11]

The team used Claude 3.5 Sonnet[12] and designed a prompt suited for this task. Specifically, they instructed the LLM to output translations along with the estimated counterparts and translation labels in JSON format. Because of this structured output design, the following processes were performed successfully.

## 7 Results and Discussion

### 7.1 Human Evaluation

We summarize the human evaluation scores of all systems in Table 4.[13] The Waseda Riko system achieved the best results in both translation adequacy and style control. Focusing on the drop from the adequacy score to the filtered adequacy score, the baseline system lost 1.8 points, whereas the

Waseda Riko system only decreased by 0.5 points. This difference highlights that the Waseda Riko team successfully controlled the translation style without compromising translation quality. The SYSTRAN system achieved an adequacy score competitive with that of the baseline system, but passed more source sentences in a repetitive style. By contrast, the baseline system was the worst in terms of the percentage of incorrect instances. Considering the difference between the SYSTRAN and baseline systems, a trade-off existed between style control and translation adequacy.

The basic idea of the Waseda Riko system is similar to that of the baseline system: possible repeated words in the source sentence were automatically detected using a third-party tokenizer and the model was explicitly informed about them. (Wang et al. (2024) also reported that it was still difficult for LLMs to consistently identify repeated words in the input sentence.) Although the baseline system was trained on parallel sentences that were (possibly) translated in a non-repetitive style, the percentage of test instances in the desired style was 50%. Although the results of the Waseda Riko team were also supported by the high performance of the commercial LLM, their proposed prompt design and pipeline configuration were equally important. The key was how to provide the instruction to "translate in a non-repetitive style," which is (probably) new and complex for many LLMs. We attempted to instruct GPT-3.5 turbo[14] to solve this task using a simple prompt, such as *"Translate the following Japanese news text into English using as few of the same content words as possible,"* in our preliminary experiments, but this did not work well.

---

## 7.2 Automatic Evaluation

We also summarize the automatic evaluation scores of all systems in Table 5. In contrast to the human evaluation, the baseline and SYSTRAN systems achieved a better BLEU score than the Waseda Riko system. This gap depended on whether the systems used the Jiji 2020 data for training. Although the Waseda Riko team analyzed these data and then built the heuristic rules to detect repeated words (Wang et al., 2024), the team did not fully train the LLM on these data. The LLM learned the several translations from the Jiji data using the few-shot prompt, whereas the baseline and SYS-TRAN models adapted the output translations more directly to the target domain. Although we configured the primary results of this task based on the human evaluation, the motivation for this task was to adapt the lexical choice of MT systems to the target domain; thus, it should be noted that BLEU scores were also important metrics in our task.

Regarding the repetition rate, the trend was coincident with the human evaluation results. Specifically, the accuracy as a binary classifier (i.e., "<REP>" or not) between automatic and human evaluations was 93.4% in the baseline system, 92.1% in the SYSTRAN system, and 93.0% in the Waseda Riko system. Importantly, this metric had a certain degree of reliability independent of the success rate of style control and the degree of matching with the target domain.

## 8 Conclusion and Future Work

In this paper, we presented an overview of the WMT2024 Shared Task on non-repetitive translation. Particularly, the experimental results revealed the effectiveness of the LLM in controlling translation. We believe that our task will encourage further research on controlling MT systems. In the future, we will address several limitations in the current task settings. First, the test instances were limited to a comparatively short content. It would be an interesting challenge to address repetition observed in longer documents. Second, we will make both human and automatic evaluations more established. Currently, (1) evaluation relies heavily on human evaluation, and (2) the human evaluation is prone to variance. Regarding (2), specifically, although the percentage of test instances where the three translators voted for all different labels was limited, that of the test instances where the three translators voted for the same label was approxi-

mately 69%. These were partially because of (1) the difficulty of automatically detecting mistranslations and undertranslations, and (2) the difficulty of defining the correct answer for a translation output using substitution or reduction, respectively. Thus, we will develop more reliable evaluation guidelines in collaboration with translators. It would also be interesting to introduce automatic evaluation using LLMs.

## Acknowledgments

## References

Marko Avila and Josep Crego. 2024. Systran @ wmt24 non-repetitive translation task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, Florida, USA. Association for Computational Linguistics.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Mervin Block. 1994. *Broadcast Newswriting: The RT-NDA Reference Guide*. Bonus Books.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.

Jill Burstein and Magdalena Wolska. 2003. Toward evaluation of writing style: Overly repetitious word use. In *10th Conference of the European Chapter of the Association for Computational Linguistics*,

pages 35–42, Budapest, Hungary. Association for Computational Linguistics.

Rene J. Cappon. 2019. *The Associated Press Guide to News Writing*, fourth edition. Peterson's.

Marine Carpuat. 2009. One translation per discourse. In *Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions (SEW-2009)*, pages 19–27, Boulder, Colorado. Association for Computational Linguistics.

Zi-Yi Dou and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2112–2128, Online. Association for Computational Linguistics.

Shuhao Gu, Jinchao Zhang, Fandong Meng, Yang Feng, Wanying Xie, Jie Zhou, and Dong Yu. 2020. Token-level adaptive training for neural machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1035–1046, Online. Association for Computational Linguistics.

Liane Guillou. 2013. Analysing lexical consistency in translation. In *Proceedings of the Workshop on Discourse in Machine Translation*, pages 10–18, Sofia, Bulgaria. Association for Computational Linguistics.

Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.

Nitish Shirish Keskar, Bryan McCann, Lav Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL - A Conditional Transformer Language Model for Controllable Generation. *arXiv preprint arXiv:1909.05858*.

Shaohui Kuang, Deyi Xiong, Weihua Luo, and Guodong Zhou. 2018. Modeling coherence for neural machine translation with dynamic and topic caches. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 596–606, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Taku Kudo, Kaoru Yamamoto, and Yuji Matsumoto. 2004. Applying conditional random fields to Japanese morphological analysis. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 230–237, Barcelona, Spain. Association for Computational Linguistics.

Surafel Melaku Lakew, Mattia Di Gangi, and Marcello Federico. 2019. Controlling the output length of

neural machine translation. In *Proceedings of the 16th International Conference on Spoken Language Translation*, Hong Kong. Association for Computational Linguistics.

Xinglin Lyu, Junhui Li, Zhengxian Gong, and Min Zhang. 2021. Encouraging lexical translation consistency for document-level neural machine translation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3265–3277, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Xinglin Lyu, Junhui Li, Shimin Tao, Hao Yang, Ying Qin, and Min Zhang. 2022. Modeling consistency preference via lexical chains for document-level neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6312–6326, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Makoto Morishita, Katsuki Chousa, Jun Suzuki, and Masaaki Nagata. 2022. JParaCrawl v3.0: A large-scale English-Japanese parallel corpus. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6704–6710, Marseille, France. European Language Resources Association.

Toshiaki Nakazawa, Kazutaka Kinugawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida, Makoto Morishita, Ondřej Bojar, Akiko Eriguchi, Yusuke Oda, Chenhui Chu, and Sadao Kurohashi. 2023. Overview of the 10th workshop on Asian translation. In *Proceedings of the 10th Workshop on Asian Translation*, pages 1–28, Macau SAR, China. Asia-Pacific Association for Machine Translation.

Toshiaki Nakazawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida, Anoop Kunchukuttan, Makoto Morishita, Ondřej Bojar, Chenhui Chu, Akiko Eriguchi, Kaori Abe, Yusuke Oda, and Sadao Kurohashi. 2022. Overview of the 9th workshop on Asian translation. In *Proceedings of the 9th Workshop on Asian Translation*, pages 1–36, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Shohei Higashiyama, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, Ondřej Bojar, Chenhui Chu, Akiko Eriguchi, Kaori Abe, Yusuke Oda, and Sadao Kurohashi. 2021. Overview of the 8th workshop on Asian translation. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 1–45, Online. Association for Computational Linguistics.

Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Shohei Higashiyama, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, Ondřej Bojar, and Sadao Kurohashi. 2020. Overview of the 7th workshop on

Asian translation. In *Proceedings of the 7th Workshop on Asian Translation*, pages 1–44, Suzhou, China. Association for Computational Linguistics.

Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.

Myle Ott, Michael Auli, David Grangier, and Marc'Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3956–3965. PMLR.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Robert A. Papper. 2021. *Broadcast News and Writing Stylebook*, seventh edition. Routledge.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Xiao Pu, Laura Mascarell, and Andrei Popescu-Belis. 2017. Consistent translation of repeated nouns using syntactic and semantic cues. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 948–957, Valencia, Spain. Association for Computational Linguistics.

Andrea Schioppa, David Vilar, Artem Sokolov, and Katja Filippova. 2021. Controlling machine translation for multiple attributes with additive interventions. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6676–6696, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Controlling politeness in neural machine translation via side constraints. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 35–40, San Diego, California. Association for Computational Linguistics.

William Strunk and E. B. White. 1999. *The Elements of Style*, fourth edition. Pearson.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.

Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2015. Rethinking the inception architecture for computer vision. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826.

Zhaopeng Tu, Yang Liu, Shuming Shi, and Tong Zhang. 2018. Learning to remember translation history with a continuous cache. *Transactions of the Association for Computational Linguistics*, 6:407–420.

Masao Utiyama and Hitoshi Isahara. 2007. A Japanese-English patent parallel corpus. In *Proceedings of Machine Translation Summit XI: Papers*, Copenhagen, Denmark.

Eva Vanmassenhove, Dimitar Shterionov, and Matthew Gwilliam. 2021. Machine translationese: Effects of algorithmic bias on linguistic complexity in machine translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2203–2213, Online. Association for Computational Linguistics.

Eva Vanmassenhove, Dimitar Shterionov, and Andy Way. 2019. Lost in translation: Loss and decay of linguistic richness in machine translation. In *Proceedings of Machine Translation Summit XVII: Research Track*, pages 222–232, Dublin, Ireland. European Association for Machine Translation.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Qiao Wang, Yixuan Huang, and Zheng Yuan. 2024. Reducing redundancy in japanese-to-english translation: A multi-pipeline approach for translating repeated elements. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, Florida, USA. Association for Computational Linguistics.

Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2020. Neural text generation with unlikelihood training. In *International Conference on Learning Representations*.

Zhen Zhang, Junhui Li, Shimin Tao, and Hao Yang. 2023. Lexical translation inconsistency-aware document-level translation repair. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12492–12505, Toronto, Canada. Association for Computational Linguistics.

725

# A  Examples of Non-Repetitive Trasnlations

Table 6 shows examples of non-repetitive translations from the task website.[15]

| Reduction | Ja | 耐震化を済ませていない４９４団体に今後の対応を尋ねたところ、改修するのは７０団体、建て替えは２６５団体、移転が１１団体だった。 |
|---|---|---|
| | En(trans) | When the 494 organizations that had not yet completed earthquake proofing were asked about their future measures, 70 organizations opted for retrofitting, 265 chose rebuilding, and 11 selected relocation. |
| | En | Of the 494 unprepared municipalities, 70 are set to carry out repairs, 265 will construct new buildings and 11 are planning relocation. |
| | Note | In the original English sentence, a noun ellipsis occurs, e.g., "70 municipalities" is expressed as "70." |
| Reduction | Ja | 開発費を参加国間で分担できるため、国産開発に比べて費用を安く抑えることが可能となる。 |
| | En(trans) | Since development expenses can be shared among participating countries, it will be possible to keep costs lower than domestic development. |
| | En | It will allow the government to cut spending compared with full domestic development by sharing costs with partner countries. |
| | Note | "Costs" is used instead of "development costs" in the original English sentence probably because it is contextually inferable. |
| Reduction | Ja | 同社はニューヨーク州のヨンカース工場と中西部ネブラスカ州のリンカーン工場で車両の製造や試験を行う。 |
| | En(trans) | The company will manufacture and test vehicles at its Yonkers, New York, factory and its Lincoln, Nebraska, factory in the Midwest. |
| | En | Kawasaki Rail Car will build and test the subway cars at its facilities in Yonkers and in Lincoln, Nebraska. |
| | Note | The two nouns ("facility") are merged into one and the noun head is shared by the two prepositional phrases. Although strictly they are not reduced, we also consider these examples to be a type of reduction. |
| Substitution | Ja | 農作物への影響が心配されるが、農林水産省は「（首都圏などでは）積雪が長引かなかったので大きな影響はない」（園芸作物課）とみている。 |
| | En(trans) | There are concerns about the impact on crops, but an official at the Horticultural Crops Division of the Ministry of Agriculture, Forestry and Fisheries (MAFF) said, "the snowfall (in the Tokyo metropolitan area and other regions) was not prolonged, so there will be no major impact." |
| | En | Although many people are worried about the effects of harsh cold on crops, an official of Japan's agricultural ministry predicted that there will be no significant impact, as the snow did not stay for long in areas such as the Tokyo metropolitan area. |
| | Note | Words with similar meaning such as synonyms and hypernyms are typically used for substitution. |
| Substitution | Ja | 物質を構成する素粒子の振る舞いは「標準理論」で説明されるが、宇宙の全質量の４分の１を占める「暗黒物質」など説明できない部分もある。 |
| | En(trans) | The Standard theory explains the behavior of elementary particles, which make up matter, but it cannot explain some things, such as dark matter, which makes up one quarter of the mass of the universe. |
| | En | The so-called Standard Model explains the behavior of elementary particles, the fundamental building blocks of matter. But the theory leaves some mysteries, such as dark matter which is thought to make up about a quarter of the mass of the universe. |
| | Note | Repeated words are sometimes translated in a non-literal manner. |
| Substitution | Ja | 当時、テニス部の生徒６人とコーチがコートで練習をしており、生徒の１人がボールを拾おうとしたところ、隣のコートにパラシュート状の物があることに気付いたという。 |
| | En(trans) | At the time, six students and the coach from the tennis club were reportedly practicing on the court when one of the students went to pick up a ball and noticed a parachute-like object on the adjacent court. |
| | En | At the time, the student was practicing tennis with five other students and one coach at another court next to the one where the parachute was found. |
| | Note | Repeated words are sometimes substituted with pronouns or pro-verbs, such as "it" and "do so." |

Table 6: Examples of non-repetitive translations from Jiji Japanese–English news articles. "Ja" and "En" indicate the original parallel sentences from the articles. "En(trans)" indicates consistent translations by humans, which are listed for comparison.

---

[15] https://www2.statmt.org/wmt24/non-repetitive-translation-task.html

# B Detailed Statistics of the Human Evaluation Results

Table 7 shows detailed statistics of the human evaluation results.

| Model | | | Translation Style | | | Total |
|---|---|---|---|---|---|---|
| | | | <NON-REP> | <REP> | <INCORRECT> | |
| Waseda Riko | Adequacy (bin) | $[5, 5]$ | 127 | 20 | 0 | 147 |
| | | $[4, 5)$ | 280 | 17 | 3 | 300 |
| | | $[3, 4)$ | 15 | 1 | 7 | 23 |
| | | $[2, 3)$ | 0 | 0 | 0 | 0 |
| | | $[1, 2)$ | 0 | 0 | 0 | 0 |
| | | **Total** | **422** | 38 | 10 | 470 |
| SYSTRAN | Adequacy (bin) | $[5, 5]$ | 32 | 45 | 0 | 77 |
| | | $[4, 5)$ | 71 | 121 | 5 | 197 |
| | | $[3, 4)$ | 32 | 66 | 25 | 123 |
| | | $[2, 3)$ | 16 | 21 | 29 | 66 |
| | | $[1, 2)$ | 1 | 0 | 6 | 7 |
| | | **Total** | 152 | 253 | 65 | 470 |
| Baseline | Adequacy (bin) | $[5, 5]$ | 66 | 46 | 0 | 112 |
| | | $[4, 5)$ | 108 | 53 | 6 | 167 |
| | | $[3, 4)$ | 43 | 21 | 41 | 105 |
| | | $[2, 3)$ | 16 | 9 | 40 | 65 |
| | | $[1, 2)$ | 3 | 0 | 18 | 21 |
| | | **Total** | 236 | 129 | 105 | 470 |

Table 7: Statistics of the human evaluation results.

# A3-108 Controlling Token Generation in Low Resource Machine Translation Systems

**Saumitra Yadav Ananya Mukherjee Manish Shrivastava**
MT-NLP Lab
LTRC, KCIS
IIIT Hyderabad, India
`saumitra.yadav@research.iiit.ac.in`
`ananya.mukherjee@research.iiit.ac.in`
`m.shrivastava@iiit.ac.in`

## Abstract

Translating for languages with limited resources poses a persistent challenge due to the scarcity of high-quality training data. To enhance translation accuracy, we explored controlled generation mechanisms, focusing on the importance of control tokens. In our experiments, while training, we encoded the target sentence length as a control token to the source sentence, treating it as an additional feature for the source sentence. We developed various NMT models using transformer architecture and conducted experiments across 8 language directions (English ⟺ Assamese, Manipuri, Khasi, and Mizo), exploring four variations of length encoding mechanisms. Through comparative analysis against the baseline model, we submitted two systems for each language direction. We report our findings for the same in this work.

## 1 Introduction

Developing Machine Translation solutions for low-resource language pairs is one of the most interesting areas under the umbrella of Machine Translation. There have been many ways of adapting Machine Translation for low-resource language pairs, like,

- Using statistical models instead of neural-based ones to build a system. (Koehn and Knowles, 2017)

- Using multiple combinations of word segmentation to tackle data sparsity in this setting. (Sennrich et al., 2016b; Mujadia and Sharma, 2021; Yadav and Shrivastava, 2021)

- Using monolingual data to create synthetic bitext and train an improved system. (Sennrich et al., 2016a; Burchell et al., 2022; Fadaee et al., 2017)

- Using a pivot language as a bridge between high and low resource language pairs. (Kunchukuttan et al., 2017)

- Using transfer learning (Zoph et al., 2016) by transferring the knowledge from a high language pair setting to a related low language pair setting.

- Multilingual NMT extended on transfer learning by sharing learning space between multiple languages, with the goal of low-resource pair learning from the high-resource pair in a system with decent success. (Johnson et al., 2017)

For low-resource languages, the scarcity of high-quality, extensive datasets necessitates carefully utilising available resources. To maximize the extraction of information from these limited data, we plan to append the target length at the end of the source sentences. This approach draws inspiration from previous research, where incorporating the target length significantly enhanced performance in subtitle generation (Lakew et al., 2019) and current work is adapted from Fan et al. (2018) work on summarization.

In the current work, we consider target token length, length of target sentence **after** subword segmentation, as an additional feature for the source sentence. Intuition is that the system will learn to produce translations subjected to target length. There is an issue of accurately predicting the number of target language tokens in test cases or real-world scenarios. To predict target length, we used multiple methods,

- Neural network to predict target length given source sentence.

- Mean token length ratio of target to source sentence from validation set (Lakew et al., 2019) to predict target length given a source sentence.

- Sampling from a normal distribution, where the mean and standard deviation are calculated based on the ratios observed in the validation dataset.

- And for comparison, we also used the actual target length from the test set provided in Pal et al. (2023).

Systems for translating between English and the languages Assamese, Manipuri, Khasi, and Mizo (collectively referred to as IL in the rest of the paper) were developed in this study. It was observed that utilizing the average ratio of target-to-source token lengths from the validation set proved to be an effective method for obtaining control tokens for translation in a low-resource context.

We summarize the contribution of our work as follows.

- Using the number of tokens as a controlling token to improve system performance in a low-resource environment.

- Viable strategy to get control tokens for unseen data.

## 2 Related Work

Lakew et al. (2019) biased the output length with a transformer architecture using i) target-source length ratio and ii) enriching the transformer positional embedding with length information.

Fan et al. (2018) added the number of tokens to be generated in abstract summarization during training and observed an improvement in the ROUGE score. However, replicating the same for machine translation has been challenging. As Stahlberg (2020) noted, length information can be provided as additional input to the decoder network (Fan et al., 2018; Liu et al., 2018) at each time step as the number of remaining tokens (Kikuchi et al., 2016), or by modifying Transformer positional embeddings (Takase and Okazaki, 2019). Nonetheless, these methods are not directly applicable to machine translation due to the difficulty in accurately predicting translation length.

Additionally, Lakew et al. (2019) biased the output length with a transformer architecture using i) the target-source length ratio and ii) enriching the transformer positional embedding with length information.

## 3 Approach

This section describes our strategies for computing the control token number, datasets used for training and testing, model architecture, evaluation and systems submitted in shared task.

### 3.1 Control Tokens

Predicting target length accurately in machine translation remains a complex task, influenced by various factors such as language pair characteristics, sentence structure, and context. To address this challenge, we use a few straightforward heuristics to leverage insights from training and validation data to estimate control tokens effectively. These heuristics aim to give additional information about output length for generation to MT systems. Control tokens (CT) were generated using the following methods (Figure 1):

- **Actual Control Token** refers to the exact count of tokens in the target sentence, derived from a reference or gold standard.

- **Predicted Control Token** is obtained by training a transformer model to predict the number of target tokens given source sentences, where the model learns to estimate the length of the target sentence based on the features extracted from the source sentence. We did this to leverage the self-attention mechanism of the transformer to capture contextual dependencies effectively, making it suitable for tasks requiring an understanding of sentence structure and length prediction.

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V \tag{1}$$

- **Ratio Control Token** is the target-to-source token length ratio of the validation dataset for each language pair. Here, we utilize the relationship between the

Figure 1: Illustration of our Control Token Generation (CTG) approach

lengths of target sentences and their corresponding source sentences from the validation dataset. For $i^{th}$ source sentence, the control token (CT) is,

$$CT = R_{avg} * len_{source_i} \quad (2)$$

where $R_{avg} = \frac{\sum_{\forall j} len_{target_j}/len_{source_j}}{\text{number of sentence pairs}}$ and, $len_{sent}$ is length of token length of *sent*.

- **Sampled Control Token** is achieved by sampling from a normal distribution where the mean and standard deviation are derived from the ratios observed in the validation dataset.

$$CT = \mathcal{N}(R_{avg}, \sigma^2) * len_{source_i} \quad (3)$$

where $R_{avg} = \frac{\sum_{\forall j} len_{target_j}/len_{source_j}}{\text{number of sentence pairs}}$, $\sigma^2$ is standard deviation in ratios and, $len_{sent}$ is token length of sentence *sent*.

### 3.2 Datasets

We used the Dataset from Pal et al. (2023), Pakray et al. (2024) for English ⟺ Assamese, Manipuri, Khasi, and Mizo. Table 1 gives the statistics for each language pair and merge operations (mergeOps) used for the Byte Pair Encoding model of both source and target sentences (Sennrich et al., 2016b).

Figure 2 gives the distribution of IL sentence length with English sentence length ratio for all language pairs. Some sentences in training data have very high ratios compared to validation or test sets. This is where our method can induce learning correspondence between the number of tokens generated and the Control token.

| Language Pair | Train | Validation | Test | mergeOps |
|---|---|---|---|---|
| English Assamese | 50 K | 2000 | 2000 | 16K |
| English Mizo | 50K | 2000 | 1500 | 16K |
| English Khasi | 24K | 1000 | 1000 | 4K |
| English Manipuri | 21K | 1000 | 1000 | 16K |

Table 1: Dataset with merge operation for respective language pair

### 3.3 Architecture

For all the models, we trained machine translation models with the Transformers architecture(Vaswani et al., 2017) using fairseq(Ott et al., 2019) tool[1]. During the training, each source sentence was appended with a 'control token number', the count of target tokens.

## 4 Experiments

To select the systems as primary and contrastive output, we carried out experiments for English ⟺ Assamese, Khasi, Manipuri and Mizo and evaluated the translations of the test set from Pal et al. (2023) using lexical-based metrics, CHRF++ (Popović, 2017).

### 4.1 Results and Analysis

Table 2 summarises the performance of translation systems for EN-IL and IL-EN using CHRF++. We found statistically significant improvement in translation performance by adding a control token as an additional feature. We observed that,

- In most of the cases, scores improve when the Actual CT is added to the source.

---

[1]We used basic configuration of transformer architecture

730

Figure 2: Distribution of Sentence Length Ratio (IL/English) across Train, Validation, and Test Datasets. The orange line denotes the average ratio. The X-axis indicates the number of sentences, and the Y-axis depicts the Sentence Length Ratio (IL/English).

This is expected since it is a gold reference, ensuring the number of target tokens is precise.

- Predicting CT is a challenging problem, as mentioned earlier. While there were improvements in systems, there is also a significant drop in CHRF++ scores for English to Mizo and English to Khasi.

- Utilizing the ratio from validation to determine CT appears to be an optimal choice, which becomes more apparent

when examining the distribution of sentence length ratios of Validation and Test in Figure 2. Here, a clear similarity is observed between the two distributions regarding the range of sentence ratios. The inclusion of Ratio CT led to improved performance for English to IL, and vice-versa following actual CT.

- Sampled CT demonstrated strong performance for English to IL, but it did not exhibit the same level of effectiveness for

| Language Direction | Baseline | Actual CTG | Predicted CTG | Ratio CTG | Sampling CTG |
|---|---|---|---|---|---|
| EN-MZ | **38.11** | 37.72 | 31.44 | 36.37 | 37.23 |
| EN-MN | 30.47 | **32.27** | **31.96** | **32.21** | **32.11** |
| EN-KH | 31.69 | **35.35** | 26.07 | **34.28** | **35.53** |
| EN-AS | 18.44 | 18.62 | 18.63 | 17.67 | 17.49 |
| MZ-EN | 31.03 | **32.77** | 29.25 | **31.5** | **32.68** |
| MN-EN | **35** | 34.24 | 32.62 | 34.03 | 34.45 |
| KH-EN | 26.23 | **27.58** | 26.42 | **27.09** | **27.39** |
| AS-EN | 22.76 | **24.12** | **23.76** | **23.98** | 22.79 |

Table 2: CHRF++ scores of EN-IL and IL-EN Translation system. Scores in Bold are statistically significant improvements compared to the baseline scores with p<0.05.

IL to English. Despite this, it performed comparably well to Ratio CTG in the English to IL direction.

We further analyzed the target-to-source length ratio for EN-IL direction for all 4 language pairs in Figure 3. Examining the average sentence length ratio between Train (●) and baseline systems (■) in comparison to the reference length ratio (▲), sheds light on the behaviour of baseline systems and highlights the advantage of employing CT. In the case of English to Manipuri and Khasi, where significant improvements in CHRF++ scores were noted, the length ratio for baseline systems fell short of the test set. Conversely, when considering the Ratio CTG (◆), we observe their proximity to the Reference Ratio. This supports the idea of using control tokens as an additional feature in the source sentence. It also explains the impact of a poorer prediction system; as seen in the English-to-Mizo ratio, which overshoots by a large margin, there is also a significant drop in translation performance.

Based on these observations, we conclude that if the target sentence length is predictable, leveraging it as an additional feature with the source sentence proves to be a great choice for training a translation model in a low-resource setting.

### 4.2 Submission

For Translation submission, we preprocessed unseen testset shared by Organizers and submitted translations from the following two systems,

- Primary System: is a model trained using transformer architecture with source sen-

| Language Direction | System | TER | RIBES | METEOR | ChrF |
|---|---|---|---|---|---|
| English to Assamese | Baseline | 100.46 | 0.0347 | 0.0587 | 0.1817 |
| | Ratio | 99.79 | 0.0243 | 0.05134 | 0.1773 |
| English to Manipuri | Baseline | 101.73 | 0.0084 | 0.0179 | 0.1401 |
| | Ratio | 101.55 | 0.0072 | 0.0166 | 0.1415 |
| English to Mizo | Baseline | 92.32 | 0.0406 | 0.0978 | 0.18 |
| | Ratio | 92.84 | 0.0328 | 0.0906 | 0.173 |
| English to Khasi | Baseline | 92.92 | 0.087 | 0.1209 | 0.1905 |
| | Ratio | 87.69 | 0.0873 | 0.1589 | 0.2296 |
| Assamese to English | Baseline | 96.44 | 0.0378 | 0.0677 | 0.1803 |
| | Ratio | 96.19 | 0.0322 | 0.0671 | 0.1883 |
| Manipuri to English | Baseline | 96.45 | 0.029 | 0.0615 | 0.1865 |
| | Ratio | 96.5 | 0.0271 | 0.0635 | 0.1889 |
| Mizo to English | Baseline | 97.75 | 0.0195 | 0.0544 | 0.1633 |
| | Ratio | 96.18 | 0.0181 | 0.0587 | 0.1826 |
| Khasi to English | Baseline | 105.76 | 0.0094 | 0.0403 | 0.1358 |
| | Ratio | 107.7 | 0.0071 | 0.0359 | 0.1348 |

Table 3: Performance on Unseen Testset

tence and target output length predicted using average **Ratio** of source and target sentences in the validation dataset.

- Contrastive System: is a model trained using transformer architecture without adding CT (Baseline).

## 5 Performance on Unseen Testset

Despite the promising results in test sets with training datasets, on the Unseen test set (3) provided by the shared task organizer (Pakray et al., 2024), our approach only gave a slight increase in score compared to the baseline in English to Manipuri, English to Khasi, Assamese to English, Manipuri to English and Mizo to English.

## 6 Conclusion and Future Work

We address the challenge of translating languages with limited resources by enhancing translation accuracy using target sentence length as an additional feature in the source sentence. We experimented using transformer architecture across 8 language directions (English ⟺ Assamese, Manipuri, Khasi, and Mizo). Evaluation against baseline models on a shared test set revealed that our approach significantly improves translation quality in some language directions, demonstrating its effectiveness in improving translation for low-resource languages. However, for the unseen dataset, even though there was an improvement, it wasn't that huge. Overall, we also found that the baseline systems themselves were not promising. Hence, we would be repli-

Figure 3: Distribution of Average Sentence length Ratio (IL/English) for Train, Validation and Test Dataset for all language pairs in English to IL direction.

cating this work with other datasets and language pairs to check the validity of this outcome.

# References

Laurie Burchell, Alexandra Birch, and Kenneth Heafield. 2022. Exploring diversity in back translation for low-resource machine translation. In *Proceedings of the Third Workshop on Deep Learning for Low-Resource Natural Language Processing*, pages 67–79, Hybrid. Association for Computational Linguistics.

Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2017. Data augmentation for low-resource neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 567–573, Vancouver, Canada. Association for Computational Linguistics.

Angela Fan, David Grangier, and Michael Auli. 2018. Controllable abstractive summarization. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 45–54, Melbourne, Australia. Association for Computational Linguistics.

Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.

Yuta Kikuchi, Graham Neubig, Ryohei Sasano, Hiroya Takamura, and Manabu Okumura. 2016. Controlling output length in neural encoder-decoders. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1328–1338, Austin, Texas. Association for Computational Linguistics.

Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.

Anoop Kunchukuttan, Maulik Shah, Pradyot Prakash, and Pushpak Bhattacharyya. 2017. Utilizing lexical similarity between related, low-resource languages for pivot-based SMT. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 283–289, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Surafel Melaku Lakew, Mattia Di Gangi, and Marcello Federico. 2019. Controlling the output length of neural machine translation. In *Proceedings of the 16th International Conference on Spoken Language Translation*, Hong Kong. Association for Computational Linguistics.

Yizhu Liu, Zhiyi Luo, and Kenny Zhu. 2018. Controlling length in abstractive summarization using a convolutional neural network. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4110–4119, Brussels, Belgium. Association for Computational Linguistics.

Vandan Mujadia and Dipti Misra Sharma. 2021. English-Marathi neural machine translation for LoResMT 2021. In *Proceedings of the 4th Workshop on Technologies for MT of Low Resource Languages (LoResMT2021)*, pages 151–157, Virtual. Association for Machine Translation in the Americas.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Partha Pakray, Santanu Pal, Advaitha Vetagiri, Reddi Mohana Krishna, Arnab Kumar Maji, Sandeep Kumar Dash, Lenin Laitonjam, Lyngdoh Sarah, and Riyanka Manna. 2024. Findings of wmt 2024 shared task on low-resource indic languages translation. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the WMT 2023 shared task on low-resource Indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694, Singapore. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016b. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Felix Stahlberg. 2020. Neural machine translation: A review. *Journal of Artificial Intelligence Research*, 69:343–418.

Sho Takase and Naoaki Okazaki. 2019. Positional encoding to control output sequence length. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3999–4004, Minneapolis, Minnesota. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Saumitra Yadav and Manish Shrivastava. 2021. A3-108 machine translation system for LoResMT shared task @MT summit 2021 conference. In *Proceedings of the 4th Workshop on Technologies for MT of Low Resource Languages (LoResMT2021)*, pages 124–128, Virtual. Association for Machine Translation in the Americas.

Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1568–1575, Austin, Texas. Association for Computational Linguistics.

# Samsung R&D Institute Philippines @ WMT 2024 Indic MT Task

**Matthew Theodore Roque**[a]    **Carlos Rafael Catalan**[a]    **Dan John Velasco**[a]
**Manuel Antonio Rufino**[a]    **Jan Christian Blaise Cruz**[a,b]
[a]Samsung R&D Institute Philippines    [b]MBZUAI
{roque.mt,c.catalan,dj.velasco,ma.rufino}@samsung.com
jan.cruz@mbzuai.ac.ae

## Abstract

This paper presents the methodology developed by the Samsung R&D Institute Philippines (SRPH) Language Intelligence Team (LIT) for the WMT 2024 Shared Task on Low-Resource Indic Language Translation. We trained standard sequence-to-sequence Transformer models from scratch for both English-to-Indic and Indic-to-English translation directions. Additionally, we explored data augmentation through backtranslation and the application of noisy channel reranking to improve translation quality. A multilingual model trained across all language pairs was also investigated. Our results demonstrate the effectiveness of the multilingual model, with significant performance improvements observed in most language pairs, highlighting the potential of shared language representations in low-resource translation scenarios.

## 1   Introduction

This paper details our primary submission for the WMT 2024 Shared Task on Low-Resource Indic Language Translation. Our submission covers the following language pairs: English ↔ Assamese (en-as), English ↔ Mizo (en-mz), English ↔ Khasi (en-kh), and English ↔ Manipuri (en-mn). Our approach builds upon the methodology used in Samsung R&D Philippines' WMT23 entry (Cruz, 2023). We employed a standard sequence-to-sequence Transformer architecture (Vaswani et al., 2023), combined with data augmentation through backtranslation (Sennrich et al., 2016), noisy channel reranking (Yee et al., 2019), and additionally experiment with a multilingual model trained on all language pairs.

---

[b]Work done while at Samsung R&D Institute Philippines

## 2   Methodology

### 2.1   Environment

For preprocessing, training, and generation, we utilized PyTorch 2.0 and fairseq 0.12.2. All training was conducted on NVIDIA P100 GPUs.

### 2.2   Data Analysis

We used the Indic dataset provided from WMT 2023 for all language pairs. First, we conducted an exploratory data analysis for all the languages to see if there were noteworthy patterns that could guide us in our translation in the Indic and English languages. We used various methods in this data analysis such as finding N-most common words, generating N-grams, and histograms of lengths of sentences.

An interesting pattern emerged when generating the histograms of sentence lengths as seen in Figure 1. For the English-Mizo pair, the distributions almost completely overlap. However, for the English-Assamese, English-Khasi, and English-Manipuri pairs, the Indic languages generally exhibit slightly longer sequences. We hypothesize that these longer sequences may cause translation errors in the Indic to English language directions. The models might be driven to provide translations that are driven more by length alignment, and so may attempt to fill in additional tokens to produce longer sequences even if it may not necessarily be semantically accurate.

### 2.3   Data Preprocessing

We exclusively used the task dataset for all language pairs. For the parallel data, we first removed exact duplicates, then detokenized the text to correct spacing around punctuation. The statistics of parallel data are summarized in Table 1. Following this, we trained a BPE tokenizer (Sennrich et al., 2015), applied BPE tokenization, and binarized the data for use with fairseq. Each language pair

| source↔target | Pairs | Words (source) | Words (target) | Vocab Size |
|---|---|---|---|---|
| en↔as | 50,000 | 969,626 | 825,063 | 31,448 |
| en↔kh | 21,000 | 729,930 | 875,545 | 9,312 |
| en↔mz | 50,000 | 981,468 | 1,062,414 | 30,432 |
| en↔mn | 21,687 | 390,730 | 330,319 | 30,736 |

Table 1: Statistics of parallel training data. Note that "Words" refers to word count estimated using the wc command on the plaintext files.

| source→target | Unfiltered Pairs | Filtered Pairs | Words (source) | Words (target) |
|---|---|---|---|---|
| en→as | 2,624,715 | 279,956 | 3,200,053 | 3,444,809 |
| en→mz | 1,900,848 | 1,637,838 | 21,534359 | 26,367139 |
| en→kh | 160,128 | 19,358 | 363,441 | 490,257 |
| en→mn | 298,608 | 10,837 | 97,418 | 145,928 |

Table 2: Statistics of generated backtranslated parallel data. Note that "Words" refers to word count estimated using the wc command on the plaintext files.

has a shared vocabulary between English and the respective Indic language. The preprocessed parallel data was used to train our translation models. The same preprocessing steps were applied to the monolingual data for training the language models. As no monolingual data was provided for English, we used the combined English sides of the parallel data to train the English language model.

## 2.4 Augmenting Data with Backtranslation

Due to time and data constraints, data augmentation via backtranslation was applied only in the English-to-Indic direction. Backtranslated data was generated by translating the monolingual Indic data into English using the trained Indic-to-English models. After generating the backtranslations, we applied ratio-based filters (Cruz, 2023) to remove low-quality parallel data, filtering based on sentence length, token length, character-to-token ratio, pair token ratio, and pair length ratio. For more details, please refer to the original paper. The dataset statistics for the backtranslated data are presented in Table 2.

## 2.5 Model Training

For each of the four language pairs, we trained four models: two **Translation Models**, one for each translation direction, and two **Language Models**, one for each language. The specifics of these models are described in the following subsections. Three of these four models were combined for noisy channel reranking in one direction, as detailed in Section 2.7. Additionally, we experimented with a **Multilingual Model** using the same

architecture as our translation models, but trained across all language pairs.

### 2.5.1 Translation Models

For the translation models (English→Indic, Indic→English), we trained encoder-decoder Transformer architectures (Vaswani et al., 2023) from scratch using parallel data. Separate models were trained for each language pair and for each translation direction. We used the large variant of the Transformer model with 213M parameters, training for 100,000 steps, with the first 10,000 being warmup steps (Gotmare et al., 2018), with a maximum of 8,000 tokens per step. The learning rates varied across language directions, as follows: en→as (9e-5), en→kh (5e-4), en→mizo (9e-5), en→mn (9e-5), as→en (5e-4), kh→en (5e-4), mizo→en (5e-4), and mn→en (5e-4). All other hyperparameters are detailed in Table 3.

These translation models were not only used as direct translation models but also served as channel translation models for noisy channel reranking, further discussed in Section 2.7.

### 2.5.2 Language Models

We trained monolingual language models for each language from scratch using the decoder-only component of the original Transformer architecture, as described by (Vaswani et al., 2023). We used the base variant of the Transformer, which contains 65M parameters. For the Indic language models (Assamese, Mizo, Khasi, Manipuri), we trained on the provided monolingual data. For the English language model, we concatenated the English side of the parallel data for training.

Figure 1: Histogram of Sentence Lengths

All models were trained using the Adam optimizer (Kingma and Ba, 2017) with $\beta_1 = 0.90$ and $\beta_2 = 0.98$. Training was conducted for a maximum of 100,000 steps, with the first 10,000 steps as a warmup (Gotmare et al., 2018). The learning rate started at 1e-7, peaked at 5e-4, and decayed following an inverse square root learning rate schedule. The batch size was set to 32,000 tokens, and a dropout rate of 0.1 was applied. These models were later used in noisy channel reranking, as detailed in Section 2.7.

## 2.6 Multilingual Model

We trained a large variant of the Transformer model with 213M parameters on all four language pairs, in both the English-to-Indic and Indic-to-English directions, following the approach of last year's entries (Zhang, 2023). Given the low-resource nature of each individual pair, we aimed to enable the language pairs to leverage cross-linguistic knowledge (Aharoni et al., 2019). The training process spanned 50,000 steps, with the first 5,000 steps serving as warmup (Gotmare et al., 2018; Neubig and Hu, 2018). We used 8 P100 GPUs for a maximum of 51,200 tokens per step and a learning rate of 1e-4. The remaining hyperparameters were consistent with those used in the other translation models as shown in Table 3.

Curriculum learning has been shown to improve generalization by introducing tasks progressively, allowing the model to build on prior knowledge (Wang et al., 2019). For our multilingual translation model, we aimed to apply a form of curriculum learning by training on different language pairs one at a time. We prepended source and target language tokens and trained the model sequentially on one language pair at a time. This structured training approach, inspired by Bengio et al. (2009), could help the model learn each language faster and transfer learned knowledge across language pairs. Similar to the benefits seen in multi-task learning by Niehues and Cho (2017), we hypothesized that this sequential training will enhance the model's ability to share representations across languages, ultimately leading to improved performance.

## 2.7 Noisy-Channel Reranking (NCR)

We experimented with Noisy Channel Reranking (Yee et al., 2019) to reevaluate and improve the translations. For brevity, we refer to this as NCR. This method utilizes three different models: a direct translation model (source→target), a channel model (target→source), and a monolingual language model (target only). These models are combined to rescore each candidate translation token during beam search decoding. The score for a candidate token $\hat{y}_i^{(T)}$ at timestep $T$ is recomputed using the linear combination of the outputs from all three models:

| Training Hyperparameters | |
|---|---|
| Vocab Size | 31,960 |
| Tied Weights | Yes |
| Dropout | 0.3 |
| Attention Dropout | 0.1 |
| Weight Decay | 0.0 |
| Label Smoothing | 0.1 |
| Optimizer | Adam |
| Adam Betas | $\beta_1$=0.90, $\beta_2$=0.98 |
| Adam $\epsilon$ | $\epsilon$=1e-6 |
| LR Schedule | Inverse Sqrt |
| Batch Size | 8,000 tokens |

Table 3: Fixed hyperparameters for direct translation models.

$$P(\hat{y}_i^{(T)}|x;\hat{y}^{(T-1)})' = \frac{1}{t}log(P(y|\hat{x}^{(T-1)})$$
$$+\frac{1}{s}[\delta_{ch}log(P(x|\hat{y}^{(T-1)}) \quad (1)$$
$$+\delta_{lm}log(P(\hat{y}^{(T-1)}))]$$

Here, $t$ represents the length of the target sentence $y$, and $s$ represents the length of the source sentence $x$, both of which serve as debiasing terms. The weights $\delta_{ch}$ and $\delta_{lm}$ control the influence of the channel model and the language model, respectively, on the final score.

## 2.8 Decoding and Noisy-channel Reranking Hyperparameter Tuning

We determined the optimal length penalty values by sweeping across four values: 0.5, 1.0, 1.5, and 2.0. This was done for each language direction, and the length penalty that resulted in the highest BLEU score on the provided test data was selected. The optimal length penalties for each direction are as follows: en→as (1.5), en→kh (2.0), en→mizo (1.0), en→mn (1.5), as→en (2.0), kh→en (1.5), mizo→en (0.5), and mn→en (2.0). These values were then used to tune the channel and language model weights for NCR.

We applied a similar approach to find the optimal values for the channel weight, $\delta_{ch}$, and the language model weight, $\delta_{lm}$. For the English-to-Indic models, we fixed $\delta_{ch}$ at 0.1 and varied $\delta_{lm}$ across 0.1, 0.2, 0.3, 0.4, 0.5, and 0.6. For the Indic-to-English models, we reversed the setup by fixing $\delta_{lm}$ at 0.1 and varying $\delta_{ch}$ across the same values.

These configurations were chosen due to time constraints, which limited our ability to perform

more exhaustive evaluations of various combinations of $\delta_{lm}$ and $\delta_{ch}$. Additionally, based on our initial evaluation of the translation models, the English-to-Indic models exhibited stronger performance, so we focused on evaluating their impact as channel models in NCR. Conversely, the Indic language models were trained on significantly more data than the English language model, making it essential to assess their influence in the reranking process.

## 3 Results and Discussion

In this section, we present the results of our experiments and provide an in-depth discussion of our findings. Tables 4 and 5 summarize the BLEU and chrF scores for each model and method on last year's test set, respectively. Table 6 summarizes the hyperparameters used for training the models.

### 3.1 Baseline

Our baseline consists of standard Transformer models trained from scratch for each language pair, without any backtranslation, length penalty tuning, noisy channel reranking, or multilingual setup. These models were trained using the parallel data provided, with a shared BPE vocabulary between English and each respective Indic language.

For the English-Khasi pair in particular, we set the target vocabulary size to 10,000, while for the other three language pairs, we retained a target of 32,000. Initially, we aimed for a 32,000 vocabulary size across all language pairs, but English-Khasi's vocabulary only reached approximately 20,000. Given that this was our worst-performing pair, we reduced the target size to 10,000, resulting in a BLEU score improvement of about 3 points.

As shown in Table 4, the baseline models performed adequately for most language pairs, with BLEU scores ranging from 4.2 (Khasi→English) to 34.1 (English→Manipuri). Notably, the English-to-Indic models generally outperformed the Indic-to-English models across all language pairs.

### 3.2 Data Augmentation Using Backtranslation

The number of pairs in the backtranslated data, as shown in Table 2, was greatly reduced after filtering. This reduction most likely stems from the poor performance of the Indic-to-English models used for backtranslation. These models may have produced low-quality translations, leading to a substantial number of backtranslated pairs being dis-

| source→target | BLEU Scores | | | | |
| --- | --- | --- | --- | --- | --- |
| | Baseline | w/ BT Data | Tuned lenpen | NCR | Multilingual |
| en→as | 13.8 | 3.0 | **14.0** | **14.0** | 3.5 |
| as→en | 9.5 | - | 9.8 | 9.8 | **10.1** |
| en→mz | 29.7 | 18.6 | 29.7 | 28.5 | **30.2** |
| mz→en | 19.9 | - | **21.3** | 19.2 | 19.1 |
| en→kh | 8.8 | 6.5 | 9.6 | 10.2 | **16.1** |
| kh→en | 4.2 | - | 4.4 | 4.3 | **7.9** |
| en→mn | 34.1 | 1.1 | 35.2 | 34.9 | **36.4** |
| mn→en | 16.7 | - | 17.0 | 17.0 | **21.4** |

Table 4: BLEU Scores on WMT2023 Indic MT test data. The use of **BT Data** (training on backtranslated data) showed a decline in performance. The **tuned lenpen** (length penalty) generally improves BLEU score while **NCR** (Noisy Channel Reranking) yielded mixed results. The multilingual setting outperforms all other settings in all language pairs except en→as, where tuned lenpen and NCR showed the same score, and mz→en, where tuned lenpen was best.

| source→target | chrF Scores | | | | |
| --- | --- | --- | --- | --- | --- |
| | Baseline | w/ BT Data | Tuned lenpen | NCR | Multilingual |
| en→as | **26.1** | 14.8 | 25.2 | 25.2 | 6.9 |
| as→en | 27 | - | 26.3 | 26.8 | **28.4** |
| en→mz | 44.4 | 34.5 | 44.4 | 42.8 | **45.2** |
| mz→en | 34.4 | - | 35.4 | 34.1 | **35.6** |
| en→kh | 29.9 | 27.9 | 30 | 30.4 | **34.7** |
| kh→en | 23.7 | - | 23.4 | 23.3 | **27.8** |
| en→mn | 45 | 11 | 43.8 | 43.5 | **45.1** |
| mn→en | 35.3 | - | 34.2 | 34.7 | **44.2** |

Table 5: chrF Scores from the WMT2023 Indic MT test data. The use of **BT Data** (training on backtranslated data) showed a decline in performance. The **tuned lenpen** (length penalty) and **NCR** (Noisy Channel Reranking) was tuned for the BLEU scores and yielded mixed results for chrF. The multilingual setting outperformed all other settings in all language pairs except en→as, where the baseline was best.

carded during the filtering process. The pairs that remained after filtering likely were still not of the best quality, which diminished the overall quality of the training. As a result, the models trained on this backtranslated data performed worse, as reflected in their BLEU scores in Table 4 and their chrF scores in Table 5.

### 3.3 Length Penalty

Our tuning of the length penalty, as shown in Table 6, revealed that most language directions, with the exception of English-to-Mizo and Mizo-to-English, preferred shorter translation sequences. As shown in Figure 1, the distribution of sentence lengths across the language pairs indicates a reasonable amount of overlap, though the Indic languages tend to have slightly longer sequences.

This preference for shorter sequences coincides with a known issue in Neural Machine Translation (NMT) models when handling long input se-

quences. NMT models typically rely on absolute positional encodings, which use fixed sine and cosine functions to assign vector positions. This approach tends to struggle with longer sequences due to the limitations of these fixed encodings, resulting in less precise representations as sentence length increases (Neishi and Yoshinaga, 2019). This is likely contributing to the models' difficulty in generating coherent longer translations, particularly for underperforming language pairs like English-Assamese and English-Khasi. As sequence length increases, the models are more prone to generating irrelevant or erroneous tokens, leading to a degradation in translation quality.

It is interesting to note that despite the Indic languages generally having longer sequences, a length penalty greater than one was found to be optimal for both directions, even in English-to-Indic translation. This indicates that the models may be biased

| source→target | Hyperparameters | | |
|---|---|---|---|
| | lenpen | ch_wt | lm_wt |
| en→as | 1.5 | 0.1 | 0.2 |
| as→en | 2.0 | 0.6 | 0.1 |
| en→mz | 1.0 | 0.1 | 0.1 |
| mz→en | 0.5 | 0.4 | 0.1 |
| en→kh | 2.0 | 0.1 | 0.1 |
| kh→en | 1.5 | 0.2 | 0.1 |
| en→mn | 1.5 | 0.1 | 0.1 |
| mn→en | 2.0 | 0.3 | 0.1 |

Table 6: Final length penalty (lenpen), channel model weight (ch_wt), and language model weight (lm_wt).

towards shorter outputs across most language pairs, potentially as a safeguard against these positional encoding limitations. While this behavior aligns with our expectations for the English-Assamese pair based on its performance, the similar tendencies in the English-Khasi pair were more surprising, given the closer alignment of sentence lengths between these languages.

### 3.4 Noisy Channel Reranking

The BLEU scores obtained with NCR, as shown in Table 4, yielded mixed results. After tuning the length penalty, we observed that NCR improved performance for only one model out of eight, specifically English-to-Khasi. The chrF scores, as shown in Table 5, also indicate slightly improved performance with NCR solely for the English-to-Khasi pair. For all other language pairs, there was either no change in BLEU and chrF scores or a slight decrease. It is crucial to highlight that these results reflect the best combination of hyperparameters we identified; alternative hyperparameter settings would have resulted in even more pronounced variations in scores.

One notable finding is that the optimal language model weight was consistently around 0.1 across most language pairs. This suggests that the language model contributed minimally to improving translation quality. This issue may stem from either data quality or data quantity limitations. Investigating data quality issues would be valuable, but addressing them poses a significant challenge due to the already low-resource nature of the Indic languages. Further filtering could exacerbate data scarcity, making it difficult to maintain sufficient training data.

Conversely, the channel model weights were found to be more effective, with optimal values

varying by language pair but generally falling in the mid-range. For the best-performing Indic-to-English pairs with NCR, specifically Mizo-to-English and Manipuri-to-English, the channel model weights were 0.4 and 0.3, respectively. These language pairs also had the best direct translation models and channel models, suggesting a stronger alignment between model quality and channel model effectiveness for these particular languages.

### 3.5 Multilingual Model

The multilingual model trained on all language pairs demonstrated considerable improvements over the baseline models, achieving the best performance in 6 out of the 8 language pairs. We attribute this success to the model's ability to learn from a broader context across all five languages, allowing for the creation of shared language representations. This approach is especially beneficial given the small size of the training datasets, as the multilingual model can leverage cross-linguistic knowledge to enhance translation quality.

However, due to time constraints, we were unable to explore the potential of using the multilingual model as a channel model within NCR. This remains a promising avenue for future research. Further studies could also investigate pre-training on the available monolingual data before fine-tuning for translation tasks. Additionally, fine-tuning the multilingual model for language modeling could further improve its utility in NCR, potentially acting out all three functions in NCR, leveraging shared linguistic knowledge on all languages and tasks, enhancing performance in low-resource language pairs.

## 4 Conclusion

In this paper, we presented our approach to the WMT 2024 Shared Task on Low-Resource Indic Language Translation. Our experiments demonstrated that the multilingual model trained across all language pairs performed exceptionally well, particularly in comparison to the baseline models, achieving the highest BLEU scores in 6 out of 8 language pairs and the highest chrF scores in 7 out of 8 language pairs. This indicates that leveraging shared language representations, especially when dealing with small datasets, can significantly enhance translation performance by utilizing cross-linguistic knowledge.

Despite some success, our attempts to improve results through data augmentation using backtranslation and noisy channel reranking yielded mixed outcomes. The poor quality of the Indic-to-English backtranslated data led to performance degradation, emphasizing the importance of both data quality and quantity in low-resource scenarios. Additionally, while noisy channel reranking provided benefits in isolated cases, its overall impact was limited, potentially due to suboptimal language model and channel model contributions.

The promising performance of our multilingual model suggests that further research could explore its integration within noisy channel reranking, possibly utilizing it as both a translation and a channel model. Additionally, future work should focus on enhancing the quality of backtranslated data and investigating pre-training strategies on monolingual data to boost the performance of low-resource language pairs.

# References

Roee Aharoni, Melvin Johnson, and Orhan Firat. 2019. Massively multilingual neural machine translation. *arXiv preprint arXiv:1903.00089*.

Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48.

Jan Christian Blaise Cruz. 2023. Samsung R&D institute Philippines at WMT 2023. In *Proceedings of the Eighth Conference on Machine Translation*, pages 103–109, Singapore. Association for Computational Linguistics.

Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation. *arXiv preprint arXiv:1810.13243*.

Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization.

Masato Neishi and Naoki Yoshinaga. 2019. On the relation between position information and sentence length in neural machine translation. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 328–338, Hong Kong, China. Association for Computational Linguistics.

Graham Neubig and Junjie Hu. 2018. Rapid adaptation of neural machine translation to new languages. *arXiv preprint arXiv:1808.04189*.

Jan Niehues and Eunah Cho. 2017. Exploiting linguistic resources for neural machine translation using multi-task learning. *arXiv preprint arXiv:1708.00993*.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need.

Wei Wang, Ye Tian, Jiquan Ngiam, Yinfei Yang, Isaac Caswell, and Zarana Parekh. 2019. Learning a multi-task curriculum for neural machine translation. *ArXiv, abs/1908.10940*.

Kyra Yee, Yann Dauphin, and Michael Auli. 2019. Simple and effective noisy channel modeling for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5696–5701, Hong Kong, China. Association for Computational Linguistics.

Wenbo Zhang. 2023. Iol research machine translation systems for wmt23 low-resource indic language translation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 978–982.

# DLUT-NLP Machine Translation Systems for WMT24 Low-Resource Indic Language Translation

**Chenfei Ju[1], Junpeng Liu[1], Kaiyu Huang[2]** and **Degen Huang[1]**
[1]Dalian University of Technology, Dalian, China
[2]Beijing Jiaotong University, Beijing, China
{845110184, liujunpeng_nlp}@mail.dlut.edu.cn, kyhuang@bjtu.edu.cn
huangdg@dlut.edu.cn

## Abstract

This paper describes the submission systems of DLUT-NLP team for the WMT24 low-resource Indic language translation shared task. We participated in the translation task of four language pairs, including en↔as, en↔mz, en↔kha, en↔mni. We used a transformer-based neural network architecture to train the model. Our system used the following methods: First, data processing was performed, and then we used monolingual data for pre-training. Next, we used parallel data for fine-tuning to obtain a multilingual translation model, and then we used this model for back-translation. We merged the back-translated data with the official parallel data and used the upsampling method to train a multilingual translation model from scratch. In order to improve the translation ability of the model for each translation direction, we fine-tuned the model for each language pair and used model averaging to obtain the best model for each language pair. Finally, we used $k$NN-MT and established a datastore using the official parallel data to assist translation in the inference stage. Experimental results show that our method greatly improves the BLEU scores for translation of these four language pairs.

## 1 Introduction

This paper introduces our system for WMT24 low-resource Indic language translation shared task. We participated in 4 language pairs, including English↔Assamese (en↔as), English↔Mizo (en↔mz), English↔Khasi (en↔kha) and English↔Manipuri (en↔mni).

The main methods used by our system are denoising language model pre-training (Lample and Conneau, 2019; Song et al., 2019; Lewis et al., 2020), back-translation (Sennrich et al., 2016a) and $k$NN-MT (Khandelwal et al., 2020). Neural machine translation is the first choice for machine translation systems nowadays, but it requires a large amount of parallel data. Therefore, low-resource translation is a major challenge due to its lack of data. In this task, the organizers provided a large amount of monolingual data in addition to a small amount of parallel data. So we considered using some pre-training methods to improve the performance of the model. At the same time, back-translation is a commonly used method in the field of machine translation, which is effective in many scenarios. Therefore, we used the back translation method to obtain pseudo-parallel data to train a strong baseline model. To obtain the best model for each translation direction, we fine-tuned the baseline model for each language pair using the official parallel data. During this process, we used model averaging technology to improve the translation quality of the model. In addition to parametric methods, a large number of non-parametric methods have recently emerged to help models generate translations. We adopted the $k$NN-MT method and built a datastore for each translation direction to assist the model in the inference phase.

The rest of the paper is organized as follows: In Section 2 we describe our data processing methods; In Section 3 we describe the implementation process of our translation systems; In Section 4, we describe the experimental settings; In Section 5, we discuss about the results; Finally, in Section 6, the conclusion is drawn.

## 2 Data

For bilingual data, we only used official bilingual data. For monolingual data, in addition to the official monolingual data for Assamese, Mizo, Khasi and Manipuri (Pal et al., 2023; Pakray et al., 2024), we obtained English monolingual data from the WMT24 general task. Specifically, we used the English side of bilingual data (English↔German) in the WMT24 general task as English monolingual data.The statistics of the dataset is shown in Table 1.

| | as | kha | mni | mz | en |
|---|---|---|---|---|---|
| train (mono) | 2.6M | 0.2M | 2.1M | 1.9M | 2.5M |
| train (para) | 50k | 24k | 22k | 50k | - |
| dev | 2k | 1k | 1k | 1.5k | - |
| test | 2k | 1k | 1k | 2k | - |

Table 1: The number of sentences in the training, dev and test sets.

Since the quality of official data is relatively high, we did not perform additional preprocessing. For the English monolingual data, we performed some additional preprocessing steps. During pre-processing, we deleted sentences that were too long or repeated. And then we filtered out sentences in other languages by applying language identification. Finally we used an n-gram language model trained with KenLM (Heafield, 2011)[1] to calculate the perplexity of English monolingual data and removed sentences with high perplexity (>7,000). We used the Sentencepiece (Kudo and Richardson, 2018) tool to train a multilingual BPE (Sennrich et al., 2016b) model for subword segmentation. The training data includes all the parallel training data and monolingual data. The vocabulary size is set to 32,000.

## 3 System Overview

### 3.1 Pre-training

Using monolingual data for pre-training tasks is an effective solution for low-resource situations (Raffel et al., 2020). To this end, we first performed BART-style pre-training (Lewis et al., 2020) with all the available monolingual data and then fine-tuned the pretrained model with bilingual data. Following Lewis et al. (2020), we masked words with a probability of 0.15 and we randomly swapped words in the input sentences with a probability of 0.5.

After pre-training, we used all the bilingual data to fine-tune the pre-trained model. The bilingual data contains 4 language pairs in 8 translation directions.

### 3.2 Back-translation

To improve our translation pipeline, we explored the integration of back-translation as a potential enhancement. Back-translation involves using a trained model to translate from the target language back to the source language, effectively creating a

synthetic parallel dataset. We used the approach inspired by Sennrich et al. (2016a) to generate pseudo-parallel corpus.

Specifically, we used the model fine-tuned in the pre-training phase. We used this model to translate all non-English monolingual data into English as pseudo-parallel data. Then we mixed all the pseudo-parallel data with the official bilingual data. We used this data to train a multilingual translation model from scratch. During training, we used up-sampling method and the official parallel data was upsampled until it reached to a ratio of 1:1 with the synthetic data.

### 3.3 Language-specific Fine-tuning

Although multilingual translation models have made great progress, there is still the problem of inconsistent convergence of different language pairs in joint training (Wu et al., 2021; Huang et al., 2022). That is, different language pairs reach convergence in various training stages. We hope to get the best model for each language pair. Due to the low quality of pseudo-parallel data, we used the official bilingual data of each language pair to fine-tune the model trained using pseudo-parallel data.

During fine-tuning, we used the model averaging technology. Through model averaging, we combined the advantages of various models into a unified translation model. This process can not only improve the stability of the translation output, but also help improve the overall translation quality. We kept the three models with the lowest loss on the validation set for each language pair. We then used these three models to get the best model for each language pair.

### 3.4 $k$NN-MT

Non-parametric, $k$-nearest-neighbor algorithms have recently made inroads to assist generative models such as language models and machine translation decoders. Khandelwal et al. (2020) introduced $k$-nearest-neighbor machine translation

---

[1] https://github.com/kpu/kenlm

($k$NN-MT): a simple non-parametric method for machine translation via nearest-neighbor retrievals was proposed and has been verified its effectiveness. According to his method, we constructed a datastore to store the translation examples to be accessed during decoding with the official parallel data. When decoding, we used the current translation context to retrieve the $k$-nearest-neighbors in the datastore. Let $\boldsymbol{x} = \left(x_1, \ldots, x_{|\boldsymbol{x}|}\right) \in \mathcal{V}_X^{|\boldsymbol{x}|}$ and $\boldsymbol{y} = \left(y_1, \ldots, y_{|\boldsymbol{y}|}\right) \in \mathcal{V}_Y^{|\boldsymbol{y}|}$ denote a source sentence and target sentence, respectively, where $|\cdot|$ represents the length of the sentence, and $\mathcal{V}_X$ and $\mathcal{V}_Y$ are the vocabularies of the source language and target language, respectively. Each target token $y_t$ from the translation examples is stored in the datastore with a $d$-dimensional key ($\in \mathbb{R}^d$), which is the representation of the translation context $(\boldsymbol{x}, \boldsymbol{y}_{<t})$ obtained from the decoder of the pre-trained NMT model. The datastore $\mathcal{M} \subseteq \mathbb{R}^d \times \mathcal{V}_Y$ is formally defined as a set of tuples as follows:

$$\mathcal{M} = \left\{ \left( f\left(\boldsymbol{x}, \boldsymbol{y}_{<t}\right), y_t \right) \mid (\boldsymbol{x}, \boldsymbol{y}) \in \mathcal{D}, 1 \le t \le |\boldsymbol{y}| \right\} \tag{1}$$

The size of the datastore for each translation direction is shown in Table 2. During decoding, $k$NN-MT retrieves the $k$-nearest-neighbor key–value pairs $\{(\boldsymbol{k}_i, v_i)\}_{i=1}^k \subseteq \mathbb{R}^d \times \mathcal{V}_Y$ from the datastore $\mathcal{M}$ using the query vector $f\left(\boldsymbol{x}, \boldsymbol{y}_{<t}\right)$ at timestep $t$. $f : \mathcal{V}_X^{|x|} \times \mathcal{V}_Y^{t-1} \to \mathbb{R}^d$ represents the intermediate representation of the final decoder layer from the source sentence and prefix target tokens. In our system, the value of $k$ is set to 32 for all translation directions. In order to speed up the retrieval during translation, we used FAISS (Johnson et al., 2019). We then obtained the output probability for each token by interpolating the $k$NN-MT probability and the probability from the translation model. The formula for calculating the $k$NN-MT probability is:

$$p_{k\text{NN}}\left(y_t \mid \boldsymbol{x}, \boldsymbol{y}_{<t}\right) \\ \propto \sum_{i=1}^k \mathbb{1}_{y_t = v_i} \exp \frac{-\left\| \boldsymbol{k}_i - f\left(\boldsymbol{x}, \boldsymbol{y}_{<t}\right) \right\|_2^2}{\tau} \tag{2}$$

The formula for calculating the output probability is as follows:

$$P\left(y_t \mid \boldsymbol{x}, \boldsymbol{y}_{<t}\right) \\ = \lambda p_{k\text{NN}}\left(y_t \mid \boldsymbol{x}, \boldsymbol{y}_{<t}\right) + (1-\lambda) p_{\text{NMT}}\left(y_t \mid \boldsymbol{x}, \boldsymbol{y}_{<t}\right). \tag{3}$$

For all translation directions, we set $\lambda = 0.3$ and $\tau = 100$ in the $k$NN-MT decoding.

| datastore | size |
|---|---|
| en→as | 1,212,711 |
| en→kha | 1,024,451 |
| en→mni | 574,142 |
| en→mz | 1,404,832 |
| as→en | 1,253,490 |
| kha→en | 878,620 |
| mni→en | 524,002 |
| mz→en | 1,263,000 |

Table 2: Datastore size for all translation directions.

## 4 Experiments

All of our translation models were implemented based on fairseq (Ott et al., 2019) and trained on 8 NVIDIA 3090 GPUs. All models use the same structure of 12 transformer layers (Vaswani et al., 2017). During training, we used the Adam (Kingma, 2014) optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.98$, the learning rate scheduling strategy of inverse sqrt, the number of warmup step set to 4000, the maximum learning rate set to 0.0005 and FP16 to accelerate the training process. We trained our models till convergence with early stopping criteria with a patience of 5. The dropout ratio is set to 0.5. We used a fixed beam size of 4 and a length penalty of 0.8 when doing back-translation.

All experiments were evaluated using the sacrebleu (Post, 2018) tool to calculate BLEU (Papineni et al., 2002) scores on the official validation sets.

## 5 Results

As shown in Table 3, each method can bring certain improvements to the model. However, pre-training and back-translation did not bring much improvement. For example, pre-training leads to an improvement of 0.82 BLEU on average, while back-translation brings BLEU improvements of 0.41. In particular, back-translation has caused some damage to the performance of the model on some translation directions. The BLEU in en→mni direction dropped from 25.17 to 24.04. This may be caused by the low quality of pseudo-parallel data. We believe that fine-tuning the model separately using the data of each language pair is necessary for a multilingual translation model. And it achieves 1.03 BLEU improvement on average. Doing so can alleviate the problem of inconsistent convergence of different language pairs in joint training, although it does not benefit all translation directions. It can be seen that all translation directions

| System | en→as | en→kha | en→mz | en→mni | as→en | kha→en | mz→en | mni→en |
|---|---|---|---|---|---|---|---|---|
| M2M Baseline | 8.75 | 17.84 | 22.26 | 24.49 | 15.69 | 13.15 | 22.45 | 32.41 |
| Pre-training | 9.24 | 17.77 | 22.72 | 25.17 | 17.70 | **14.05** | 22.89 | 34.08 |
| Back-translation | 11.51 | 18.24 | 23.29 | 24.04 | 17.98 | 13.22 | 23.36 | 35.25 |
| Fine-tuning | 12.50 | 18.29 | 24.17 | 26.93 | 18.55 | 13.33 | 24.32 | 37.06 |
| $k$NN-MT | **12.82** | **18.78** | **29.39** | **28.99** | **19.69** | 13.82 | **31.27** | **39.02** |

Table 3: BLEU scores of all translation direction on validation sets

are further improved with $k$NN-MT (+2.33 BLEU). The four translation directions of the two language pairs en↔mni and en↔mz can even get an average improvement of 4.05 BLEU. This shows the great potential of $k$NN-MT in improving data utilization efficiency, inspiring more research on $k$NN-MT in low-resource scenarios. Finally, from the overall perspective, some translation directions do not benefit much from our methods. The translation performance of the model in these translation directions may be most limited by the size of the data. However, the results in most translation directions still achieve significant improvements over the baseline, which demonstrates the effectiveness of our approach for low-resource machine translation.

## 6 Conclusion

In this paper, we describe DLUT-NLP's submission to the WMT24 low-resource Indic language translation shared task. We participated in four subtasks with a total of eight translation directions. We leveraged methods ranging from pre-training, back-translation, language-specific fine-tuning and $k$NN-MT. Experimental results show that we achieved large improvements in all directions.

## Limitations

We found that our system still has the following limitations:

- We did not perform effective filtering on the pseudo-parallel corpus, and we did not perform iterative back-translation. This may be the reason why our back-translation did not achieve the expected results.

- We believe that we have not made enough use of monolingual data. Next, we need to explore other ways to use monolingual data, such as using other pre-training tasks.

- We did not leverage any existing LLMs because we were not sure whether they were

trained on languages other than English included in the task. This will also be a future exploration mission.

## References

Kenneth Heafield. 2011. Kenlm: Faster and smaller language model queries. In *Proceedings of the sixth workshop on statistical machine translation*, pages 187–197.

Yichong Huang, Xiaocheng Feng, Xinwei Geng, and Bing Qin. 2022. Unifying the convergences in multilingual neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6822–6835, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547.

Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Nearest neighbor machine translation. *arXiv preprint arXiv:2010.00710*.

DP Kingma. 2014. Adam: a method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Guillaume Lample and Alexis Conneau. 2019. Cross-lingual language model pretraining. *arXiv preprint arXiv:1901.07291*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Partha Pakray, Santanu Pal, Advaitha Vetagiri, Reddi Mohana Krishna, Arnab Kumar Maji, Sandeep Kumar Dash, Lenin Laitonjam, Lyngdoh Sarah, and Riyanka Manna. 2024. Findings of wmt 2024 shared task on low-resource indic languages translation. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the wmt 2023 shared task on low-resource indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016b. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Minghao Wu, Yitong Li, Meng Zhang, Liangyou Li, Gholamreza Haffari, and Qun Liu. 2021. Uncertainty-aware balancing for multilingual and multi-domain neural machine translation training. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7291–7305, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

# SRIB-NMT's Submission to the Indic MT Shared Task in WMT 2024

**Pranamya Patil , Raghavendra HR , Aditya Raghuwanshi** and **Kushal Verma**
Samsung Research India Bangalore, Bangalore, India
{pran.patil,raghav.hr,aditya.r1,kushal.verma}@samsung.com

## Abstract

In the context of the Indic Low Resource Machine Translation (MT) challenge at WMT-24 ((Pakray et al., 2024)), we participated in four language pairs: English-Assamese (en-as), English-Mizo (en-mz), English-Khasi (en-kh), and English-Manipuri (en-mn). To address these tasks, we employed a transformer-based sequence-to-sequence architecture (Vaswani et al., 2017). In the PRIMARY system, which did not utilize external data, we first pretrained language models (low resource languages) using available monolingual data before finetuning them on small parallel datasets for translation. For the CONTRASTIVE submission approach, we utilized pretrained translation models like Indic Trans2 (Gala et al., 2023) and applied LoRA Fine-tuning (Hu et al., 2021) to adapt them to smaller, low-resource languages, aiming to leverage cross-lingual language transfer capabilities (CONNEAU and Lample, 2019). These approaches resulted in significant improvements in SacreBLEU scores(Post, 2018) for low-resource languages.

## 1 Introduction

With increasing digital connectivity there is huge demand for good translations systems for people to access wide array of digital information in their native local languages. This gives people flexibility and ease of access. For any machine learning task, the quality and quantity of data is of paramount importance. There are multiple languages which have either negligible or zero digital footprint. On top of that presence of good quality parallel/ bitext data is even more rare event.

Due to increased demand and intangible benefits from translation systems there has been lot of research in the field of machine translation. One such area is the low resource machine translation system. In the above statement "resource" refers to data resource. We are working on handling translations for languages which have very less data available

(both monolingual and parallel).

Some major works for multilingual Machine translation (approx 200*200 languages) is NLLB (Costa-jussà et al., 2022). Here the authors use Mixture of Experts (MoE) to train single large multilingual model capable of handling approx 200+ languages as source(src) and target(tgt) languages. One of the objectives behind NLLB is handling low resource languages.

The details of WMT23 Indic MT findings can be found here (Pal et al., 2023). For our PRIMARY approach (no additional data apart from what was shared), we pretrained language model (using monolingual data). We explored 2 pretraining objectives namely Causal Language Modeling (Radford et al., 2019) and denoising (Lewis et al., 2020). Using these Pretrained Language Models as initial model weights we finetuned for tranlstaion task using available parallel corpus.

For the CONTRASTIVE submission approach, we utilized pretrained translation models like Indic Trans2 (Gala et al., 2023) and applied LoRA Fine-tuning (Hu et al., 2021) to adapt them to smaller, low-resource languages, aiming to leverage cross-lingual language transfer capabilities (CONNEAU and Lample, 2019).

## 2 Related Work

Our submissions use the concepts like transfer learning, denoising pretraining (Lewis et al., 2020), Causal Language Modeling (Radford et al., 2019). We use denoising pretraining and causal language modeling as pretraining tasks. Using transfer learning we use the pretrained model for new task like translation. We use pretrained indic translation model like Indic Trans2 (Gala et al., 2023) for contrastive submission and adapted them to low resource languages using LoRA Fine-tuning (Hu et al., 2021). We aim to take advantage of cross-lingual language transfer capabilities (CONNEAU and Lample, 2019) and hence used IndicTrans2

| Monolingual Data | # lines |
|---|---|
| Assamese | 2,624,715 |
| Manipuri | 2,144,897 |
| Mizo | 1,909,823 |
| Khasi | 182,737 |
| **Parallel Train Data** | **# lines** |
| English <-> Assamese | 50,000 |
| English <-> Manipuri | 21,687 |
| English <-> Mizo | 50,000 |
| English <-> Khasi | 24,000 |

Table 1: Dataset Sizes shared as part of Indic MT Task

machine translation model.

# 3 System Description

We have submitted in 2 categories i) PRIMARY ii) CONTRASTIVE

## 3.1 PRIMARY System

For PRIMARY System we trained the model from scratch using only the shared data as part of workshop. We trained separate model for each direction.

### 3.1.1 Tokenizer

First we need to train tokenizer for each english <-> language pair (ie:- one of Assamese, Manipuri, Mizo, Khasi). We use sentence piece[1] tokenizer library to train joint dictionary (combined vocab) for each of english <-> (Assamese, Manipuri, Mizo, Khasi) language pairs. We use mono data(each from as,mn,kh,mz languages) + almost equal amount of english data to train sentence piece tokenizer (subset choosen from AI4Bharat Samanantaar dataset (Ramesh et al., 2022)).

### 3.1.2 Pretraining

We Pretrain the model using monolingual data. We explored 2 subtasks for the same

1. Causal Language Modeling (Radford et al., 2019) - Here we train decoder only model for causal language modeling ie:- predicting the next word using context words. The pretrained decoder from this task is utilized to initialize the decoder component of our seq2seq architecture employed for the translation task.

2. Denoising Task (Lewis et al., 2020) - We integrate a seq2seq transformer model that takes a noisy version of the input (such as perturbed

---

[1] https://github.com/google/sentencepiece

mono data, added tokens, or shuffled data) and expects the output from the decoder to be the original, unperturbed input. By utilizing this denoising objective task, we aim for the model to understand language patterns and structures. The pretrained model from this task can be used for translation task.

Both of the above pretraining tasks are explored independently and we used each tasks pretrained checkpoints for finetuning separately.

### 3.1.3 Finetuning

Using the pretrained checkpoint we finetune the models for translation task (with small amount of parallel data). Pretraining helps the model to understand the language nuances and leads to faster converging of models for translation tasks.

## 3.2 CONTRASTIVE System

For CONTRASTIVE Submission (where external data etc is allowed). We use translation model of other languages eg:- IndicTrans2 (Gala et al., 2023). As the 4 low resource languages ie:- Mizo, Manipuri, Khasi and Assamese are near to Indic Languages supported by (Gala et al., 2023) we believe the the model will benefit from shared parameters, vocabs and hence map Cross Lingual language references (CONNEAU and Lample, 2019). We use the same tokenizer as used by IndicTrans2 Model.

LoRA (Hu et al., 2021) adaptation is a lightweight and resource-friendly technique for customizing pretrained models. It involves adding small adapter weights (to certain layers) alongside the existing model weights. During training, the original model weights remain unchanged while only the adapter weights are updated. At test time, the adapter weights and the original model weights are combined to generate predictions. This approach allows for efficient customization without requiring extensive modifications to the original model. Since only a small number of parameters are updated during training, the overall training time is reduced. Additionally, this approach helps mitigate the issue of catastrophic forgetting to some degree.

| Parameter | Value |
|---|---|
| `encoder_layers` | 4 |
| `decoder_layers` | 4 |
| `attention_heads` | 8 |
| `embedding_dimension` | 512 |
| `ffn_embedding_dimension` | 4096 |

Table 2: PRIMARY Submission Model Architecture details

| Parameter | Value |
|---|---|
| `lora_rank` | 32 |
| `lora_alpha` | 32 |
| `lora_dropout` | 0.1 |
| `device_batch_size` | 16 |
| `device_grad_accumulation_steps` | 2 |
| `max_steps` | 100,000 |
| `eval_steps` | 5,000 |
| `patience` | 10 |

Table 3: CONTRASTIVE Submission Model details

# 4 Experiments

## 4.1 Implementation

### 4.1.1 PRIMARY Submission

For Primary submission we use fairseq[2] framework for both pretraining and finetuning stage. The model architecture details can be found in Table 2. We experimented with lesser #encoder, #decoder layers as compared to standard (6 encoder and 6 decoder layers) to reduce model complexity and hence training time.

### 4.1.2 CONTRASTIVE Submission

For CONTRASTIVE submission we use Indic-Trans2(Gala et al., 2023) and use huggingface peft library for LoRA finetuning. The model details can be found in Table 3

# 5 Results

The results as shared by conference committee are attached below. PRIMARY submission results Table 4 and CONTRASTIVE submission results in Table 5.

We use SacreBLEU (Post, 2018) for validation evaluation. The validation scores (development set) reported during training for Primary system are attached in Table 6

| Direction | BLEU | TER | RIBES | METEOR | ChrF |
|---|---|---|---|---|---|
| en -> as | 1.32 | 101.83 | 7.1 | 7.44 | 22.15 |
| en -> mn | 0 | 101.83 | 1.91 | 3.07 | 18.89 |
| en -> mz | 0 | 102.98 | 3.61 | 6.20 | 16.46 |
| en -> kh | 0.54 | 103.72 | 8.21 | 9.69 | 17.78 |

Table 4: PRIMARY Submission Scores on test suite

| Direction | BLEU | TER | RIBES | METEOR | ChrF |
|---|---|---|---|---|---|
| as -> en | 29.59 | 34.92 | 35.05 | 74.09 | 64.88 |
| mn -> en | 18.89 | 53.05 | 29.17 | 59.43 | 57.1 |
| mz -> en | 11.27 | 64.94 | 20.26 | 47.84 | 44.82 |
| kh -> en | 4.2 | 80.89 | 12.05 | 32.83 | 31.8 |

Table 5: CONTRASTIVE Submission Scores on test suite

## 5.1 Learnings

Following are the learnings from our Experiments

1. Transfer Learning benefits translation task. We saw it in PRIMARY submissions in which language models are pretrained on denoising/ causal language modeling(CLM) task and then transferred for translation task. Especially its evident from initial bleu score and loss. We saw denoising task led to faster converging of models (lower initial loss) relative to CLM task objective.

2. Languages that share a common linguistic ancestor or follow similar word order patterns (such as SVO or SOV) can benefit from using the same vocabulary and sharing parameters during initialization. This allows for more efficient training and better performance across related languages.

3. Using Translation for related language benefits from cross lingual language reference.(Bleu scores of CONTRASTIVE submission)

4. LoRA finetuning is effective for adapting a translation model to new low resource language (lesser training time and resources).

| Direction | SacreBLEU |
|---|---|
| en -> as | 8 |
| en -> mn | 16.9 |
| en -> mz | 22.3 |
| en -> kh | 11.1 |

Table 6: PRIMARY Submission Scores on development suite shared along with training data

## 5.2 Conclusion

The adpatation of another language translation model to similar but low resource language is benefitted by sharing params, vocabs etc across languages (due to cross lingual language learning). LoRA finetuning leds to quicker converging for low resource languages (18-19 hours on A100 GPU with 40GB of RAM).

We have described our submission to WMT2024 Indic Translation Task, leveraging various concepts like Denoising task(Lewis et al., 2020), Cross Lingual Transfer Learning(CONNEAU and Lample, 2019), IndicTrans2 Model(Gala et al., 2023), LoRA adaptation(Hu et al., 2021) etc.

## Limitations

1. Exploring impact of Iterative Backtranslation(Hoang et al., 2018) benefits using intermediate models in PRIMARY setting.

2. Exploring more pretraining task objectives for PRIMARY System.

3. Exploring multi task learning impact for PRIMARY systems.

4. Exploring the difference in scores, training resources for full precision finetuning vs LoRA finetuning.

## References

Alexis CONNEAU and Guillaume Lample. 2019. Cross-lingual language model pretraining. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Jay Gala, Pranjal A Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, et al. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307*.

Vu Cong Duy Hoang, Philipp Koehn, Gholamreza Haffari, and Trevor Cohn. 2018. Iterative back-translation for neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 18–24, Melbourne, Australia. Association for Computational Linguistics.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Partha Pakray, Santanu Pal, Advaitha Vetagiri, Reddi Mohana Krishna, Arnab Kumar Maji, Sandeep Kumar Dash, Lenin Laitonjam, Lyngdoh Sarah, and Riyanka Manna. 2024. Findings of wmt 2024 shared task on low-resource indic languages translation. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the WMT 2023 shared task on low-resource Indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694, Singapore. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Gowtham Ramesh, Sumanth Doddapaneni, Aravinth Bheemaraj, Mayank Jobanputra, Raghavan AK, Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Mahalakshmi J, Divyanshu Kakwani, Navneet Kumar, Aswin Pradeep, Srihari Nagaraj, Kumar Deepak, Vivek Raghavan, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh Shantadevi Khapra. 2022. Samanantar: The largest publicly available parallel corpora collection for 11 Indic languages. *Transactions of the Association for Computational Linguistics*, 10:145–162.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

# MTNLP-IIITH: Machine Translation for Low-Resource Indic Languages

**Abhinav P M**[*]**, Ketaki Shetye**[*]**, Parameswari Krishnamurthy**
Machine Translation and NLP Lab, LTRC
International Institute of Information Technology, Hyderabad, India
pmabhinav20@gmail.com, ketaki.shetye@research.iiit.ac.in, param.krishna@iiit.ac.in

## Abstract

Machine Translation for low-resource languages poses significant challenges, primarily due to the limited availability of data.The WMT24 Low-Resource Indic Neural Machine Translation task challenges us to employ innovative techniques to improve machine translation for low-resource Indian languages. Our proposed solution leverages advancements in neural machine translation, focusing on methodologies such as back-translation and fine-tuning. By fine-tuning pretrained models like mBART, we achieved significant progress in translating languages such as Manipuri and Khasi. The best score was achieved for the English-to-Khasi (en-kh) primary model, with the highest BLEU score of 0.0492, chrF score of 0.3316, and METEOR score of 0.2589 (on scale of 0 to 1) and comparable scores for other language pairs.

## 1 Introduction

Machine translation is a sub-field of computational linguistics that focuses on developing systems capable of automatically translating text or speech from one language to another. The WMT24 task enables us to perform machine translation on those languages which are considered low-resource that is with limited data availability due to their lesser prevalence or documentation. Our work focuses on translating the 'En-X' language pair in both directions, where 'En' stands for English and 'X' includes Manipuri, a Tibeto-Burman language, and Khasi, which belongs to the Austroasiatic language family.

In recent years, Neural Machine Translation (NMT) has emerged as a powerful approach within machine translation, leveraging deep learning to achieve state-of-the-art results. Although, the NMT models being the data-hungry lead to peformance

degradation when it comes to low resource languages. To tackle this problem, we employed fine-tuning and utilised the mBART (mbart-large-50-many-to-many-mmt) (Tang et al., 2020) experimenting with different configurations and settings for both preprocessing and training. mBART (Liu et al., 2020) is a multilingual sequence-to-sequence model trained on extensive monolingual datasets using a denoising autoencoder approach. It builds on the BART framework (Lewis et al., 2019) by combining a bidirectional encoder with a left-to-right autoregressive decoder, making it suitable for various translation tasks across multiple languages. Even if most of our final systems did not reach a satisfactory or competitive performance, we argue that our experiments brought up some interesting points that deserve more attention.

## 2 Related Work

In a comprehensive study, (Gaikwad et al., 2023) examined fine-tuning-based techniques to improve translation capabilities for low-resource languages by harnessing the multilingual IndicTrans2 model and achieved significant results.

In 2023, (Suman et al., 2023) utilized IndicBART (Dabre et al., 2022) and mBART-large-50, adapting them to specific language pairs and this method led to substantial performance gains for the Assamese and Manipuri languages.

Another 2023 study, (Jha et al., 2023), proposed and evaluated a multilingual neural machine translation system for Indian languages using the mT5 transformer. This system, trained on the modified Asian Language Treebank (ALT) dataset, demonstrated strong performance in translations between English, Hindi, and Bengali, achieving BLEU scores above 20 for five out of the six language pairs.

(Saini and Vidhyarthi, 2023) evaluated various pretrained models for English-to-Marathi translation, developing a bidirectional system. The

---

[0]* These authors contributed equally to this work.

findings indicated that fine-tuning significantly enhanced the mBART model's performance.

(Signoroni and Rychly, 2023) addressed the challenges of neural machine translation (NMT) for low-resource language pairs by using supervised NMT systems. They experimented with different configurations and settings for both preprocessing and training, delving into the complexities of translating these languages.

## 3 Dataset

### 3.1 Languages

**Manipuri**, also known as Meitei or Meiteilon, is predominantly spoken in the northeastern Indian state of Manipur and is one of India's 22 scheduled languages having about 1.8 million native speakers (Signoroni and Rychly, 2023). It is distinguished by its rich morphological features, including a complex phonological system with tones, an agglutinative structure, and a Subject-Object-Verb (SOV) word order, as shown in Figure 1. As a tonal language, Manipuri uses various tones and pronunciations to convey meaning and employs primarily two scripts: Meitei and Bengali. Despite being a scheduled language, Manipuri is often considered a low-resource language in natural language processing, presenting valuable opportunities for transfer learning and the development of multilingual models.

**Khasi** primarily spoken in the northeastern Indian state of Meghalaya by the Khasi people, is one of the major languages of the region spoken by over 1 million individuals (Signoroni and Rychly, 2023). Belonging to the Austroasiatic language family, Khasi is more commonly written using the Latin alphabet. Structurally, Khasi typically follows a subject-verb-object (SVO) order, similar to English, but differs from most Indian languages, which generally use a subject-object-verb (SOV) order, as shown in Figure 1.

| Language | Sentence |
|---|---|
| Manipuri (Subject-Object-Verb) | নুংঙবা চাক চবা। |
| Khasi (Subject-Verb-Object) | U khynnah u bam ïa ka soh. |

Figure 1: Translations of "The boy eats an apple" showing word order in Manipuri and Khasi.

### 3.2 Composition

In this study, we used WMT 2024 (Pal et al., 2023) (Kakum et al., 2023) to fine-tune which includes both bilingual and monolingual data. For the bilingual data, we used the language pairs English-Khasi (en ↔ kh) and English-Manipuri (en ↔ mn). The compositions of these datasets are presented in Table 1 and 2.

| Lang. Pair | Train | Test | Validation | Monolingual |
|---|---|---|---|---|
| en ↔ kh | 24,000 | 1000 | 1000 | 182,737 |
| en ↔ mn | 21,687 | 1000 | 1000 | 2,144,897 |

Table 1: Number of lines in the dataset for the language pairs used in the task. The Monolingual column refers to the size of the non-English side.

| Lang. Pair | Type:Token Ratio | Avg. Sentence Length |
|---|---|---|
| en ↔ kh | 0.019 (en) | 30.41 tokens (en) |
|  | 0.0093 (kh) | 36.48 tokens (kh) |
| en ↔ mn | 0.0573 (en) | 18.02 tokens (en) |
|  | 0.0083 (mn) | 15.23 tokens(mn) |

Table 2: Training Dataset Statistics for Language Pairs: Type-Token Ratio and Average Sentence Length

## 4 System Overview

### 4.1 Initial Fine-Tuning

We begin by fine-tuning the mBART model (mbart-large-50-many-to-many-mmt) for the language pairs: English to Khasi (en → kh), Khasi to English (kh → en), English to Manipuri (en → mn), and Manipuri to English (mn → en). To ensure the quality and consistency of the bilingual data, we perform several preprocessing steps, including the removal of HTML tags, invisible characters, newline tabs, and duplicate entries.

For machine translation preparation, we tokenize both the input and the target texts. Truncating techniques are applied to standardize the texts by setting the maximum length of the tokenized sequences to 512 tokens. This ensures uniformity across all the examples in the dataset. This serves as our baseline model.

### 4.2 Data Augmentation

#### 4.2.1 Backtranslation

A back-translation strategy (Sennrich et al., 2016) is employed to augment the training dataset with more data. Specifically, we back-translate 100,000 monolingual Khasi and Manipuri sentences into English using the baseline model. However, it is likely that the backtranslated data contains a significant portion of low-quality translations. To remove

| Lang. Pair | Filtered Data |
|---|---|
| kh ↔ en | 534 |
| mn ↔ en | 662 |

Table 3: Count of high-quality sentence pairs after cosine similarity filtering (threshold 0.84) for Khasi-English (kh ↔ en) and Manipuri-English (mn ↔ en).

low-quality data and ensure high-quality translation pairs, we employ a filtering process using the LaBSE model and cosine similarity.

### 4.2.2 Data Filtering

**LaBSE Fine-tuning** The Language-agnostic BERT Sentence Embedding (LaBSE) model (Feng et al., 2022) is not originally trained in the Khasi or Manipuri languages. Therefore, to generate accurate sentence embeddings for these language pairs, we fine-tune the LaBSE model specifically for en ↔ kh and en ↔ mn pairs, despite the limited size of available bilingual data. This fine-tuned model is then employed to produce sentence embeddings for the back-translated Khasi-English and Manipuri-English pairs.

**Embedding and Similarity Calculation** To ensure the accuracy of the back-translated data, we use cosine similarity, a metric that measures the cosine of the angle between two vectors in multidimensional space, to compare sentence embeddings. We apply a threshold of 0.84, effectively filtering out low-quality translations, and retaining only those pairs that meet our quality standards. Consequently, only a small portion of the original 100,000 back-translated sentences remain after filtering using this threshold. The data retained after filtering are presented in Table 3.

### 4.2.3 Further Filtering and cleaning

Despite filtering, some sentences with continuous symbols or non-English characters remain. To address this, we conduct an additional data cleaning round, removing sentences with continuous symbols and residual Manipuri or Khasi words in the English translations. The cleaned data is then combined with the original training set to create the augmented dataset as the final dataset.

### 4.3 Training with Augmented Data

Subsequently, we fine-tune the mBART model (mbart-large-50-many-to-many-mmt) using the augmented dataset, which includes both the original training data and the filtered back-translated

data. The same data-preprocessing steps are employed for the augmented dataset as applied for the baseline model to maintain uniformity. The fine-tuning process incorporates this augmented data to enhance the model's performance and robustness.

## 5 Results and Analysis

Table 4 shows WMT24 evaluation results, highlighting a more challenging test set compared to last year. The low average semantic similarity score of 0.0253 as found using the (Reimers and Gurevych, 2019) sentence-transformer model that maps sentences to a 384 dimensional dense vector space to calculate semantic similarity between train and test data indicating reduced model performance too.

Among the primary models, the English-to-Khasi (en → kh) model, trained on both original and filtered back-translated data, performed the best across most metrics. It achieved the highest BLEU score of 0.0492, a chrF score of 0.3316, and a METEOR score of 0.2589, indicating strong performance for this language pair. The high chrF score shows effective capture of character-level nuances, while the lowest TER score of 84.79 reflects fewer required edits to match reference translations.

It is important to note that the Khasi-to-English (kh-en) primary model is excluded from this evaluation because of issues encountered during the evaluation process. Meanwhile, the English-to-Manipuri (en → mn) model shows a BLEU score of 0, highlighting the difficulty in translating from English to Manipuri. This may be partly due to the smaller size of the training data for this pair compared to the en-kh pair. Despite some fluctuations, the overall performance of both models is comparable in the test data.

When considering all metrics, the primary model shows slight improvements over the baseline model, indicating that the additional filtered back-translated data enhanced translation quality. However, the baseline model also performs competitively. The filtered back-translated data includes only 534 sentences for the kh ↔ en pair and 662 for the mn ↔ en pair (Table 3), with the small dataset likely due to stringent filtering.

Table 5 shows that the primary model slightly outperforms the baseline model in BLEU, chrF, and TER metrics for the English-to-Khasi (en → kh) and English-to-Manipuri (en → mn) models, sug-

| Lang. Pair | BLEU (↑) | chrF (↑) | TER (↓) | METEOR (↑) | RIBES (↑) |
|---|---|---|---|---|---|
| **Baseline** | | | | | |
| en-kh | 0.0359 | 0.2333 | 103.49 | 0.1649 | 0.1106 |
| kh-en | 0.0060 | 0.1731 | 106.60 | 0.1020 | 0.0487 |
| en-mn | 0.0064 | 0.3191 | 96.46 | 0.0724 | 0.0628 |
| mn-en | 0.0484 | 0.2662 | 101.76 | 0.1940 | 0.1087 |
| **Primary** | | | | | |
| en-kh | 0.0492 | 0.3316 | 84.79 | 0.2589 | 0.1595 |
| en-mn | 0.0000 | 0.3325 | 94.77 | 0.0822 | 0.0737 |
| mn-en | 0.0362 | 0.2777 | 94.79 | 0.1873 | 0.1136 |

Table 4: Results of Primary and Baseline models evaluated on WMT24 evaluation test data (scores calculated on the scale from 0 to 1)

| Lang. Pair | BLEU (↑) | chrF (↑) | TER (↓) |
|---|---|---|---|
| **Baseline** | | | |
| en-kh | 0.1748 | 0.3964 | 0.75699 |
| kh-en | 0.1274 | 0.3566 | 0.8791 |
| en-mn | 0.2089 | 0.5676 | 0.6537 |
| mn-en | 0.3265 | 0.5709 | 0.6522 |
| **Primary** | | | |
| en-kh | 0.1867 | 0.4126 | 0.7275 |
| kh-en | 0.1234 | 0.3570 | 0.8845 |
| en-mn | 0.2097 | 0.5726 | 0.6495 |
| mn-en | 0.3224 | 0.5698 | 0.6483 |

Table 5: Results of Baseline and Primary models evaluated on WMT23 validation data (scores calculated on the scale from 0 to 1)

gesting that filtered back-translated data improves translation quality. However, for the Khasi-to-English (kh → en) and Manipuri-to-English (mn → en) models, the primary model experiences a slight performance drop. This decrease is likely due to the filtered back-translated English sentences lacking coherence and contextual appropriateness, which affects the model's effectiveness. Despite these variations, the primary model still performs slightly better than the baseline model when considering all metrics.

## 6 Conclusion

Improving machine translation for low-resource languages remains a critical focus in the field. In this paper, we develop a system for translating low-resource Indic languages, specifically Manipuri and Khasi, in both English-to-Indic and Indic-to-English language pairs. We use back-translation and then apply cosine similarity for data filtering. While effective, their success depends on the quality of the back-translation and fine-tuned LaBSE models.

The morphological complexity of the Indic languages along with inability of capturing cultural and context specific meanings also poses as a challenge which the model could not solve. We further encounter challenges including data scarcity and

high computational requirements which we believe can help produce better results.

## 7 Future Work

For future work, we aim to focus on enhancing machine translation for low-resource languages by leveraging language-specific properties such as part-of-speech (POS) tags and dependency parsing. By integrating POS tagging one can enable the model to better understand the syntactic roles of words, leading to more accurate and contextually appropriate translations. Dependency parsing can also capture the grammatical structure and relationships between words, allowing the model to manage complex sentence structures more effectively. Additionally, the use of more filtered backtranslated data can provide a richer training dataset, further improving translation quality. Combining these linguistic techniques with extensive backtranslation, so that we can capture the nuances of the individual low-resource languages, we can significantly address the current challenges in machine translation.

## References

Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh Khapra, and Pratyush Kumar. 2022. IndicBART: A pre-trained model for indic

natural language generation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1849–1863, Dublin, Ireland. Association for Computational Linguistics.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic bert sentence embedding.

Pranav Gaikwad, Meet Doshi, Sourabh Deoghare, and Pushpak Bhattacharyya. 2023. Machine translation advancements for low-resource Indian languages in WMT23: CFILT-IITB's effort for bridging the gap. In *Proceedings of the Eighth Conference on Machine Translation*, pages 950–953, Singapore. Association for Computational Linguistics.

Abhinav Jha, Hemprasad Yashwant Patil, Sumit Kumar Jindal, and Sardar M N Islam. 2023. Multilingual indian language neural machine translation system using mt5 transformer. In *2023 2nd International Conference on Paradigm Shifts in Communications Embedded Systems, Machine Learning and Signal Processing (PCEMS)*, pages 1–5.

N. Kakum, S.R. Laskar, K. Sambyo, and et al. 2023. Neural machine translation for limited resources english-nyishi pair. *Sādhanā*, 48:237.

Ivana Kvapilíková and Ondřej Bojar. 2023. Low-resource machine translation systems for Indic languages. In *Proceedings of the Eighth Conference on Machine Translation*, pages 954–958, Singapore. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *CoRR*, abs/1910.13461.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the WMT 2023 shared task on low-resource Indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694, Singapore. Association for Computational Linguistics.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Lekhraj Saini and Deepti Vidhyarthi. 2023. Bidirectional english-marathi translation using pretrained models: A comparative study of different pre-trained models. pages 1–8.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Edoardo Signoroni and Pavel Rychly. 2023. MUNI-NLP systems for low-resource Indic machine translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 959–966, Singapore. Association for Computational Linguistics.

Dhairya Suman, Atanu Mandal2, Santanu Pal3, and Sudip Kumar Naskar. 2023. Iacs-lrilt: Machine translation for low-resource indic languages.

Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. Multilingual translation with extensible multilingual pretraining and finetuning.

# Exploration of the CycleGN Framework for Low-Resource Languages

**Sören Dréano**
ML-Labs
Dublin City University
soren.dreano2@mail.dcu.ie

**Derek Molloy**
School of Electronic Engineering
Dublin City University
derek.molloy@dcu.ie

**Noel Murphy**
School of Electronic Engineering
Dublin City University
noel.murphy@dcu.ie

## Abstract

CycleGN is a Neural Machine Translation framework relying on the Transformer architecture. Its approach is similar to a Discriminatorless CycleGAN, specifically tailored for non-parallel text datasets.

The foundational concept of our research posits that in an ideal scenario, retro-translations of generated translations should revert to the original source sentences. Consequently, a pair of models can be trained using a Cycle Consistency Loss only, with one model translating in one direction and the second model in the opposite direction.

One of the main advantages of such an approach is that it makes it possible to learn with non-parallel datasets, which are by definition rare and short for low-resource languages. In order to verify this hypothesis and as a contribution to the WMT24 challenge, CycleGN models were trained for both the "Translation into Low-Resource Languages of Spain" and "Low-Resource Indic Language Translation" tasks. These submissions fall under the "constrained" category, as no pre-trained translation model was used, and the models were trained using the provided datasets.

Given that the CycleGN architecture demonstrated its capacity to learn from non-parallel datasets, the authors anticipated that it would similarly be effective in learning from low-resource languages. However, preliminary results indicate that, for most low-resource language pairs, the models did not exhibit significant learning ability. This study explores this lack of learning.

## 1 Introduction

The introduction of the Transformer architecture (Vaswani et al., 2017) marked a significant advancement in Machine Translation, rapidly gaining widespread adoption. Its parallelized structure enhanced computational efficiency, allowing for the integration of a larger number of parameters.

Neural Machine Translation (NMT) relies on extensive text corpora, structured as aligned pairs, where sentences of equivalent meaning are available in at least two different languages. This alignment is crucial for initiating model training to establish linguistic connections. Ongoing efforts, such as OPUS (Tiedemann and Thottingal, 2020) and Tatoeba (Tiedemann, 2012), focus on providing public access to these datasets. However, parallel datasets represent only a small fraction of the total data available in monolingual datasets.

While large parallel corpora exist for many language pairs, the ability to utilize monolingual datasets alone would greatly increase the available training data. This approach is particularly advantageous for low-resource languages, with limited parallel text corpora.

Back-translation (Sennrich et al., 2016) is a technique that enhances training data by using a pretrained machine translation (MT) model to translate sentences from a monolingual dataset, creating synthetic parallel pairs. This method allows for the generation of additional training examples in situations where parallel corpora are scarce.

This research builds on the concept that translating a sentence from a source language to a target language, and then back-translating it to the source language, provides a means to evaluate the effectiveness of the translation models. By comparing the original sentence with the machine-generated back-translation, the discrepancy is then quantified using a Cycle Consistency Loss, which serves as a metric for model performance and guides the backpropagation of gradients within the neural networks. This approach is analogous to techniques used in Image-to-Image Translation, such as the CycleGAN framework proposed by Zhu et al. (2017).

## 2 Previous work

The TextCycleGAN model (Lorandi et al., 2023), although not based on the Transformer architecture

or focused on Machine Translation (MT), introduced a novel approach for text style transfer. This method applied a CycleGAN to the Yelp dataset, enabling the model to learn mappings between positive and negative textual styles without the need for paired examples.

Shen et al. (2017) demonstrated the potential of training two encoder-decoder networks in an unsupervised manner, allowing for the sharing of a latent space and facilitating style transfer. Similarly, Lample et al. (2018) extended this technique to the MT domain, proving that effective translation can be achieved without relying on parallel datasets.

## 3 Definitions

Machine Translation models are most commonly trained using "parallel" datasets, which are structured collections of text pairs. Each pair comprises a segment of text in a source language and its translation in the target language. A non-parallel dataset on the other hand does not consist in pairs of text segments, consequently the source and target sentences do not share any explicit correspondence.

In the context of this study, the datasets are "permuted". A permuted dataset is defined as a parallel dataset wherein the sentences of one language have been systematically rearranged. Consequently, this results in a non-parallel corpus where it is guaranteed that each sentence has a corresponding translation located at an unspecified index within the dataset.

## 4 Datasets

The PILAR dataset (Galiano-Jiménez et al., 2024) has been used exclusively for the low-resource languages of Spain. Using a parallel curated dataset as a starting point ensures that the dataset is non-parallel by permuting the sentences. For each Iberian language, both a literary and a crawled versions were available in the PILAR datasets and have been merged for training. The development sets of the PILAR dataset are translations of the development sets of the FLORES dataset (NLLB Team et al., 2022), which is an evaluation benchmark for multilingual machine translation.

The Low-Resource Indic Language Translation task was also part of the WMT23 (Pal et al., 2023). The datasets were kept the same between the two editions.

Table 1 references the number of sentences used for each language-pair.

| Language Pair | Number of lines | Number of epochs |
|---|---|---|
| Spanish-Aragonese | 84,703 | 10 |
| Spanish-Asturian | 38,869 | 10 |
| English-Assamese | 2,624,715 | 1 |
| English-Khasi | 182,737 | 3 |
| English-Manipuri | 2,144,897 | 1 |
| English-Mizo | 1,909,823 | 1 |

Table 1: Number of sentences for each language pair and number of epochs during training

## 5 Training

For clarity and consistency, the mathematical notations from the original CycleGAN framework will be adopted in this study. The objective is to develop two Neural Machine Translation (NMT) models for two languages, $\mathcal{X}$ and $\mathcal{Y}$, using their respective datasets. Specifically, we aim to construct models $\mathcal{G} : \mathcal{X} \mapsto \mathcal{Y}$ and $\mathcal{F} : \mathcal{Y} \mapsto \mathcal{X}$ such that, in the ideal scenario of perfect translation, the relationships $\mathcal{G}(\mathcal{F}(y)) = y$ and $\mathcal{F}(\mathcal{G}(x)) = x$, with $x \in \mathcal{X}$ and for $y \in \mathcal{Y}$.

To achieve this, the Cross-Entropy Loss (CEL) (Zhang and Sabuncu, 2018) is utilised as the Cycle Consistency Loss (CCL), which measures the distance between the original sentence and its doubly translated counterpart, thereby guiding the computation of gradients.

Furthermore, similar to the original CycleGAN implementation, our study also incorporates an Identity Loss (IL) to enhance training stability. This loss, also based on CEL, ensures that when the model $\mathcal{G}$, which maps $\mathcal{X} \mapsto \mathcal{Y}$, receives an input $y \in \mathcal{Y}$, the output remains unchanged, i.e., $\mathcal{G}(y) = y$. The same loss function is applied to $\mathcal{F}$, ensuring that $\mathcal{F}(x)$ remains equal to $x$, as illustrated in Figure 1.

Further details of the training process, including the specific methodologies, vocabulary organization and pretraining, are comprehensively discussed in the CycleGN submission for the WMT24 main translation task. Readers interested in the full technical details are encouraged to refer to that publication for a more complete understanding of the training framework.

### 5.1 Model architecture

The architecture used for both models, $\mathcal{G}$ and $\mathcal{F}$, is the Marian framework (Junczys-Dowmunt et al., 2018) implemented by Huggingface's Transformers library (Wolf et al., 2020), which is licensed under the Apache Licence. While most parameters

Figure 1: CycleGN training process

follow the default configuration, Table 2 references the changes that were made in order to reduce the computational cost of the architecture.

| Parameter | Huggingface | Current work |
|---|---|---|
| Vocabulary size | 58,101 | 32,000 |
| Encoder layers | 12 | 6 |
| Decoder layers | 12 | 6 |
| Encoder attention heads | 16 | 8 |
| Decoder attention heads | 16 | 8 |
| Encoder feed-forward | 4096 | 2048 |
| Decoder feed-forward | 4096 | 2048 |
| Position embeddings | 1024 | 128 |
| Activation function | GELU | ReLU |

Table 2: Non-default parameters in the configuration of Marian Transformer models

## 6 Results

Even if tracking the CCL is an inexpensive manner to estimate the progress of the training of the CycleGN architecture, a low loss value can also hide an absence of translation. Indeed, as there is no Discriminator to ensure that $\hat{x}$ belongs to $\mathcal{X}$ and $\hat{y}$ belongs to $\mathcal{Y}$, $\mathcal{G}$ and $\mathcal{F}$ will converge towards $x = \hat{y} = \hat{\hat{x}}$ and $y = \hat{x} = \hat{\hat{y}}$, as this approach achieves an optimal outcome on the CCL function, registering a value of zero. This is why an evaluation metric such as COMET is crucial to assess the progression of the CycleGN framework. To measure the performances of CycleGN, every $1,000^{th}$ batch the CCL was averaged.

### 6.1 Indic Languages

Tracking the evolution of the CDC clearly shows the absence of learning in the four language pairs examined. The evolution of the CCL is particularly chaotic, which is partly due to an imbalance of class. Table 3 displays the average number of tokens in the Indic datasets depending on the language. In 3 of the 4 cases, the difference is large, i.e. sentences where the difference in the number of tokens is more than 10%.

| Language pair | Length of source | Length of target |
|---|---|---|
| English-Assamese | 33.10 | 22.81 |
| Encoder layers | 24.09 | 75.27 |
| English-Manipuri | 24.09 | 26.07 |
| English-Mizo | 32.05 | 17.55 |

Table 3: Average number of tokens in sentences

Figures 2, 3, 4 and 5 display the respective evolution of the Cycle Consistency Loss during the training of the language-pairs English-Assamese, English-Khasi, English-Manipuri and English-Mizo.

Contrary to what the authors had hoped for on the basis of previous results obtained for the main task of the WMT24 challenge, no model followed the expected learning curve, i.e. $G$ and $F$ models with a close and slowly decreasing Cycle Consistency Loss.

To reduce this imbalance of class, it may be necessary to manually adjust the size of the sentences.

Figure 2: Evolution of the Cycle Consistency Loss during the training of the English-Assamese model



Figure 3: Evolution of the Cycle Consistency Loss during the training of the English-Khasi model



Figure 4: Evolution of the Cycle Consistency Loss during the training of the English-Manipuri model

This can be done by choosing another tokenization method, selectively choosing phrases to keep only those of a similar size, or by trimming sentences to lengthen or shorten them as required.



Figure 5: Evolution of the Cycle Consistency Loss during the training of the English-Mizo model

## 6.2 Iberian Languages

As with Indic Languages, CycleGN was unable to learn from the datasets provided. However, it was not due to an imbalance of classes in this case, but rather because the classes were too close together, as the Iberian languages are very close to the source language, Spanish.



Figure 6: Evolution of the Cycle Consistency Loss during the training of the Spanish-Aragonese model

Rather than translating directly from Spanish into Aranese or Asturian, it is possible that translation can be achieved by using a different intermediate language such as English. Thus, two CycleGN models would have to be trained, the first to translate from Aranese or Asturian into English, and the second from English into Spanish. This would double the training time for an already expensive framework.

Figure 7: Evolution of the Cycle Consistency Loss during the training of the Spanish-Asturian model

## 7 Conclusion

In conclusion, while the training process demonstrated significant progress and effective translation capabilities in the main study, the results presented in this paper reveal several challenges that prevented similar success. The issues identified, particularly in relation to both class imbalance and class proximity, indicate that further refinement and investigation are necessary. Future research should focus on addressing these challenges, with the aim of optimizing the training process and overcoming the outlined issues. Resolving these problems is crucial for realizing the full potential of the framework within the context discussed in this paper.

## 8 Future Work

Further investigations will benefit from the incorporation of a more extensive dataset and an exploration of larger model architectures. Future work also include methods discussed in Section 6 to allow translation training.

### 8.1 Large dataset

The current work has been trained on a small dataset compared to MT standards. Future work should try to see how convergence progresses with more iterations. Further computational optimizations are probably necessary to shorten the training time required.

### 8.2 Larger models

The current architecture relies on a total of 158,769,152 parameters, which is only about a third of the size of the default in the Huggingface library.

## 9 Source Code

The source code of CycleGN is available at https://github.com/SorenDreano/CycleGN.

## Limitations

The investigation acknowledges certain inherent limitations which may impact the generalizability and applicability of the findings.

### Language diversity

Another issue that arises from the computing cost of CycleGN is the lack in language diversity. Indeed, our current work only used the English-German and Chinese-English language pairs. Consequently, it cannot be certain that the approach presented can be applied to other languages and all alphabets. This is why CycleGN is taking part in WMT24, to explore the framework's performance on a wide range of language pairs.

### Training limitations

Due to time constraints and the fact that CycleGN is a computationally expensive architecture, it was not possible to train the Spanish-Aranese pair. Similarly, the training of all models was stopped early, before reaching performance stagnation.

## Ethics Statement

This study, focusing on the training of NMT models using non-parallel datasets, adheres to the highest ethical standards in research. We recognize the critical importance of ethical considerations in computational linguistics and machine learning, especially as they pertain to data sourcing, model development, and potential impacts on various linguistic communities.

Our research utilizes publicly available, non-parallel linguistic datasets. We ensure that all data is sourced following legal and ethical guidelines, respecting intellectual property rights and privacy concerns.

In our commitment to scientific integrity, we maintain transparency in our research methodologies, model development, and findings. We aim to make our results reproducible and accessible to the scientific community, contributing positively to the field of machine translation.

## Acknowledgements

# References

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Pilar.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018. Unsupervised machine translation using monolingual corpora only.

Michela Lorandi, Maram A.Mohamed, and Kevin McGuinness. 2023. Adapting the CycleGAN Architecture for Text Style Transfer. *Irish Machine Vision and Image Processing Conference*.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the WMT 2023 shared task on low-resource Indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694, Singapore. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data.

Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment.

Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-MT — Building open translation services for the World. In *Proceedings of the 22nd Annual Conferenec of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey. European Language Resources Association (ELRA).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface's transformers: State-of-the-art natural language processing.

Zhilu Zhang and Mert R. Sabuncu. 2018. Generalized cross entropy loss for training deep neural networks with noisy labels.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks.

# The SETU-ADAPT Submissions to the WMT24 Low-Resource Indic Language Translation Task

**Neha Gajakos, Prashanth Nayak**[a]
**Rejwanul Haque**[b], **Andy Way**
ADAPT Centre, Dublin City University, Dublin, Ireland
[a]KantanAI, Dublin, Ireland
[b]South East Technological University, Carlow, Ireland
neha.gajakos@adaptcentre.ie,pnayak@kantanai.io
rejwanul.haque@setu.ie,andy.way@adaptcentre.ie

## Abstract

This paper presents the SETU-ADAPT's submissions to the WMT 2024 Low-Resource Indic Language Translation task. We participated in the unconstrained segment of the task, focusing on the Assamese-to-English and English-to-Assamese language pairs. Our approach involves leveraging Large Language Models (LLMs) as the baseline systems for all our MT tasks. Furthermore, we applied various strategies to improve the baseline systems. In our first approach, we fine-tuned LLMs using all the data provided by the task organisers. Our second approach explores in-context learning with few-shot prompting. In our final approach we explore an efficient data extraction technique based on a fuzzy match-based similarity measure for fine-tuning. We evaluated our systems using BLEU, chrF, WER, and COMET. The experimental results showed that our strategies can effectively improve the quality of translations in low-resource scenarios.

## 1   Introduction

Advances in deep learning have led to major improvements in present-day MT systems. However, developing reasonable-quality MT systems for low-resource languages, especially those from the Indic language family, remains a challenge (Pal et al., 2023). India, home to numerous ancient and morphologically rich languages, presents unique obstacles for MT development due to the intricate morphology, syntax, and scarcity of parallel data for many regional languages (Suman et al., 2023; Ahmed et al., 2023). This motivated us to participate in the WMT 2024 Low-Resource Indic Language Translation task and contribute to the advancements in indic MT systems.

Large-pre-trained models are becoming the norm in MT due to their accuracy, scalability, and usage flexibility. Hence, for our experiments we chose LLMs as our baseline MT systems. More specifically, we used IndicTrans2[1] as the baseline for building all our MT systems. We carried out our experiments for Assamese-to-English and English-to-Assamese language pairs.

We conducted experiments applying different methodologies for improving the performance of our MT systems. Our primary approach involves fine-tuning LLMs using all the available data. However, Assamese is a very low-resource language, and obtaining good quality data is challenging. Since there is limited availability of domain-specific parallel data, in our second approach we generated synthetic data by retrieving a large corpus of monolingual data from OPUS[2]. We then performed similarity search in order to identify domain-specific sentences of target language from the generic data and back-translated them into the source language. Our third approach involves investigating in-context learning using few-shot prompting. We augmented the prompt with samples whose source-side is similar to the source sentence to be translated.

The rest of the paper is organised as follows: we discuss related works in Section 2. We detail the data sets used in Section 3. Our models and experimental setups are described in Sections 4 and 5. The results are reported and findings are discussed in Section 6. Section 7 concludes this work and discusses avenues for future work.

---

[1]https://github.com/AI4Bharat/IndicTrans2?tab=readme-ov-file#indictrans2
[2]https://opus.nlpl.eu/NLLB/as&en/v1/NLLB

## 2 Related Work

In this section, we discuss the papers that are related to our work. Burchell et al. (2022) introduced a framework that differentiates between lexical and syntactic diversity in back translation. Their research highlights that while both types of diversity improve Neural MT (NMT) performance, lexical diversity is more critical. They also demonstrated that nucleus sampling, a method that balances diversity with adequacy, provides superior results for low-resource and mid-resource language pairs.

Ahmed and Buys (2024) introduced the concept of "Synthetic Pivoting" to address the limitations of traditional pivot-based methods, which often face challenges due to structural mismatches between the pivot and low-resource languages. Synthetic Pivoting generates synthetic pivot sentences that better align with the structure of both the source and target languages, resulting in more accurate translations. This method has substantially improved translation quality, particularly for Southern African languages, by simplifying the translation process and effectively utilising high-quality synthetic data.

Suman et al. (2023) focused on improving the translation quality for low-resource Indic languages: Manipuri and Assamese. They leveraged linguistic and scriptural similarities between these languages and Bengali to improve translation outcomes. By utilising pre-trained models on Bengali and incorporating transliteration techniques, they were able to overcome the challenges posed by the limited resources available for Manipuri and Assamese. Their experiments showed that their approaches were effective in improving translation.

Moslem et al. (2023) explored using LLMs for adaptive translation. Their research demonstrated that in-context learning with LLMs enables real-time adaptation to specific terminology and stylistic preferences during inference. They showed that this eliminated the need for extensive fine-tuning. They found that few-shot in-context learning, especially when combined with fuzzy matches from translation memories, can outperform traditional encoder-decoder models regarding translation quality, particularly for high-resource language pairs.

Zhang et al. (2023) investigated the potential of fine-tuning LLMs for MT, focusing on decoder-based models that had not been extensively studied before. They evaluated 15 publicly available LLMs using methodologies such as zero-shot prompting,

few-shot learning, and fine-tuning, with a particular emphasis on the QLoRA (Dettmers et al. (2023)) fine-tuning method. QLoRA proved a highly effective technique, reducing memory usage by quantising the model to 4-bit precision and limiting the number of trainable parameters. Their findings showed that fine-tuning LLMs, especially using QLoRA, significantly outperformed zero-shot and few-shot approaches, particularly in document-level translation tasks.

## 3 Data

We utilised the data provided by WMT organisers for our experiments. The data statistics are detailed in Table 1.

| Assamese ↔ English | |
|---|---|
| Files | Sentences |
| Train | 50,000 |
| Valid | 2,000 |
| Test (2023) | 2,000 |
| Test (2024) - Blind Test | 500 |

Table 1: *Statistics of the datasets used.*

## 4 Models used

### 4.1 IndicTrans2

We used IndicTrans2, a Transformer-based (Vaswani et al., 2023) Multilingual NMT model trained on the BPCC dataset,[3] as our baseline MT system. We used the `ai4bharat/indictrans2-indic-en-1B` and `ai4bharat/indictrans2-en-indic-1B` checkpoints for our systems. For building our MT systems we set the following hyperparameters:

- the data was tokenised to a fixed length of 128 tokens, where sequences longer than 128 tokens were truncated and shorter ones were padded to ensure uniform length across batches,

- the learning rate: $2 \times 10^{-5}$,

- the batch size: 16,

- the training ran for 3 epochs satisfying our stopping criterion,

- a weight decay of 0.01 for improving the model's generalisation capabilities.

---

[3] https://ai4bharat.iitm.ac.in/bpcc/

We fine-tuned the model in order to adapt it to the domain and styles of data of Assamese-to-English translation task.

## 4.2 GPT-4o

GPT-4o (OpenAI et al., 2024) is a language model from OpenAI based on Transformer, which serves as the foundation for many language models today. It comprises multiple layers of self-attention mechanisms and feed-forward neural networks, enabling the model to efficiently process and generate text sequences. The model has been trained on a diverse and extensive dataset, allowing it to capture various linguistic patterns and contextual knowledge. We used GPT-4o for our in-context learning strategy. We used the following set of hyper-parameters for our experiment: (i) the temperature was set to 0.2, which controls the randomness of the output, ensuring more deterministic responses, and (ii) all other hyperparameters were not explicitly set and were set to the default values.

## 5 Experiments

In this Section we discuss our experiments. As discussed in Section 4, we used IndicTrans2 as our baseline model. We selected this model as the baseline due to its superiority as far as translation performance on low-resource Indian languages like Assamese is concerned (cf. Figure 1). We evaluated our MT models using the test data described in Section 3. We used BLEU (Papineni et al. (2002)), chrF (Popović (2015)), WER, and COMET[4] (Rei et al. (2020)) metrics for evaluation. The following subsections describes our MT systems.

### 5.1 Assamese-to-English

#### 5.1.1 Primary

Our primary MT system for the Assamese-to-English translation task is the fine-tuned Indic-Trans2 model (cf. Section 4). In other words, we fine-tuned the baseline model on the domain data provided by the organisers. Our data sets were detailed in Section 3. We used the same set of hyperparameters that we described in Section 4.

#### 5.1.2 Contrastive System One

as for our second system, we implemented an in-context few-shot learning approach, using which we generated English translations of Assamese sentences using OpenAI's GPT-4o.

More specifically, for few-shot learning we create prompts for the model with a few samples of translation pairs (source and target) whose source-side is similar to the source sentence we want to translate. We will now explain how we obtained training instances, whose source-side is similar to the sentence to be translated. We first convert all the Assamese training set sentences into dense vector embeddings using `sentence-transformers/all-MiniLM-L6-v2` [5]. The resulting embeddings were then indexed using FAISS,[6] enabling efficient similarity searches to retrieve the most relevant examples.

Furthermore, for each Assamese sentence of the test set, we used FAISS to retrieve the top five closest sentences from the training data based on the cosine similarity of their embeddings. Then, we constructed a detailed prompt for the GPT-4o model using the sentence-pairs that were retrieved from the training set. In Figure 2, we show an example of prompt used for in-context learning.

#### 5.1.3 Contrastive system two

For building our second Contrastive system we used our primary MT model (see Section 5.1.1) as our baseline. We adapted this MT system by fine-tuning it with a synthetic data. In order to create the synthetic data, we used a large English corpus comprising 5,000K sentences from the OPUS repository's NLLB project. We took 500k sentences for that large corpus for our experiment. We further filtered the sentences to include only those whose lengths are of 100 to 500 characters. With this, we omitted very short and excessively long sentences.

To extract domain-similar sentences from the now filtered corpus, we performed a semantic search on it using the validation set. All the corpus sentences were first converted to 768-dimensional dense vector embeddings using the `sentence-transformers-qa-mpnet-base-dot-v1` [7] model. We chose the `qa-mpnet-base-dot-v1` model over the `all-MiniLM-L6-v2` model used in 5.1.2 because it can store more detailed information about a sentence and capture semantic relationships across a wide range of contexts,

---

[4]COMET version 3.19.1 supports Assamese language.

[5]sentence-transformers/all-MiniLM-L6-v2: https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

[6]FAISS: https://github.com/facebookresearch/faiss

[7]sentence-transformers-qa-mpnet-base-dot-v1: https://huggingface.co/sentence-transformers/multi-qa-mpnet-base-dot-v1

Figure 1: *A visual representation of the advancements in machine translation systems for Indic languages using the IN22-Gen Evaluation set in the En-Indic direction. IT1, IT2 refers to IndicTrans1 and IndicTrans2 respectively. Negative chrF++ values indicate poor translation quality or situations where the translation system fails to generate meaningful or accurate translations. Adapted from (Gala et al., 2023)*

essential for extracting accurate and richer sentence representations from OPUS. These sentence embeddings were then indexed using FAISS. Later, we performed similarity searches by querying the FAISS index with embeddings of the validation sentences. We retrieved the top five nearest neighbours from the corpus for each validation sentence based on cosine similarity. We then removed sentences with similarity scores below 0.2. This process ensured that only the most relevant and contextually similar sentences were selected.

The final fuzzy matching English sentences were then back-translated into Assamese using our primary checkpoint. These new English-Assamese sentence pairs were used to create a new checkpoint by fine-tuning the primary system checkpoint and translation capabilities.

### 5.2 English to Assamese

#### 5.2.1 Primary

We build out primary systems using an MT approach similar to the one we used in Constrative system one (5.1.2) of the Assamese-to-English translation section, where we utilised OpenAI's GPT-4o model. The primary difference lies in the prompt structure. In Figure 3, we show the sample

prompt that was modified to treat English sentences as inputs and Assamese sentences as outputs.

## 6 Results

This section presents the evaluation results of the MT systems for both the Assamese-to-English and English-to-Assamese tasks. We performed the initial evaluation using the test pairs from the WMT 2023 dataset. Additionally, we present the results of our evaluation of the blind test set provided by the organisers. The results are reported in terms of BLEU, chrF, WER, and COMET metrics.

To ensure the reliability of our findings using the 2023 dataset, we conducted a statistical evaluation across pairs of models. This involved using bootstrap resampling (Koehn, 2004), calculating BLEU scores, and performing paired t-tests. For each test, we generated 100 bootstrap samples, each containing 100 randomly selected sentences from the dataset without repetition. This method maintains the original dataset's integrity while ensuring diversity in each sample. The results of these statistical analyses are also presented in this section. In all comparisons, we tested the null hypothesis that there is no difference in performance between the systems by calculating p-values. A low p-value (less than 0.05) indicates that we can reject the

```
Give only the final English
   sentence in a single line.
Context:
Assamese 1: <Assamese sentence 1>
Translation in English 1: <
   English translation 1>
Assamese 2: <Assamese sentence 2>
Translation in English 2: <
   English translation 2>
...
Assamese 5: <Assamese sentence 5>
Translation in English 5: <
   English translation 5>

What is the English translation
   for Assamese: <input sentence
   >?
```

Figure 2: *Prompt structure for GPT-4o model: Assamese-to-English*

```
Give only the final Assamese
   sentence in a single line.
Context:
English 1: <English sentence 1>
Translation in Assamese 1: <
   Assamese translation 1>
English 2: <English sentence 2>
Translation in Assamese 2: <
   Assamese translation 2>
...
English 5: <English sentence 5>
Translation in Assamese 5: <
   Assamese translation 5>

What is the Assamese translation
   for English: <input sentence>?
```

Figure 3: *Prompt structure for GPT-4o model: English to Assamese*

null hypothesis, suggesting that the observed differences are statistically significant and not due to random variation.

Four models were evaluated for the Assamese-to-English translation task: the baseline, primary, contrastive model one and contrastive model two. The evaluation results are summarised in Table 2.

| Model | BLEU ↑ | chrF ↑ | WER ↓ | COMET ↑ |
|-------|--------|--------|-------|---------|
| B | 0.2946 | 0.5646 | 0.7000 | 0.8064 |
| P | 0.3418 | 0.5748 | 0.6455 | 0.8086 |
| C1 | 0.3110 | 0.5690 | 0.7035 | 0.8157 |
| C2 | 0.3221 | 0.5724 | 0.6556 | 0.8075 |

Table 2: *Evaluation Results for Assamese-to-English Translation using WMT2023 test pair.*
*B = Baseline, P = Primary (5.1.1), C1 = Contrastive 1 (5.1.2), C2 = Contrastive 2 (5.1.3). ↑ indicates higher is better, and ↓ indicates lower is better.*

As shown in Table 2, the primary model (P) outperforms the baseline (B) in all metrics except COMET, where Contrastive system one (C1) achieves slightly higher scores than Contrastive system two (C2) . The BLEU and WER improvements suggest that the primary MT model provides more accurate and fluent translations compared to those by the baseline and contrastive models.

The statistical analysis further supports these findings. When comparing model **B** and **P**, the BLEU score of **P** (0.3418) was higher than that of

**B** (0.2946), with a t-statistic of -10.71 and a p-value of 1.72e-09. Similarly, when comparing **P** to **C1** (0.3110), the t-statistic was -10.17 with a p-value of 7.70e-20. In the comparison with **C2** (0.3221), the t-statistic was -8.64, and the p-value was 5.25e-08. Across all comparisons, the null hypothesis was rejected, indicating that **p** consistently performed better than the other models.

For the English-to-Assamese translation task, two models were evaluated: the baseline and the primary model. The results are summarised in Table 3.

| Model | BLEU ↑ | chrF ↑ | WER ↓ | COMET ↑ |
|-------|--------|--------|-------|---------|
| B | 0.1432 | 0.4948 | 0.8105 | 0.8263 |
| P | 0.1768 | 0.4815 | 0.7457 | 0.8220 |

Table 3: *Evaluation Results for English-to-Assamese Translation using WMT2023 test pair.*
*B = Baseline, P = Primary (5.2.1). ↑ indicates higher is better, and ↓ indicates lower is better.*

In Table 3, the primary model shows a noticeable improvement over the baseline in BLEU and WER, indicating better translation accuracy and reduced word errors. However, the chrF and COMET scores are slightly lower than those of the baseline.

The statistical significance tests compares Baseline and Primary (BLEU scores of 0.1432 for the Baseline and BLEU scores of 0.1768 for the Primary) with a t-statistic of -53.11 and a p-value of 1.45e-74. These results clearly indicate that Pri-

mary produces better translations than those by the Baseline. The null hypothesis, which assumes no difference in performance between the two systems, was rejected, confirming that the Primary system outperforms the Baseline.

We now present our results on the blind test set provided by the WMT organisers. The results for Assamese-to-English translation in WMT24 Low-Resource Indic Language Translation Task are summarised in Table 4. We observe that contrastive system two generally achieves the best results, leading to 0.3227 BLEU, 0.7563 METEOR, and 0.6573 chrF points, indicating better overall translation quality and semantic accuracy. The primary system closely follows the best-performing system (contrastive system two), performing slightly better in TER (33.56 points) and RIBES (0.3778 points), suggesting that translations require fewer edits, though it falls slightly behind contrastive system two on other metrics (0.3180 BLEU, 0.7537 METEOR, and 0.6551 chrF points). Contrastive system one consistently underperforms the other two systems, with lower scores across all metrics, particularly 39.03 TER and 0.7219 METEOR points.

| Model | BLEU ↑ | TER ↓ | RIBES ↑ | METEOR ↑ | chrF ↑ |
|-------|--------|-------|---------|----------|--------|
| P  | 0.3180 | 33.56 | 0.3778 | 0.7537 | 0.6551 |
| C1 | 0.2981 | 39.03 | 0.3713 | 0.7219 | 0.6437 |
| C2 | 0.3227 | 33.63 | 0.3720 | 0.7563 | 0.6573 |

Table 4: *Evaluation Results for Assamese-to-English Translation (2024).*
***P** = Primary (5.1.1), **C1** = Contrastive 1 (5.1.2), **C2** = Contrastive 2 (5.1.3). ↑ indicates higher is better, and ↓ indicates lower is better.*

The results for English-to-Assamese translation in WMT24 Low-Resource Indic Language Translation Task are summarised in Table 5. We only had one system for this direction, where we obtained 0.1612 BLEU, 65.96 TER, 0.2641 RIBES, 0.3927 METEOR, and 0.5673 chrF points on the test set.

| Model | BLEU ↑ | TER ↓ | RIBES ↑ | METEOR ↑ | chrF ↑ |
|-------|--------|-------|---------|----------|--------|
| P | 0.1612 | 65.96 | 0.2641 | 0.3927 | 0.5673 |

Table 5: *Evaluation Results for English-to-Assamese. **P** = Primary (5.2.1). ↑ indicates higher is better, and ↓ indicates lower is better.*

## 7 Conclusion

In this work, we presented our MT models developed for the WMT 2024 Low Resource Indic Translation Task, focusing on the Assamese-to-English

and English-to-Assamese language pairs. We conducted a comparative analysis using experimental setups to explore strategies such as fine-tuning, back translation, and in-context learning with few-shot prompting. All of these methods demonstrated significant performance improvements in translation.

For our future work, we intend to investigate synthetic pivoting methods for Indic languages and implement QLoRA technique to improve our current in-context learning approach, both discussed in Section 2. We believe that these techniques hold the potential to address the challenges associated with low-resource language translation and further improve the performance of our models.

## References

Khalid Ahmed and Jan Buys. 2024. Neural machine translation between low-resource languages with synthetic pivoting. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 12144–12158, Torino, Italia. ELRA and ICCL.

Mazida Ahmed, Kuwali Talukdar, Parvez Boruah, Prof. Shikhar Kumar Sarma, and Kishore Kashyap. 2023. GUIT-NLP's submission to shared task: Low resource Indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 935–940, Singapore. Association for Computational Linguistics.

Laurie Burchell, Birch, and Alexandra Kenneth Heafield. 2022. Exploring diversity in back translation for low-resource machine translation. In *Proceedings of the Third Workshop on Deep Learning for Low-Resource Natural Language Processing*. Association for Computational Linguistics.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms.

Jay Gala, Pranjal A. Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages.

Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pages 388–395.

Yasmin Moslem, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. Adaptive machine translation with large language models.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the wmt 2023 shared task on low-resource indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Dhairya Suman, Atanu Mandal, Santanu Pal, and Sudip Naskar. 2023. IACS-LRILT: Machine translation for low-resource Indic languages. In *Proceedings of the Eighth Conference on Machine Translation*, pages 972–977, Singapore. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need.

Xuan Zhang, Navid Rajabi, Kevin Duh, and Philipp Koehn. 2023. Machine translation with large language models: Prompting, few-shot learning, and fine-tuning with QLoRA. In *Proceedings of the Eighth Conference on Machine Translation*, pages 468–481, Singapore. Association for Computational Linguistics.

# SPRING Lab IITM's submission to Low Resource Indic Language Translation Shared Task

**Hamees Sayed, Advait Joglekar, Srinivasan Umesh**
SPRING Lab,
Indian Institute of Technology Madras
hameessayed71@gmail.com, advaitjoglekar@gmail.com, umeshs@ee.iitm.ac.in

## Abstract

We develop a robust translation model for four low-resource Indic languages: Khasi, Mizo, Manipuri, and Assamese. Our approach includes a comprehensive pipeline from data collection and preprocessing to training and evaluation, leveraging data from WMT task datasets, BPCC, PMIndia, and OpenLanguage-Data. To address the scarcity of bilingual data, we use back-translation techniques on monolingual datasets for Mizo and Khasi, significantly expanding our training corpus. We fine-tune the pre-trained NLLB 3.3B model for Assamese, Mizo, and Manipuri, achieving improved performance over the baseline. For Khasi, which is not supported by the NLLB model, we introduce special tokens and train the model on our Khasi corpus. Our training involves masked language modelling, followed by fine-tuning for English-to-Indic and Indic-to-English translations.

## 1 Introduction

Translation of low-resource languages poses significant challenges in natural language processing. While substantial progress has been made in developing machine translation models for high-resource languages, low-resource languages often suffer from a lack of parallel corpora and digital resources (Haddow et al., 2022). Languages like Khasi, Mizo, Manipuri, and Assamese are representative of this challenge, where limited data and unique linguistic complexities hinder the development of robust translation systems.

In recent years, efforts to bridge this gap have gained momentum, driven by initiatives such as the Bharat Parallel Corpus Collection[1] (BPCC) (Gala et al., 2023) and government-supported projects like PMIndia (Haddow and Kirefu, 2020), which aim to provide bilingual data for Indic languages.

Despite these efforts, translation models for low-resource Indic languages have yet to achieve performance levels comparable to their high-resource counterparts (Suman et al., 2023), necessitating innovative approaches to model training and data utilization.

In this work, we develop a robust translation model for four low-resource Indic languages: Khasi, Mizo, Manipuri, and Assamese. Our approach involves data collection, preprocessing, training, and evaluation. We utilize datasets from WMT, BPCC, PMIndia, and OpenLanguage-Data[2] (Maillard et al., 2023), and enhance bilingual data through back-translation (Edunov et al., 2018) techniques, especially for Mizo and Khasi, significantly expanding our training corpus.

We follow Meta's data preprocessing standards and use LoRA (Low-Rank Adaptation) (Hu et al., 2021) fine-tuning on the NLLB (et al., 2022) 3.3B model to improve efficiency and performance with fewer parameters. Our model initially focuses on one-way translation from English to the Indic languages, then on reverse translations (Dabre et al., 2019). The results show improved performance over the baseline, particularly for Khasi, where we address gaps in pre-trained model support.

## 2 Dataset

In this study, we focus on four low-resource Indic languages covered in the Low Resource Indic Languages Shared Task: Khasi, Mizo, Manipuri, and Assamese. This section highlights the significance of each language, including their role in their respective regions, their linguistic and cultural importance, and the details of the datasets used. Statistics regarding language speakers are according to the 2011 Indian Census[3].

---

[1] https://ai4bharat.iitm.ac.in/bpcc/

[2] https://github.com/openlanguagedata/seed
[3] https://censusindia.gov.in/

| Language | ISO-693-3 | WMT Parallel | BPCC | PMIndia | OLD | Back-Translated | Total |
|----------|-----------|--------------|------|---------|-----|-----------------|-------|
| Assamese | asm | 50,000 | 35,354 | 9,732 | 0 | 0 | 95,086 |
| Manipuri | mni | 21,687 | 0 | 7,419 | 6,193 | 0 | 35,036 |
| Khasi | kha | 24,000 | 0 | 0 | 0 | 102,070 | 126,070 |
| Mizo | lus | 50,000 | 0 | 0 | 0 | 30,164 | 80,164 |

Table 1: Breakdown of data sources and volumes for each language. "OLD" refers to OpenLanguageData. The "Back-Translated" data was initially generated using Google Translate[4] for the first 500k characters from the monolingual WMT task data, with subsequent iterations increasing the data size using the trained model.

## 2.1 Languages

**Assamese** *(Asamiya)* is an Indo-Aryan language spoken primarily in the northeastern Indian state of Assam, where it serves as an official language and a regional lingua franca. With over 15 million native speakers, it is one of the most widely spoken languages in the region. Historically, Assamese was the court language of the Ahom kingdom. It is written in the Assamese script, an abugida system, known for its unique typographic ligatures.

**Manipuri** *(Meiteilon)* is a key Tibeto-Burman language spoken mainly in Manipur, India, where it is an official language and it is one of the constitutionally scheduled official languages of the Indian Republic. With 1.76 million speakers, it is the most widely spoken Tibeto-Burman language in India and holds the third place among the fastest-growing languages of India, following Hindi and Kashmiri. It is written in its own Meitei script as well as the Bengali script.

**Khasi** *(Ka Ktien Khasi)* is an Austroasiatic language primarily spoken by the Khasi people in Meghalaya, India, with approximately 1 million native speakers as of the 2011 census. The language holds an associate official status in certain districts of Meghalaya. Khasi is written in the Latin script. The closest relatives of Khasi are other languages in the Khasic group, such as Pnar and War.

**Mizo** *(Mizo ṭawng)* belonging to the Sino-Tibetan language family, is primarily spoken in the state of Mizoram, India, with around 800 thousand speakers. The Mizo language, also known as Lushai, has a rich oral history and was first written using the Latin script in the late 19th century. Mizo is recognized as the official language of Mizoram and is used in education, government, and media.

## 3 Methodology

This section covers the preprocessing steps and training methods used, including dataset preparation and the fine-tuning of Meta's multilingual NLLB 3.3B base pre-trained model. Detailed statistics on data distribution are presented in Table 1.

### 3.1 Preprocessing

In the preprocessing phase, we followed a series of steps to ensure the text data was clean and consistent before model training. We began by normalizing punctuation using Moses (Koehn et al., 2007), an open-source toolkit designed for preprocessing, training, and testing translation models. This step helps maintain consistency in text data, which is crucial for training robust models.

Non-printable characters, which often interfere with text processing, were replaced with a space. This choice ensures that any invisible or non-standard characters do not disrupt the tokenization process and ensures that the text is composed of standard printable characters.

We also applied Unicode normalization (NFKC) to transform characters into their canonical forms, making the text more uniform across different Unicode representations.

These preprocessing steps are aligned with those outlined by Meta for their multilingual models, and further details can be found on their GitHub[5]. This approach ensures that the text data used for training is clean, consistent, and compatible with the modelling requirements.

### 3.2 Training

For model training, we employed Meta's NLLB (No Language Left Behind) 3.3B parameter model,

---

[4] https://google.translate.com/
[5] https://github.com/facebookresearch/stopes/blob/main/stopes/pipelines/monolingual/monolingual_line_processor.py

a state-of-the-art multilingual machine translation model built to support over 200 languages, making it ideal for tasks involving low-resource languages (Tran et al., 2021; Yang et al., 2021). The NLLB 3.3B model is based on a Transformer (Vaswani et al., 2023) architecture with 3.3 billion parameters, featuring a dense encoder-decoder design. It includes the following hyperparameters:

| Hyperparams | |
| --- | --- |
| embed size | 2048 |
| ffn size | 8192 |
| attn heads | 16 |
| enc/dec layers | 48 |

Table 2: Hyperparameters for the baseline pre-trained model. 24 Encoder and 24 Decoder Layers.

To fine-tune the model, we employed LoRA, a technique that significantly reduces computational demands and training time by adapting only a small subset of the model's parameters. LoRA has been shown to match the performance of traditional fine-tuning methods while reducing the number of trainable parameters by a factor of 50 (Alves et al., 2023). This approach is especially effective for large-scale models like Meta's NLLB 3.3B, allowing efficient adaptation without significantly compromising on performance.

### 3.3 Parameters

The training process was conducted in three stages: first, the model was trained on masked language modelling (Devlin et al., 2019) to enhance its understanding of the target language by leveraging monolingual data. Next, it was fine-tuned for English-to-Indic translations, followed by further fine-tuning for Indic-to-English translations. In the case of Khasi, which was not natively supported by the NLLB model, special tokens were added to the tokenizer's vocabulary to accommodate the Khasi language. The model was subsequently trained on the Khasi corpus to ensure proper handling and integration of this language.

The training was performed across 4 Nvidia A6000 GPUs. These settings allowed us to optimize the model's performance while managing computational efficiency.

### 3.4 Inference

For inference, the trained adapter was loaded onto the NLLB 3.3B model. The model generated

| Training Args | |
| --- | --- |
| optimizer | adafactor |
| learning Rate | 1e-5 |
| epochs | 8 |
| precision | bf16 |
| $p_{mask}$ | 0.15 |
| peft type | lora |
| rank | 128 |
| lora alpha | 256 |
| lora dropout | 0.1 |
| target modules | all linear |

Table 3: Training parameters and LoRA configuration used for fine-tuning the NLLB 3.3B model.

predictions using a beam search strategy with 10 beams and a repetition penalty of 2.5 to improve the diversity and quality of the translations. We experimented with various beam and penalty configurations, ultimately finding that this particular setup produced the most accurate and linguistically coherent outputs.

## 4 Results

The evaluation of our translation model across various language pairs and directions is shown in Table 4, with performance assessed using BLEU (Papineni et al., 2002), Translation Error Rate (Snover et al., 2006), RIBES (Isozaki et al., 2010), METEOR (Banerjee and Lavie, 2005), and chrF (Popović, 2015) metrics. We found that the scores in English-to-Manipuri and English-to-Mizo direction suffered from the poor quality of backtranslated data used in our training.

**English-Assamese** The model performed relatively well, with BLEU scores of 27.26 for English-to-Assamese and 26.69 for Assamese-to-English.

**English-Manipuri** The model showed lower BLEU scores for English-to-Manipuri (2.7) compared to Manipuri-to-English (20.88). The TER score was higher for English-to-Manipuri, reflecting greater translation errors in this direction.

**English-Khasi** For Khasi, the BLEU score was 12.12 for English-to-Khasi and 10.47 for Khasi-to-English.

**English-Mizo** The performance was mixed, with a BLEU score of 6.6 for English-to-Mizo and 18.49 for Mizo-to-English. The TER score indicates a higher error rate for English-to-Mizo, while the METEOR and ChrF scores were relatively balanced across both directions.

| Language Pairs | Test Set | BLEU | TER | RIBES | METEOR | ChrF |
|---|---|---|---|---|---|---|
| English-Assamese | en_to_as_contrastive | 27.26 | 52.79 | 0.3032 | 0.513 | 65.2 |
| | as_to_en_contrastive | 26.69 | 39.08 | 0.3308 | 0.7066 | 60.48 |
| English-Manipuri | en_to_mn_contrastive | 2.7 | 84.6 | 0.1185 | 0.1567 | 44.28 |
| | mn_to_en_contrastive | 20.88 | 48.77 | 0.3031 | 0.61 | 53.64 |
| English-Khasi | en_to_kh_contrastive | 12.12 | 63.31 | 0.1864 | 0.4453 | 44.55 |
| | kh_to_en_contrastive | 10.47 | 61.43 | 0.2172 | 0.5042 | 42.71 |
| English-Mizo | en_to_mz_contrastive | 6.6 | 66.06 | 0.1746 | 0.495 | 49.79 |
| | mz_to_en_contrastive | 18.49 | 53.19 | 0.2684 | 0.588 | 50.44 |

Table 4: Translation performance metrics of our MT System reported in the final evaluation.

## 5 Conclusion

In this work, we utilized Meta's NLLB 3.3B model, a large-scale multilingual transformer with 3.3 billion parameters, to enhance translation between low-resource Indic languages and English. The training process included masked language modelling, followed by English-to-Indic and Indic-to-English translations. Special tokens were added for Khasi, and LoRA (Low-Rank Adaptation) was employed to optimize computational efficiency and reduce training time.

Conducted on 4 NVIDIA A6000 GPUs, our approach demonstrates that large-scale multilingual models, when combined with LoRA, effectively capture diverse linguistic patterns and advance translation capabilities.

## 6 Limitations

In this study, we encountered several limitations that impacted the overall effectiveness of our translation system. One major challenge was the constrained size of our dataset due to computational resource limitations. The limited dataset size may have hindered the model's ability to generalize, particularly for low-resource languages where larger and more diverse datasets would have been advantageous.

Another issue we faced was the quality of back-translated data. The process of augmenting the dataset through machine translation often resulted in lower-quality data, which negatively influenced the model's performance. This highlights the need for more robust data generation techniques in future work.

We also observed a noticeable performance gap between translations where English was the target language and those where an Indic language was the target. This suggests that while the model may understand the morphological aspects of Indic languages, it struggles to generate accurate translations in these languages. This limitation underscores the need for further refinement in handling the complexities of Indic language generation.

Finally, the potential domain mismatch between our training data and real-world applications posed a significant challenge. The training data may not fully capture the linguistic and contextual nuances found in practical scenarios, leading to reduced performance in actual use cases. Addressing this issue in future work will be crucial for improving the model's real-world applicability.

## References

Duarte M. Alves, Nuno M. Guerreiro, João Alves, José Pombal, Ricardo Rei, José G. C. de Souza, Pierre Colombo, and André F. T. Martins. 2023. Steering large language models for machine translation with finetuning and in-context learning. *Preprint*, arXiv:2310.13448.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Raj Dabre, Atsushi Fujita, and Chenhui Chu. 2019. Exploiting multilingualism through multistage finetuning for low-resource neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1410–1416, Hong Kong, China. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *Preprint*, arXiv:1810.04805.

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. *Preprint*, arXiv:1808.09381.

NLLB Team et al. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.

Jay Gala, Pranjal A. Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Preprint*, arXiv:2305.16307.

Barry Haddow, Rachel Bawden, Antonio Valerio Miceli Barone, Jindřich Helcl, and Alexandra Birch. 2022. Survey of low-resource machine translation. *Computational Linguistics*, 48(3):673–732.

Barry Haddow and Faheem Kirefu. 2020. Pmindia – a collection of parallel corpora of languages of india. *Preprint*, arXiv:2001.09907.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *Preprint*, arXiv:2106.09685.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 944–952, Cambridge, MA. Association for Computational Linguistics.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.

Jean Maillard, Cynthia Gao, Elahe Kalbassi, Kaushik Ram Sadagopan, Vedanuj Goswami, Philipp Koehn, Angela Fan, and Francisco Guzmán. 2023. Small data, big impact: Leveraging minimal data for effective machine translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2740–2756, Toronto, Canada. Association for Computational Linguistics.

Partha Pakray, Santanu Pal, Advaitha Vetagiri, Reddi Mohana Krishna, Arnab Kumar Maji, Sandeep Kumar Dash, Lenin Laitonjam, Lyngdoh Sarah, and Riyanka Manna. 2024. Findings of wmt 2024 shared task on low-resource indic languages translation. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the wmt 2023 shared task on low-resource indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Dhairya Suman, Atanu Mandal, Santanu Pal, and Sudip Naskar. 2023. IACS-LRILT: Machine translation for low-resource Indic languages. In *Proceedings of the Eighth Conference on Machine Translation*, pages 972–977, Singapore. Association for Computational Linguistics.

Chau Tran, Shruti Bhosale, James Cross, Philipp Koehn, Sergey Edunov, and Angela Fan. 2021. Facebook ai wmt21 news translation task submission. *Preprint*, arXiv:2108.03265.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *Preprint*, arXiv:1706.03762.

Jian Yang, Shuming Ma, Haoyang Huang, Dongdong Zhang, Li Dong, Shaohan Huang, Alexandre Muzio, Saksham Singhal, Hany Hassan, Xia Song, and Furu Wei. 2021. Multilingual machine translation systems from Microsoft for WMT21 shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 446–455, Online. Association for Computational Linguistics.

# Machine Translation Advancements of Low-Resource Indian Languages by Transfer Learning

**Bin Wei, Jiawei Zheng, Zongyao Li, Zhanglin Wu, Daimeng Wei,**
**Jiaxin Guo, Zhiqiang Rao, Shaojun Li, Yuanchang Luo, Hengchao Shang,**
**Jinlong Yang, Yuhao Xie, Hao Yang**

Huawei Translation Service Center, Beijing, China
{weibin29, zhengjiawei15, lizongyao, wuzhanglin2, weidaimeng,
guojiaxin1,raozhiqiang, lishaojun18, luoyuanchang1, shanghengchao,
yangjinlong7, xieyuhao2, yanghao30}@huawei.com

## Abstract

This paper introduces the submission by Huawei Translation Center (HW-TSC) to the WMT24 Indian Languages Machine Translation (MT) Shared Task. To develop a reliable machine translation system for low-resource Indian languages, we employed two distinct knowledge transfer strategies, taking into account the characteristics of the language scripts and the support available from existing open-source models for Indian languages. For Assamese(as) and Manipuri(mn), we fine-tuned the existing IndicTrans2(Gala et al., 2023) open-source model to enable bidirectional translation between English and these languages. For Khasi (kh) and Mizo (mz), We trained a multilingual model as a baseline using bilingual data from these four language pairs, along with an additional about 8kw English-Bengali bilingual data, all of which share certain linguistic features. This was followed by fine-tuning to achieve bidirectional translation between English and Khasi, as well as English and Mizo. Our transfer learning experiments produced impressive results: 23.5 BLEU for en→as, 31.8 BLEU for en→mn, 36.2 BLEU for as→en, and 47.9 BLEU for mn→en on their respective test sets. Similarly, the multilingual model transfer learning experiments yielded impressive outcomes, achieving 19.7 BLEU for en→kh, 32.8 BLEU for en→mz, 16.1 BLEU for kh→en, and 33.9 BLEU for mz→en on their respective test sets. These results not only highlight the effectiveness of transfer learning techniques for low-resource languages but also contribute to advancing machine translation capabilities for low-resource Indian languages.

## 1 Introduction

In the realm of machine translation, Neural Machine Translation (NMT) has become the dominant technology, as confirmed by previous research. However, training NMT models requires large amounts of data, which presents a significant challenge when dealing with low-resource languages.

To tackle this challenge, we employed transfer learning, a well-established approach that enhances model performance by transferring knowledge gained from one task to other related tasks. To improve translation capabilities for low-resource languages, we faced the challenge of limited bilingual resources for Indian languages. To overcome this issue, we trained a multilingual model using not only all the bilingual data provided for the task but also additional Bengali data. Additionally, we examined the languages supported by the existing IndicTrans2(Gala et al., 2023) open-source model and conducted a comparative analysis. Based on our findings, we selected different baseline models for knowledge transfer depending on the language pair: for Assamese and Manipuri, we used the IndicTrans2 model as the baseline, while for Khasi and Mizo, we trained multilingual model by ourselves as the baseline. This approach enabled us to effectively leverage existing resources while addressing the specific challenges associated with each language pair.

IndicTrans2 is the first open-source transformer-based multilingual NMT model that supports high-quality translations across all the 22 scheduled Indic languages. It was trained on the extensive Bharat Parallel Corpus Collection (BPCC), a publicly accessible repository encompassing both pre-existing and freshly curated data for all 22 scheduled Indian languages, this model boasts a comprehensive understanding of the linguistic diversity within the Indian subcontinent. To enhance its linguistic prowess, IndicTrans2 has undergone auxiliary training utilizing the rich resource of back-translated monolingual data. The model was then trained on human-annotated data to achieve further improvements. We used this model in the first two subtasks and fine-tuned it on the training data provided by WMT24. By adopting this approach, we aim to capitalize on the acquired knowledge during training to significantly bolster the performance of

the model in the specific translation task at hand. The fine-tuned IndicTrans2 achieves good scores, so we are using it for our final submission in the first two subtasks.

For the multilingual model, we first utilized resources from Bengali. The choice of Bengali was based on its belonging to the Indo-Aryan branch, its linguistic feature similarities with some of the target low-resource languages, and its relatively rich available data. By introducing Bengali data, we aimed to enable the model to learn features potentially shared with the target languages, thereby laying a foundation for processing other related languages. Next, we integrated all available bilingual data from Indic language MT track. This included parallel corpora between various Indian languages and English. Although the data for each language pair might be limited individually, the combined dataset offered diverse learning samples. We believe that this integration of multilingual data helps the model capture both the commonalities and differences among different Indian languages. Based on this carefully selected and integrated data, we trained a multilingual model. The design goal of this model was to handle translation tasks for multiple Indian languages simultaneously, using the commonalities between languages to compensate for the scarcity of data in any single language. Through this approach, we expect the model to learn more generalized language representations and translation knowledge under resource constraints, leading to improved performance on Khasi and Mizo translation tasks.

Ultimately, we adopted a differentiated strategy for knowledge transfer. This approach thoroughly considered the characteristics of each language to achieve optimal transfer effects. In Section 2, we will discuss the details of the data, the methods and processes used for data pre-processing. Section 3 will cover the overall architecture and training strategies of the NMT system, including a detailed account of the various optimization methods. In Section 4, we will present the experimental parameters, results, and their analysis. The final section will summarize the key findings of the paper.

## 2 Data

### 2.1 Data Details

We have fine-tuned the model using the WMT24 corpus. Additionally, we used 2M monolingual english dataset to do BT and FT. The amount of

data we used is shown in Tables 1.

| language pairs | bitext data | monolingual data |
|:---:|:---:|:---:|
| en-as | 50K | en: 2M, as: 2.62M |
| en-mn | 21K | en: 2M, mn: 2.14M |
| en-kh | 24K | en: 2M, kh: 182K |
| en-mz | 50K | en: 2M, mz: 1.9M |

Table 1: Bilingual and monolingual used for training NMT models.

### 2.2 Data Pre-processing

Our data pre-processing methods for NMT include:

- Remove duplicate sentences or sentence pairs.

- Convert full-width symbols to half-width.

- Use fasttext[1] (Joulin et al., 2016) to filter other language sentences.

- Use mosesdecoder[2] (Koehn et al., 2007) to normalize English punctuation.

- Filter out sentences with more than 150 words.

- Use fast-align (Dyer et al., 2013) to filter sentence pairs with poor alignment.

- Sentencepiece[3] (SPM) (Kudo and Richardson, 2018) is used to perform subword segmentation, and the vocabulary size is set to 32K.

Since there may be some semantically dissimilar sentence pairs in bilingual data, we use LaBSE[4] (Feng et al., 2022) to calculate the semantic similarity of each bilingual sentence pair, and exclude bilingual sentence pairs with a similarity score lower than 0.75 from our training corpus.

## 3 NMT System

### 3.1 System Overview

We use Transformer (Vaswani, 2017) as our neural machine translation (NMT)(Bahdanau et al., 2014) model architecture. For the first two subtasks(en-as, en-mn), we use the IndicTrans2(Gala et al., 2023)

---

[1] https://github.com/facebookresearch/fastText
[2] https://github.com/moses-smt/mosesdecoder
[3] https://github.com/google/sentencepiece
[4] https://huggingface.co/sentence-transformers/LaBSE

Figure 1: The overall training flow of NMT system.

model as our baseline model, which is a deep Transformer architecture with 18-layers encoder and 18-layers decoder. with the latter two subtasks(en-kh, en-mz), we trained a multilingual model as our baseline model, which is a deep Transformer architecture with 35-layers encoder and 3-layers decoder.

Fig. 1 shows the overall training flow of NMT system. Referred to previous work (Wei et al., 2021, 2022; Wu et al., 2023), We use training strategies such as regularized dropout (R-Drop) (Wu et al., 2021), data diversification (DD) (Nguyen et al., 2020), forward translation FT) (Abdulmumin, 2021), back translation (BT) (Sennrich et al., 2016), denoise, Transfer learning(TL) and transductive ensemble learning (TEL) (Wang et al., 2020) for training.

## 3.2 Regularized Dropout

Regularized Dropout (R-Drop)[5] (Wu et al., 2021) is a simple yet more effective alternative to regularize the training inconsistency induced by dropout (Srivastava et al., 2014). Concretely, in each mini-batch training, each data sample goes through the forward pass twice, and each pass is processed by a different sub model by randomly dropping out some hidden units. R-Drop forces the two distributions for the same data sample outputted by the two sub models to be consistent with each other, through minimizing the bidirectional Kullback-Leibler (KL) divergence (Van Erven and Harremos, 2014) between the two distributions. That is, R-Drop regularizes the outputs of two sub models randomly sampled from dropout for each data sample

---

[5] https://github.com/dropreg/R-Drop

in training. In this way, the inconsistency between the training and inference stage can be alleviated.

## 3.3 Data Diversification

Data Diversification (DD) (Nguyen et al., 2020) is a data augmentation method to boost NMT performance. It diversifies the training data by using the predictions of multiple forward and backward models and then merging them with the original dataset which the final NMT model is trained on. DD is applicable to all NMT models. It does not require extra monolingual data, nor does it add more parameters. To conserve training resources, we only use one forward model and one backward model to diversify the training data.

## 3.4 Forward Translation

Forward translation (FT) (Abdulmumin, 2021), also known as self-training, is one of the most commonly used data augmentation methods. FT has proven effective for improving NMT performance by augmenting model training with synthetic parallel data. Generally, FT is performed in three steps: (1) randomly sample a subset from the large-scale source monolingual data; (2) use a "teacher" NMT model to translate the subset data into the target language to construct the synthetic parallel data; (3) combine the synthetic and authentic parallel data to train a "student" NMT model.

## 3.5 Back Translation

An effective method to improve NMT with target monolingual data is to augment the parallel training data with back translation (BT) (Sennrich et al., 2016; Wei et al., 2023). There are many published works that expand the understanding of BT and investigate methods for generating synthetic source sentences. Edunov et al. (2018) find that back translations obtained via sampling or noised beam outputs are more effective than back translations generated by beam or greedy search in most scenarios. Caswell et al. (2019) show that the main role of such noised beam outputs is not to diversify the source side, but simply to tell the model that the given source is synthetic. Therefore, they propose a simpler alternative strategy: Tagged BT. This method uses an extra token to mark back translated source sentences, which generally outperforms noised BT (Edunov et al., 2018). For better joint use with FT, we use sampling back translation (ST) (Edunov et al., 2018).

## 3.6 Denoise

In machine translation, denoising improves translation quality by removing noise from the training data, such as inaccurate translations, grammatical errors, or unnatural sentence structures, allowing the model to focus on high-quality data and produce more accurate and fluent translations. Additionally, denoising enhances the model's robustness by eliminating noisy data, which helps the model better learn the target language's patterns, reducing errors and leading to more stable and reliable performance across diverse inputs. It also optimizes training efficiency by decreasing the amount of data the model needs to process, particularly by filtering out low-quality data, which results in a cleaner and more consistent dataset and can shorten the overall training time. Moreover, denoising reduces error propagation by preventing the model from learning incorrect language patterns, thereby minimizing the accumulation and spread of errors in generated translations. Finally, it enhances the model's generalization ability, as denoised data is more representative, enabling the model to better adapt to different types of input sentences and improving its performance in real-world applications. Through denoising, machine translation models can more effectively utilize high-quality data, leading to superior translation outcomes and greater overall model stability.

## 3.7 Transductive Ensemble Learning

Ensemble learning (Garmash and Monz, 2016), which aggregates multiple diverse models for inference, is a common practice to improve the performance of machine learning models. However, it has been observed that the conventional ensemble methods only bring marginal improvement for NMT when individual models are strong or there are a large number of individual models. Transductive Ensemble Learning (TEL) (Zhang et al., 2019) studies how to effectively aggregate multiple NMT models under the transductive setting where the source sentences of the test set are known. TEL uses all individual models to translate the source test set into the target language space and then fine-tune a strong model on the translated synthetic data, which significantly boosts strong individual models and benefits a lot from more individual models.

## 3.8 Transfer Learning

Transfer learning(TL) is a machine learning technique where a model trained on one task is adapted for a second related task. Instead of starting the training of a new model from scratch, transfer learning leverages the knowledge learned from the first task to improve learning on the second task. For Assamese(as) and Manipuri(mn), We have used IndicTrans2(Gala et al., 2023), a powerful model that performs well for English-to-Indic and Indic-to-English translation for 22 scheduled Indian languages. This knowledge can be used to translate other Indian languages to and from English. Our approach entailed the fine-tuning of this model, leveraging the parallel corpus provided by the WMT24 for the Indic MT task. This fine-tuning process equipped the model with the expertise required to proficiently translate Assamese and Manipuri to and from English, ultimately yielding the most outstanding results. Similarly, for Khasi and Mizo, we trained a multilingual model as the baseline. We also applied transfer learning techniques to enhance the baseline model using data specific to these language pairs. The results on both the test and dev sets were highly encouraging.

## 4 Experiment

### 4.1 Settings

We use Transformer architecture in all the subtasks. For the first two subtasks, we use IndicTrans2 (Gala et al., 2023) as our baseline model, which is a deep Transformer architecture with 18-layers encoder and 18-layers decoder. With the latter subtasks, the model is also a Transformer architecture with 35-layers encoder and 3-layers decoder. For the first two subtasks, our models apply Adam (Kingma and Ba, 2014) as optimizer to update the parameters with $\beta_1 = 0.9$ and $\beta_2 = 0.98$. We employ a warm-up learning rate of $10^{-7}$ for 2000 update steps and a learning rate of $3 * 10^{-5}$. For normalization, we use a dropout value of 0.2 and normalize the probabilities using smoothed label cross-entropy. We use GeLU activations (Hendrycks and Gimpel, 2016) for better learning. For the latter subtasks, parameter update frequency is 2, and learning rate is 5e-4. The number of warmup steps is 4000, and model is saved every 1000 steps. R-Drop (Wu et al., 2021) is used in model training for all subtasks, and we set $\lambda$ to 5.

We use the scareBLEU library to calculate our BLEU (Papineni et al., 2002) and ChrF (Popović,

| Language-pair | Training strategies | Bleu(test) | ChrF2(test) | Bleu(dev) | ChrF2(dev) |
|---|---|---|---|---|---|
| en→as | IndicTrans2 baseline | 18.9 | 51.4 | 14.7 | 44.8 |
| | + DD, FT, BT | 22.9 | 52.5 | 21.1 | 47.7 |
| | + denoise | 23.3 | 53.1 | 22.5 | 48.9 |
| | + TEL | 23.5 | 53.2 | 22.8 | 49.0 |
| en→mn | IndicTrans2 baseline | 11.9 | 48.5 | 11.9 | 48.5 |
| | + DD, FT, BT | 30.9 | 62.8 | 31.1 | 63.4 |
| | + denoise | 31.7 | 64.7 | 31.7 | 64.9 |
| | + TEL | 31.8 | 64.6 | 31.6 | 64.9 |
| as→en | IndicTrans2 baseline | 29.7 | 56.3 | 25.6 | 49.3 |
| | + DD, FT, BT | 35.8 | 58.6 | 35.0 | 54.5 |
| | + denoise | 36.1 | 58.6 | 34.8 | 54.6 |
| | + TEL | 36.2 | 59.4 | 33.7 | 54.2 |
| mn→en | IndicTrans2 baseline | 32.6 | 62.3 | 33.4 | 61.8 |
| | + DD, FT, BT | 47.5 | 70.8 | 47.0 | 69.7 |
| | + denoise | 47.7 | 70.8 | 47.2 | 69.7 |
| | + TEL | 47.9 | 70.8 | 47.4 | 69.8 |

Table 2: The results of en-as and en-mn language pairs on the test and dev set.

| Language-pair | Training strategies | Bleu(test) | ChrF2(test) | Bleu(dev) | ChrF2(dev) |
|---|---|---|---|---|---|
| en→kh | multilingual baseline | 17.4 | 40.4 | 17.0 | 39.7 |
| | + DD, FT, BT | 18.1 | 41.8 | 17.9 | 41.3 |
| | + denoise | 19.5 | 43.3 | 19.2 | 42.7 |
| | + TEL | 19.7 | 43.5 | 19.3 | 42.8 |
| en→mz | multilingual baseline | 25.0 | 51.6 | 22.3 | 46.6 |
| | + DD, FT, BT | 30.8 | 55.7 | 25.2 | 49.1 |
| | + denoise | 32.5 | 57.1 | 25.4 | 49.3 |
| | + TEL | 32.8 | 57.3 | 25.7 | 49.4 |
| kh→en | multilingual baseline | 15.1 | 37.7 | 15.0 | 38.1 |
| | + DD, FT, BT | 15.8 | 37.8 | 15.0 | 38.3 |
| | + denoise | 15.9 | 38.5 | 15.5 | 39.0 |
| | + TEL | 16.1 | 38.8 | 15.6 | 39.2 |
| mz→en | multilingual baseline | 26.7 | 48.2 | 22.9 | 44.0 |
| | + DD, FT, BT | 32.9 | 52.2 | 25.0 | 45.4 |
| | + denoise | 33.7 | 52.2 | 25.8 | 46.5 |
| | + TEL | 33.9 | 52.7 | 26.0 | 46.7 |

Table 3: The results of en-kh and en-mz language pairs on the test and dev set.

2015) scores with a word order of 2.

## 4.2 Results

Regarding this four language pair directions, we use Regularized Dropout, Bidirectional Training, Data Diversification, Forward Translation, Back Translation, Alternated Training, Curriculum Learning, and Transductive Ensemble Learning. The evaluation results of four language pair directions NMT system on WMT24 Indic MT test and dev set are shown in Tables 2 and Tables 3.

As shown in Table 2, IndicTrans2(Gala et al., 2023) provides a strong baseline. Fine-tuning the model with FT, BT, and bitext data leads to significant improvements, particularly in the en-mn direction, where the BLEU score increases by nearly 20 points over the baseline on the test and dev set. This improvement is largely attributed to Data Diversification. Table 3 further illustrates that FT and BT data contribute the most to model performance, especially in the en-mz direction, which sees an increase of nearly six BLEU points compared to the multilingual baseline. Even after enhancing the model with BT and FT data, adding filtered high-quality bilingual data results in an average gain of about one BLEU point, highlighting the critical role of data quality. Finally, we all use TEL technique to obtain a good result, the improvement is very small, almost less than one bleu score.

## 5 Conclusion

This paper presents the submission of HW-TSC to the WMT24 Indic MT Task. For the first two subtasks, we use IndicTrans2 as our baseline model to fine-tune it with corpus provided by WMT24 on the en-as and en-mn language pairs, which achieves remarkable performance. For the latter two subtasks, we train a multilingual model on the en-kh

and en-mz language pairs, and then use training strategies such as R-Drop, DD, FT, BT, denoise and TEL to train the NMT model based on the deep Transformer-big architecture. By applying these training strategies, our submission achieved a competitive result in the final evaluation.

## References

Idris Abdulmumin. 2021. Enhanced back-translation for low resource neural machine translation using self-training. In *Information and Communication Technology and Applications: Third International Conference, ICTA 2020, Minna, Nigeria, November 24–27, 2020, Revised Selected Papers*, volume 1350, page 355. Springer Nature.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473.

Isaac Caswell, Ciprian Chelba, and David Grangier. 2019. Tagged back-translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 53–63.

Chris Dyer, Victor Chahuneau, and Noah A Smith. 2013. A simple, fast, and effective reparameterization of ibm model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648.

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, page 489. Association for Computational Linguistics.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic bert sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891.

Jay Gala, Pranjal A Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, et al. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307*.

Ekaterina Garmash and Christof Monz. 2016. Ensemble learning for multi-source neural machine translation. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1409–1418.

Dan Hendrycks and Kevin Gimpel. 2016. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. corr abs/1606.08415 (2016). *arXiv preprint arXiv:1606.08415*.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hérve Jégou, and Tomas Mikolov. 2016. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *EMNLP 2018*, page 66.

Xuan-Phi Nguyen, Shafiq Joty, Wu Kui, and Ai Ti Aw. 2020. Data diversification: a simple strategy for neural machine translation. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pages 10018–10029.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Maja Popović. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *54th Annual Meeting of the Association for Computational Linguistics*, pages 86–96. Association for Computational Linguistics (ACL).

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958.

Tim Van Erven and Peter Harremos. 2014. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820.

Ashish Vaswani. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762*.

Yiren Wang, Lijun Wu, Yingce Xia, Tao Qin, ChengXiang Zhai, and Tie-Yan Liu. 2020. Transductive ensemble learning for neural machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6291–6298.

Daimeng Wei, Zongyao Li, Zhanglin Wu, Zhengzhe Yu, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Minghan Wang, Lizhi Lei, Min Zhang, et al. 2021. Hw-tsc's participation in the wmt 2021 news translation shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 225–231.

Daimeng Wei, Zhiqiang Rao, Zhanglin Wu, Shaojun Li, Yuanchang Luo, Yuhao Xie, Xiaoyu Chen, Hengchao Shang, Zongyao Li, Zhengzhe Yu, et al. 2022. Hw-tsc's submissions to the wmt 2022 general machine translation shared task. In *Proceedings of the Seventh Conference on Machine Translation*, pages 403–410.

Daimeng Wei, Zhanglin Wu, Hengchao Shang, Zongyao Li, Minghan Wang, Jiaxin Guo, Xiaoyu Chen, Zhengzhe Yu, and Hao Yang. 2023. Text style transfer back-translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pages 7944–7959, Toronto, Canada.

Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, Tie-Yan Liu, et al. 2021. R-drop: Regularized dropout for neural networks. *Advances in Neural Information Processing Systems*, 34:10890–10905.

Zhanglin Wu, Daimeng Wei, Zongyao Li, Zhengzhe Yu, Shaojun Li, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Yuhao Xie, Lizhi Lei, et al. 2023. The path to continuous domain adaptation improvements by hw-tsc for the wmt23 biomedical translation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 271–274.

Xuan Zhang, Pamela Shapiro, Gaurav Kumar, Paul McNamee, Marine Carpuat, and Kevin Duh. 2019. Curriculum learning for domain adaptation in neural machine translation. In *Proceedings of NAACL-HLT*, pages 1903–1915.

# NLIP_Lab-IITH Low-Resource MT System for WMT24 Indic MT Shared Task

**Pramit Sahoo    Maharaj Brahma    Maunendra Sankar Desarkar**
Natural Language and Information Processing Lab (NLIP)
Indian Institute of Technology Hyderabad
Hyderabad, India
{ai23mtech14004, cs23resch01004}@iith.ac.in, maunendra@cse.iith.ac.in

## Abstract

In this paper, we describe our system for the WMT 24 shared task of Low-Resource Indic Language Translation. We consider eng ↔ {as, kha, lus, mni} as participating language pairs. In this shared task, we explore finetuning of a pre-trained machine translation model, where the pretraining objective includes alignment of embeddings of tokens from the 22 scheduled Indian languages by a carefully constructed alignment augmentation strategy (Lin et al., 2020). Our primary system[1] is based on language-specific finetuning on this pre-trained model. We achieve chrF2 scores of 50.6, 42.3, 54.9, and 66.3 on the official public test sets for eng→as, eng→kha, eng→lus, eng→mni respectively. We also explore multilingual training with/without language grouping and freezing of encoder and/or embedding layers.

## 1 Introduction

The WMT 2024 Shared Task on "Low-Resource Indic Language Translation" (Pakray et al., 2024) extends the efforts in this direction originally initiated in WMT 2023 (Pal et al., 2023), which garnered significant participation from the global community. Recent advancements in machine translation (MT), particularly through techniques like multilingual training and transfer learning, have expanded the scope of MT systems beyond high-resource languages (Johnson et al., 2017). However, low-resource languages continue to present substantial challenges due to the scarcity of parallel data required for effective training (Siddhant et al., 2020; Wang et al., 2022). The shared task focuses on low-resource Indic languages with limited data from diverse language families: Assamese (as), Mizo (lus), Khasi (kha), and Manipuri (mni). The task aims to improve translation quality for the English⇔Assamese, English⇔Mizo, English⇔Khasi, and English⇔Manipuri given the data provided in the constrained setting.

To address the challenges inherent in translating low-resource languages, participants are encouraged to explore several strategies. First, leveraging monolingual data is essential for enhancing translation quality, especially in the absence of sufficient parallel data. Second, multilingual approaches offer the potential for cross-lingual transfer, where knowledge from high-resource languages can be applied to low-resource pairs (Sen et al., 2019). Third, transfer learning provides a mechanism for adapting pre-trained models from high-resource languages to low-resource settings (Wang et al., 2020). Lastly, innovative techniques tailored to low-resource scenarios, such as data augmentation and language-specific fine-tuning, are crucial for improving performance.

In this paper, we describe our system for the WMT 2024 shared task, focusing on fine-tuning two pre-trained models developed by us: IndicRASP and IndicRASP-Seed[2]. IndicRASP model is pre-trained with the objective of aligning embeddings inspired by alignment augmentation (Lin et al., 2020) on 22 Indic languages. Our primary approach involves language-specific fine-tuning, leveraging multilingual training setups, language grouping, and layer freezing. We set up experiments in both bilingual and multilingual settings. We achieve BLEU scores of 20.1 for English→Assamese, 19.1 for English→Khasi, 30.0 for English→Mizo, and 35.6 for English→Manipuri on the public test set, demonstrating the effectiveness of our approach. Specifically, language-specific fine-tuning yielded significant improvements in translation quality, while multilingual setups provided balanced performance across all language pairs. Language grouping and layer freezing are effective techniques

---

[1] Our code, models, and generated translations are available here: https://github.com/pramitsahoo/WMT2024-LRILT

[2] These pre-trained models were developed for WAT 2024 MultiIndicMT shared task by the authors.

for preserving pre-trained knowledge and mitigating the challenges of multilinguality. Our results highlight the importance of tailored fine-tuning strategies for low-resource languages and show the potential of using alignment-augmented pre-trained models to improve translation quality in low-resource settings.

## 2   Data

In this section, we present the details of the IndicNECorp1.0 dataset provided by the IndicMT shared task[3] organizers.

| Language pair | Script | Dataset | #parallel sents |
|---|---|---|---|
| English-Assamese | Bengali | Training | 50000 |
| | | Validation | 2000 |
| | | Test | 2000 |
| English-Khasi | Latin | Training | 24000 |
| | | Validation | 1000 |
| | | Test | 1000 |
| English-Manipuri | Bengali | Training | 21687 |
| | | Validation | 1000 |
| | | Test | 1000 |
| English-Mizo | Latin | Training | 50000 |
| | | Validation | 1500 |
| | | Test | 2000 |

Table 1: Parallel dataset details. Script refers to the writing script of the Indic language.

### 2.1   Monolingual Data

The official data also includes monolingual data for four languages. The dataset comprises approximately 2.6M sentences for Assamese, 0.1M for Khasi, 2M for Mizo, and 1M for Manipuri.

### 2.2   Parallel Data

The dataset includes four bilingual pairs between English and Indic languages[4]: English (en) - Assamese (as), English (en) - Khasi (kha), English (en) - Mizo (lus), and English (en) - Manipuri (mni). These languages are mainly spoken in the Northeastern part of India. The English-Assamese and English-Mizo training sets contain 50k parallel sentences each, while the English-Khasi and English-Manipuri training sets contain 24k and 21.6k parallel sentences, respectively. Dataset statistics are presented in Table 1.

## 3   Approach

In this section, we briefly describe our approaches. We explore transfer learning, language grouping,

and layer-freezing techniques.

### 3.1   Transfer Learning

We explore transfer learning based on two pre-trained models IndicRASP and IndicRASP-Seed. IndicRASP-Seed is a fine-tuned model of IndicRASP on small and high-quality data. Particularly, the pre-trained model is trained on agreement-based objective (Lin et al., 2020; Yang et al., 2020) for Indic languages. Specifically, words from source sentences are randomly substituted by the semantically equivalent words from other languages. The model is pre-trained in 22 scheduled Indic languages using a subset of the Bharat Parallel Corpus Collection (BPCC) dataset (Gala et al., 2023). Out of these 22 languages, two of the shared task languages, Assamese and Manipuri, are part of the pre-training. Alignment augmentation is performed using bi-lingual dictionaries from MUSE[5] (Conneau et al., 2017) and GATITOS[6].

### 3.2   Language Grouping

We explore the effect of grouping languages based on script similarity in a multilingual setup. Although our primary focus is on bilingual models, for language grouping experiments, we utilize a multilingual approach where languages sharing similar scripts are trained together. This approach is motivated by the idea that joint training with similar languages can improve translation quality due to shared vocabulary and linguistic properties (Jiao et al., 2022; Gala et al., 2023).

- **Group 1** (Bengali script): Assamese and Manipuri

- **Group 2** (Latin script): Khasi and Mizo

### 3.3   Layer Freezing

We explored layer-freezing approaches to see the impact of freezing different layers of the architecture on final translation performance.
**Frozen Encoder:** In this approach, we freeze the encoder components during the fine-tuning process to preserve their pre-trained weights from the parent model while the embedding and decoder components are updated.
**Frozen Embedding + Encoder:** In this setup, we keep the embedding and encoder frozen during

fine-tuning to preserve their pre-trained weights while updating only the parameters of the rest of the layers.

## 4 Experimental Setup

**Settings:** We fine-tune pre-trained checkpoints: IndicRASP and IndicRASP-Seed models on official parallel data using the Adam optimizer (Kingma and Ba, 2014) with $\beta_1$ set to 0.9 and $\beta_2$ set to 0.98. We set the initial warmup learning rate to 1e-07 and the learning rate to 3e-5, with a warmup step of 4000. We train the models with a dropout rate of 0.3 and a label smoothing rate of 0.1. All experiments are conducted on a single NVIDIA A100 GPU. We use a maximum token count of 512 per batch, accumulating gradients over two steps to simulate a larger batch size. The model is trained for up to 1,000,000 updates. We save checkpoints every 2500 updates. We employed a patience of 10 for early stopping.

**Evaluation Metrics:** We use the official dev and test sets of IndicNECorp1.0 for validation and evaluation. We evaluate using BLEU (Papineni et al., 2002), chrF (Popović, 2015), and chrF++ (Popović, 2017) metrics. We use the SacreBLEU toolkit (Post, 2018) to perform our evaluation[7] with a chrF word order of 2. Additionally, as per the evaluation metrics used by the organizers, we report results on TER (Snover et al., 2006), RIBES (Isozaki et al., 2010), and COMET (Rei et al., 2022) for our primary and contrastive submissions.

**Models:** We conducted our experiments in both bilingual and multilingual settings. In the bilingual setup, we fine-tuned the IndicTrans2 Distilled model (Gala et al., 2023), IndicRASP, and IndicRASP-Seed models for both English to Indic and Indic to English directions. The translation models are trained separately for each Indic language. In the multilingual setup, we fine-tuned pre-trained checkpoints of IndicRASP and IndicRASP-Seed for both directions. Inspired by Chiang et al. (2022), we initialized the bilingual model with a fine-tuned multilingual model for both English to Indic and Indic to English.

For experiments with layer freezing, we fine-tune pre-trained checkpoints of IndicTrans2 Distilled and IndicRASP-Seed models. Particularly, we perform experiments by freezing the embed-

---

[7]SacreBLEU signature:
`nrefs:1|case:mixed|eff:no|tok:13a`
`|smooth:exp|version:2.3.1`

dings and encoder and only the encoder component for both English to Indic and Indic to English directions. We conduct all layer-freezing experiments in a bilingual setup. For language grouping experiments, we fine-tune the IndicRASP and IndicRASP-Seed models based on script similarity in a multilingual setup.

## 5 Results and Discussions

In this section, we report our experimental results and describe our primary and contrastive submissions. The results for our primary and contrastive systems are shown in Table 4. Tables 2, 3, and 5 reports the chrF2, BLEU, and chrF++ scores respectively.

①  **English → Indic:** Our primary English to Indic systems are language pair-specific (bilingual models) fine-tuned on pre-trained IndicRASP-Seed, achieving chrF2 scores of 50.6, 42.3, 54.9, and 66.3 for Assamese, Khasi, Mizo, and Manipuri respectively. For the contrastive systems, we consider a bilingual model fine-tuned on a pre-trained IndicRASP checkpoint. The contrastive system achieves chrF2 scores of 49.9, 42.2, 36.5, and 65.8 for Assamese, Khasi, Mizo, and Manipuri, respectively. The detailed primary and contrastive system results are reported in Table 4.

②  **Indic → English:** Our primary Indic-to-English systems for Assamese and Manipuri are bilingual models fine-tuned on the pre-trained IndicRASP-Seed model, each achieving chrF2 scores of 52.8 and 67.9, respectively. Similarly, for Khasi and Mizo, our primary systems are bilingual models fine-tuned on a pre-trained IndicRASP checkpoint, achieving a chrF2 score of 36.1 and 49.4, respectively.

For the contrastive Indic-to-English system, we submit a multilingual system fine-tuned on the pre-trained checkpoint of the IndicRASP model, achieving chrF2 scores of 51.2, 36.0, 46.5, and 65.3 for Assamese, Khasi, Mizo, and Manipuri respectively. Table 4 shows the detailed scores in various metrics.

**Bilingual vs. Multilingual:** We observe IndicRASP-Seed outperforms the IndicRASP model for Assamese and Manipuri. This might be due to the fact that IndicRASP-Seed performs

| Models | English → Indic | | | | Indic → English | | | |
|---|---|---|---|---|---|---|---|---|
| | as | kha | lus | mni | as | kha | lus | mni |
| **BILINGUAL SETUP** | | | | | | | | |
| INDICTRANS2 DISTILLED FT ON BILINGUAL DATA | 49.5 | 24.9 | 29.1 | 60.1 | 50.9 | 21.1 | 22.0 | 61.9 |
| INDICRASP FT ON BILINGUAL DATA | 49.9 | 42.2 | 36.5 | 65.8 | 50.1 | 36.1 | 49.4 | 67.7 |
| INDICRASP-SEED FT ON BILINGUAL DATA | 50.6 | 42.3 | 54.9 | 66.3 | 52.8 | 36.1 | 25.1 | 67.9 |
| **MULTILINGUAL SETUP** | | | | | | | | |
| INDICRASP FT ON MULTILINGUAL DATA | 49.8 | 34.6 | 51.5 | 63.2 | 51.2 | 36.0 | 46.5 | 65.3 |
| INDICRASP-SEED FT ON MULTILINGUAL DATA | 48.7 | 34.6 | 50.2 | 62.2 | 52.2 | 35.3 | 44.3 | 65.1 |
| **MULTILINGUAL MODEL FT ON BILINGUAL DATA** | | | | | | | | |
| INDICRASP MULTILINGUAL MODEL FT ON BILINGUAL DATA | 49.3 | 42.4 | 54.7 | 65.8 | 50.9 | 36.3 | 46.8 | 67.4 |
| **LAYER FREEZING** | | | | | | | | |
| INDICTRANS2 DISTILLED FT WITH FROZEN ENCODER | 47.4 | 24.4 | 28.0 | 57.8 | 48.7 | 19.8 | 18.7 | 58.8 |
| INDICRASP-SEED FT WITH FROZEN ENCODER | 50.4 | 41.3 | 48.6 | 63.4 | 52.6 | 26.4 | 34.2 | 65.3 |
| INDICTRANS2 DISTILLED FT WITH FROZEN EMBEDDING & ENCODER | 46.7 | 23.1 | 9.2 | 15.9 | 48.8 | 20.2 | 19.6 | 58.1 |
| INDICRASP-SEED FT WITH FROZEN EMBEDDING & ENCODER | 50.5 | 41.2 | 45.8 | 62.4 | 52.9 | 25.9 | 29.6 | 64.1 |
| **LANGUAGE GROUPING** | | | | | | | | |
| INDICRASP FT WITH SCRIPT SIMILARITY | 50.2 | 35.0 | 52.1 | 63.3 | 52.6 | 36.4 | 46.5 | 66.0 |
| INDICRASP-SEED MODEL FT WITH SCRIPT SIMILARITY | 50.3 | 34.9 | 53.5 | 63.6 | 53.6 | 36.8 | 47.4 | 66.8 |

Table 2: chrF2 scores on IndicMT WMT24 shared task public test set.

| Models | English → Indic | | | | Indic → English | | | |
|---|---|---|---|---|---|---|---|---|
| | as | kha | lus | mni | as | kha | lus | mni |
| **BILINGUAL SETUP** | | | | | | | | |
| INDICTRANS2 DISTILLED FT ON BILINGUAL DATA | 18.0 | 9.3 | 13.6 | 21.6 | 26.3 | 2.7 | 5.0 | 36.2 |
| INDICRASP FT ON BILINGUAL DATA | 20.5 | 18.9 | 13.1 | 33.9 | 20.0 | 14.4 | 29.1 | 43.6 |
| INDICRASP-SEED FT ON BILINGUAL DATA | 20.1 | 19.1 | 30.0 | 35.6 | 27.4 | 14.1 | 6.0 | 44.1 |
| **MULTILINGUAL SETUP** | | | | | | | | |
| INDICRASP FT ON MULTILINGUAL DATA | 18.7 | 13.5 | 25.8 | 29.0 | 25.8 | 14.1 | 25.4 | 39.3 |
| INDICRASP-SEED FT ON MULTILINGUAL DATA | 17.1 | 13.2 | 24.4 | 27.2 | 26.7 | 14.1 | 23.3 | 38.3 |
| **MULTILINGUAL MODEL FT ON BILINGUAL DATA** | | | | | | | | |
| INDICRASP MULTILINGUAL MODEL FT ON BILINGUAL DATA | 19.1 | 19.0 | 29.7 | 34.7 | 25.8 | 14.8 | 26.1 | 43.5 |
| **LAYER FREEZING** | | | | | | | | |
| INDICTRANS2 DISTILLED FT WITH FROZEN ENCODER | 15.6 | 8.9 | 13.1 | 19.6 | 22.7 | 1.5 | 3.0 | 31.3 |
| INDICRASP-SEED FT WITH FROZEN ENCODER | 19.7 | 18.1 | 22.4 | 29.0 | 26.8 | 5.6 | 15.2 | 40.7 |
| INDICTRANS2 DISTILLED FT WITH FROZEN EMBEDDING & ENCODER | 14.8 | 8.3 | 2.6 | 1.3 | 22.7 | 1.9 | 3.8 | 30.5 |
| INDICRASP-SEED FT WITH FROZEN EMBEDDING & ENCODER | 19.4 | 17.7 | 19.7 | 27.2 | 26.9 | 5.4 | 10.9 | 37.9 |
| **LANGUAGE GROUPING** | | | | | | | | |
| INDICRASP FT WITH SCRIPT SIMILARITY | 19.1 | 13.8 | 26.6 | 28.9 | 26.9 | 14.6 | 25.5 | 39.8 |
| INDICRASP-SEED MODEL FT WITH SCRIPT SIMILARITY | 19.4 | 14.1 | 28.6 | 29.4 | 28.3 | 14.8 | 26.4 | 40.6 |

Table 3: BLEU scores on IndicMT WMT24 shared task public test set.

an additional pre-training on a small, high-quality dataset over IndicRASP. However, when the original pre-training dataset did not contain the languages, like the case of Mizo and Khasi languages here, the comparison shows an opposite trend.

Bilingual models perform better than multilingual models, showing a +4.1 and +7.7 chrF2 score improvement for English to Manipuri and English to Khasi, respectively.

Bilingual models initialized with the weights from multilingual models show improvement over the standalone multilingual models, achieving a +7.8 chrF2 score for English to Khasi. This suggests that initializing bilingual models can be helpful in low-resource settings.

**Language Grouping:** We observe that script-based language grouping shows improvements over a standalone multilingual model with +1.6, +0.3, +3.3, and +1.4 for English to Assamese, Khasi, Mizo, and Manipuri, respectively. It suggests that grouping languages based on script similarity can be effective in addressing the curse of multilinguality.

**Layer Freezing:** We observe that freezing only the encoder yields better chrF2 scores compared to freezing both the embedding and the encoder. However, layer freezing underperforms compared to full parameter fine-tuned bilingual models.

| | BLEU | chrF2 | TER | RIBES | COMET |
|---|---|---|---|---|---|
| **PRIMARY** | | | | | |
| en→as | 20.1 | 50.6 | 66.0 | 0.5543 | 0.8090 |
| en→kha | 19.1 | 42.3 | 63.5 | 0.6470 | 0.6817 |
| en→lus | 30.0 | 54.9 | 50.0 | 0.6764 | 0.7105 |
| en→mni | 35.6 | 66.3 | 50.5 | 0.6995 | 0.7669 |
| as→en | 27.4 | 52.8 | 65.3 | 0.6749 | 0.7854 |
| kha→en | 14.4 | 36.1 | 82.0 | 0.5601 | 0.5773 |
| lus→en | 29.1 | 49.4 | 66.7 | 0.6436 | 0.7004 |
| mni→en | 44.1 | 67.9 | 50.2 | 0.7894 | 0.8162 |
| **CONTRASTIVE** | | | | | |
| en→as | 20.5 | 49.9 | 67.2 | 0.5356 | 0.8043 |
| en→kha | 18.9 | 42.2 | 63.5 | 0.6499 | 0.6791 |
| en→lus | 13.1 | 36.5 | 73.8 | 0.4357 | 0.6462 |
| en→mni | 33.9 | 65.8 | 50.5 | 0.6972 | 0.7672 |
| as→en | 25.8 | 51.2 | 66.8 | 0.6744 | 0.7802 |
| lus→en | 25.4 | 46.5 | 69.0 | 0.6307 | 0.6882 |
| mni→en | 39.3 | 65.3 | 52.4 | 0.7806 | 0.8034 |

Table 4: Submission results on the IndicMT WMT24 public test set.

## 6 Conclusion

In this paper, we describe NLIP Lab's Indic low-resource machine translation systems for the WMT24 shared task. We explore the translation capabilities of the alignment-augmented pre-trained model, IndicRASP and IndicRASP-Seed, to enhance translation quality for low-resource Indic languages. Experimentally, we found that the IndicRASP model performs better than the IndicTrans2 Distilled model. Additionally, we experiment with layer-freezing and language grouping techniques. In the future, we will focus on refining these techniques and utilizing monolingual data to enhance MT performance for low-resource Indic languages.

## Limitations

The pre-trained models use bilingual dictionaries whose domains might differ from the shared task training corpus. Additionally, the considered pre-trained models cover only a limited number of shared task languages. Our submission does not utilize the provided monolingual data, which could further improve model performance through back-translation.

## Acknowledgements

## References

Ting-Rui Chiang, Yi-Pei Chen, Yi-Ting Yeh, and Graham Neubig. 2022. Breaking down multilingual machine translation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2766–2780, Dublin, Ireland. Association for Computational Linguistics.

Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 944–952, Cambridge, MA. Association for Computational Linguistics.

Wenxiang Jiao, Zhaopeng Tu, Jiarui Li, Wenxuan Wang, Jen-tse Huang, and Shuming Shi. 2022. Tencent's multilingual machine translation system for WMT22 large-scale African languages. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1049–1056, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Zehui Lin, Xiao Pan, Mingxuan Wang, Xipeng Qiu, Jiangtao Feng, Hao Zhou, and Lei Li. 2020. Pre-training multilingual neural machine translation by leveraging alignment information. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2649–2663, Online. Association for Computational Linguistics.

Partha Pakray, Santanu Pal, Advaitha Vetagiri, Reddi Mohana Krishna, Arnab Kumar Maji, Sandeep Kumar Dash, Lenin Laitonjam, Lyngdoh Sarah, and Riyanka Manna. 2024. Findings of wmt 2024 shared task on low-resource indic languages

| Models | English → Indic | | | | Indic → English | | | |
|---|---|---|---|---|---|---|---|---|
| | as | kha | lus | mni | as | kha | lus | mni |
| **BILINGUAL SETUP** | | | | | | | | |
| INDICTRANS2 DISTILLED FT ON BILINGUAL data | 45.8 | 25.6 | 30.3 | 55.4 | 49 | 20 | 21 | 59.6 |
| INDICRASP FT ON BILINGUAL data | 46.4 | 41.3 | 35.2 | 61.8 | 46.5 | 35.3 | 48.2 | 65.4 |
| INDICRASP-SEED FT ON BILINGUAL data | 47 | 41.4 | 53.2 | 62.3 | 50.6 | 35.3 | 24 | 65.7 |
| **MULTILINGUAL SETUP** | | | | | | | | |
| INDICRASP FT ON MULTILINGUAL data | 46.2 | 33.4 | 49.8 | 58.9 | 49.1 | 35.2 | 45.4 | 63 |
| INDICRASP-SEED FT ON MULTILINGUAL data | 45.1 | 33.4 | 48.5 | 57.9 | 50.1 | 34.6 | 43.2 | 62.6 |
| **MULTILINGUAL MODEL FT ON BILINGUAL DATA** | | | | | | | | |
| INDICRASP MULTILINGUAL MODEL FT ON BILINGUAL data | 45.7 | 41.5 | 53.1 | 61.9 | 48.8 | 35.5 | 45.7 | 65.2 |
| **LAYER FREEZING** | | | | | | | | |
| INDICTRANS2 DISTILLED FT WITH FROZEN ENCODER | 43.7 | 25.1 | 29.3 | 53 | 46.7 | 18.5 | 17.6 | 59.8 |
| INDICRASP-SEED FT WITH FROZEN ENCODER | 46.8 | 40.3 | 46.9 | 59.1 | 50.4 | 25.3 | 33.1 | 63 |
| INDICTRANS2 DISTILLED FT WITH FROZEN ENCODER & EMBEDDINGS | 43 | 24 | 11.3 | 13.1 | 46.8 | 18.9 | 18.5 | 55.6 |
| INDICRASP-SEED FT WITH FROZEN ENCODER & EMBEDDINGS | 46.8 | 40.2 | 44.1 | 58 | 50.6 | 24.9 | 28.6 | 61.7 |
| **LANGUAGE GROUPING** | | | | | | | | |
| INDICRASP FT WITH SCRIPT SIMILARITY | 46.6 | 33.8 | 50.4 | 59 | 50.4 | 35.6 | 45.4 | 63.6 |
| INDICRASP-SEED MODEL FT WITH SCRIPT SIMILARITY | 46.7 | 33.7 | 51.8 | 59.4 | 51.5 | 36 | 46.3 | 64.4 |

Table 5: chrF2++ scores on IndicMT WMT24 shared task public test set.

translation. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the WMT 2023 shared task on low-resource Indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694, Singapore. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins.

2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Sukanta Sen, Kamal Kumar Gupta, Asif Ekbal, and Pushpak Bhattacharyya. 2019. Multilingual unsupervised NMT using shared encoder and language-specific decoders. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3083–3089, Florence, Italy. Association for Computational Linguistics.

Aditya Siddhant, Ankur Bapna, Yuan Cao, Orhan Firat, Mia Chen, Sneha Kudugunta, Naveen Arivazhagan, and Yonghui Wu. 2020. Leveraging monolingual data with self-supervision for multilingual neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2827–2835, Online. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Wenxuan Wang, Wenxiang Jiao, Shuo Wang, Zhaopeng Tu, and Michael R. Lyu. 2022. Understanding and mitigating the uncertainty in zero-shot translation. *ArXiv*, abs/2205.10068.

Xinyi Wang, Yulia Tsvetkov, and Graham Neubig. 2020. Balancing training for multilingual neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*,

pages 8526–8537, Online. Association for Computational Linguistics.

Zhen Yang, Bojie Hu, Ambyera Han, Shen Huang, and Qi Ju. 2020. CSP:code-switching pre-training for neural machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2624–2636, Online. Association for Computational Linguistics.

# Yes-MT's Submission to the Low-Resource Indic Language Translation Shared Task in WMT 2024

**Yash Bhaskar[1], Parameswari Krishnamurthy[2]**
IIIT Hyderabad
yash.bhaskar@research.iiit.ac.in, param.krishna@iiit.ac.in

## Abstract

This paper presents the systems submitted by the Yes-MT team for the Low-Resource Indic Language Translation Shared Task at WMT 2024 (Pakray et al., 2024), focusing on translating between English and the Assamese, Mizo, Khasi, and Manipuri languages. The experiments explored various approaches, including fine-tuning pre-trained models like mT5 (Xue et al., 2020) and IndicBart (Dabre et al., 2021) in both Multilingual and Monolingual settings, LoRA (Hu et al., 2021) finetune IndicTrans2 (Gala et al., 2023), zero-shot and few-shot prompting (Brown, 2020) with large language models (LLMs) like Llama 3 (Dubey et al., 2024) and Mixtral 8x7b (Jiang et al., 2024), LoRA Supervised Fine Tuning (Mecklenburg et al., 2024) Llama 3, and training Transformers (Vaswani, 2017) from scratch. The results were evaluated on the WMT23 Low-Resource Indic Language Translation Shared Task's test data using SacreBLEU (Post, 2018) and CHRF (Popović, 2015) highlighting the challenges of low-resource translation and show the potential of LLMs for these tasks, particularly with fine-tuning.

## 1 Introduction

Developing robust machine translation systems for India's diverse languages is crucial given the country's growing economic importance and the increasing availability of digital content. However, a significant challenge in developing effective translation tools arises from the limited availability of data for many Indian languages, particularly those spoken in the northeastern regions. This paper describes the Yes-MT team's efforts to address this challenge by participating in the WMT 2024 Low-Resource Indic Language Translation Shared Task, focusing on English to Assamese, Mizo, Khasi, and Manipuri translation. We explored techniques like fine-tuning pre-trained models (mT5, IndicBart) and utilizing large lan-

guage models (LLMs) like Llama 3 and Mixtral for zero-shot and few-shot learning. Furthermore, we explored using the LoRA technique to fine-tune the IndicTrans2 model, and we also trained Transformer models from scratch. Our findings provide valuable insights into the strengths and weaknesses of different approaches, highlighting the potential of LLMs and fine-tuning techniques in overcoming the limitations of data scarcity.

## 2 Dataset

The dataset used in this study consists of parallel bilingual data provided by the WMT 2024 Low-Resource Indic Language Translation Shared Task organizers (Pal et al., 2023) & (Pakray et al., 2024). The training, validation, and test splits for each language pair are detailed in Table 1.

| Language Pair | Train | Val | Test |
|---|---|---|---|
| Assamese (en-as) | 50,000 | 2,000 | 2,000 |
| Mizo (en-lus) | 50,000 | 1,500 | 2,000 |
| Khasi (en-kha) | 24,000 | 1,000 | 1,000 |
| Manipuri (en-mni) | 21,000 | 1,000 | 1,000 |

Table 1: Number of Sentences in Train, Validation, and Test Sets

In addition to the bilingual data, we also had access to a significant amount of Monolingual data for each of the target languages, which included 2.60 million sentences in Assamese, 1.90 million sentences in Mizo, 0.18 million sentences in Khasi, and 2.10 million sentences in Manipuri. However, for the scope of this work, we focused exclusively on utilizing the provided bilingual data for training and evaluation, aiming to explore the capabilities of the models under truly low-resource conditions.

Limiting our study to the provided bilingual data allowed us to maintain a consistent and controlled experimental environment, ensuring the results reflected the performance of our approaches under

the typical constraints of low-resource language translation tasks. In the future, we may explore incorporating the available monolingual data, such as through back-translation, to further improve translation quality.

## 3 Experiments

This section details the experimental setup used for the various models and training strategies employed in our submission.

### 3.1 Primary Submission

Our primary submission involved training a Transformer model from scratch using the Fairseq framework (Ott et al., 2019). This model was trained for Multilingual translation, handling all four language directions (English to Assamese, Manipuri, Mizo, and Khasi) simultaneously. We utilized BPE tokenizer (Araabi et al., 2022) and Transformer architecture. The architectural details are shown in Table 2.

| Parameter | Value |
|---|---|
| Embedding Dimension | 512 |
| FFN Dimension | 1024 |
| Attention Heads | 4 |
| Encoder Layers | 6 |
| Decoder Layers | 6 |

Table 2: Transformer Architecture Details

### 3.2 Contrastive Submission

The contrastive submission explored fine-tuning pre-trained models in two settings: language-specific and Multilingual.

#### 3.2.1 Multilingual Fine-tuning:

Both mT5 and IndicBart were fine-tuned in a Multilingual setting, where a single model was trained to handle all four language directions. To enable the models to distinguish between the target languages, we added language-specific tokens to their existing vocabularies, as suggested by previous work (Johnson et al., 2017). The language-specific tokens used are shown in Table 3. A single model was trained for one-to-many translation across all four language directions for each of the indicBart, mT5-small, and IndicTrans2 systems. The results are in Table 4. IndicBart and mT5-small were fine-tuned using Full Fine-Tuning (FFT), while IndicTrans2 was fine-tuned employing the LoRA (Low-Rank Adaptation) technique (Hu et al., 2021).

| Language | Token |
|---|---|
| Assamese (asm) | '<asm_Beng>' |
| Manipuri (mni) | '<mni_Beng>' |
| Khasi (kha) | '<kha_Latn>' |
| Mizo (lus) | '<lus_Latn>' |

Table 3: Language-Specific Tokens

#### 3.2.2 Monolingual Fine-tuning:

We also trained separate models for each language pair, as these focused on a single translation direction and did not require language-specific tokens.

For each language direction, we trained four distinct models using mT5-Small and IndicBart with Full Fine-Tuning (FFT). The results are in Table 4.

### 3.3 Experiments with LLMs

Additionally, we explored the use of the Llama3 model in conjunction with the LoRA (Low-Rank Adaptation) technique.

**Zero-Shot and Few-Shot Translation Evaluation** We tested Zero Shot Translation capabilities of Llama 3-8B-8192, Llama 3-70B-8192, mixtral-8x7B-32768, Llama3-8B-instruct and Llama3.1-8B-instruct. We also tested the few-shot translation capabilities of Llama3.1-8B-instruct with 3-shot, 5-shot, and 10-shot prompting.

**Supervised Fine-Tuning with LoRA** We fine-tuned a 4-bit quantized (Liu et al., 2023) Llama3 model using the LoRA technique with Supervised Fine-Tuning (SFT), employing the LlamaFactory framework (Zheng et al., 2024). We used a prompt-based approach for translation, providing the model with a system prompt and a prompt template specifying the source and target languages.

The following template was used for fine-tuning the Large Language Models (LLMs):

```
System Prompt : You are a helpful assistant.
Prompt Template : Translate the following
English sentence to {target_language} in
{target_script} Script:\n{input_sent}
```

## 4 Results

### 4.1 Multilingual vs. Monolingual Performance

One key finding from our experiments was the performance comparison between the Multilingual and Monolingual training approaches for the mT5 and IndicBart models. As shown in Table 4, the

| Model | Training Type | en-as | en-kha | en-mz | en-mni |
|---|---|---|---|---|---|
| Transformers | Multilingual | 16.06 | 19.67 | 5.49 | 20.60 |
| IndicBart | Monolingual | 6.4 | 11.2 | 25.1 | 8.8 |
| | Multilingual | 6.5 | 11.4 | 25.3 | 9.1 |
| mT5-small | Monolingual | 14.3 | 12.9 | 31.4 | 19.2 |
| | Multilingual | 15.6 | 13.6 | 32.3 | 23.9 |
| IndicTrans2-2B | ZeroShot | 49.2 | - | - | 44.9 |
| IndicTrans2-200M | ZeroShot | 49.5 | - | - | 45.3 |
| | Multilingual | 47.27 | - | - | 49.12 |

Table 4: ChrF Scores for Monolingual : Models fine-tuned for one-to-one language translation
Multilingual : Models fine-tuned for one-to-many language translation

Multilingual versions of both mT5 and IndicBart consistently outperformed their Monolingual counterparts across the translation tasks.

- For mT5, the Multilingual model outperformed the Monolingual model across all language pairs, with ChrF score improvements ranging from 1.3 to 4.7 points. This suggests that mT5 benefits from the shared linguistic knowledge across different languages in a Multilingual setting, which enhances its ability to generalize to low-resource languages.

- Likewise, IndicBart demonstrated a slight performance boost in the Multilingual setting compared to the Monolingual models, suggesting that the Multilingual training approach provided a benefit.

The better performance of the Multilingual models is likely due to the shared linguistic knowledge they gained during training, which may have provided a richer context and improved their ability to generalize. This indicates that leveraging Multilingual data, even in limited-resource scenarios, can be a more effective approach than focusing on Monolingual training.

## 4.2 Expected Structured Output

A challenge observed during the experiments was the generation of structured output. Ideally, the output should directly provide the translated sentence without additional, unnecessary text. However, we noticed that the LLM models sometimes wrapped the translation in extraneous text, such as "The translation of the given sentence is: Translation", followed by further analysis and explaination making it difficult to extract the translation. This adds noise to the output and complicates the process of extracting the actual translation.



Figure 1: Inconsistent Output Format with Few Shot Prompting

We analyzed the percentage of outputs that were wrapped with unnecessary text across different settings:

This issue of unnecessary text in the output was more common in the zero-shot setting, where 66.% of the outputs included additional text. As the number of shots increased, the percentage of such outputs decreased significantly to 0.18% in 10 Shot Prompting, indicating that few-shot prompting can help guide the LLM to produce more structured and concise translations.

To improve the usability of LLM-based machine translation systems, it's crucial to fine-tune the models or design prompts that consistently yield clean and structured outputs, particularly in low-resource settings where post-processing resources might be limited.

## 5 WMT 2024 Results

The performance of our models on the WMT 2024 Low-Resource Indic Language Translation Shared Task dataset is summarized in the following table, focusing on the ChrF (Popović, 2015) metrics:

For the primary submissions, we utilized Transformers trained from scratch without additional data. As indicated by the scores, the primary

| Model | Inference | en-as | en-kha | en-mz | en-mni |
|---|---|---|---|---|---|
| Llama3-8B-8192 | Zero Shot | 18.56 | 14.92 | 15.57 | 13.45 |
| Llama3-70B-8192 | Zero Shot | 27.54 | 18.57 | 20.62 | 15.53 |
| mixtral-8x7B-32768 | Zero Shot | 6.79 | 15.45 | 16.57 | 2.65 |
| Llama3-8B-instruct | Zero Shot | 26.13 | 8.38 | 18.06 | 15.29 |
| | 1 Epoch | 29.82 | 33.19 | 32.72 | 37.85 |
| | 2 Epoch | 31.68 | 35.26 | 37.73 | 44.51 |
| Llama3.1-8B-instruct | Zero Shot | 22.93 | 12.03 | 15.23 | 14.47 |
| | 3 Shot | 23.26 | 13.66 | 18.89 | 15.30 |
| | 5 Shot | 23.48 | 15.11 | 18.77 | 15.29 |
| | 10 Shot | 23.89 | 16.03 | 19.39 | 15.43 |

Table 5: ChrF Scores for Various Models, Shot Types, and Language Pairs

| Language Pair | Submission Type | ChrF |
|---|---|---|
| Eng-Asm | primary | 0.1123 |
| | contrastive | 0.6518 |
| Eng-Mni | primary | 0.1102 |
| | contrastive | 0.4438 |
| Eng-Lus | primary | 0.1282 |
| | contrastive | 0.4151 |
| Eng-Kha | primary | 0.1139 |
| | contrastive | 0.3541 |

Table 6: ChrF Scores for WMT 2024 Shared Task

systems struggled significantly, yielding very low ChrF values across all language pairs.

In contrast, the models fine-tuned for the contrastive submissions demonstrated noticeable improvements. For Assamese and Manipuri, we fine-tuned IndicTrans2, achieving the highest ChrF scores in these language pairs. For Mizo and Khasi, we fine-tuned Llama3, which also resulted in enhanced performance compared to the primary systems. These findings highlight the effectiveness of fine-tuning pre-trained models, even in low-resource settings.

## 6 Potential Test Set Bias

One of the noteworthy observations in this year (2024) WMT 2024 results is the significant difference in the performance of the primary Transformers trained from scratch when evaluated on this year's (2024) test set compared to last year's (2023) test set. Specifically, we observed that the models performed better on last year's test set despite using the same training data.

This discrepancy could be indicative of a translation bias present in last year's dataset, which might have inadvertently favored the models trained on

that data. The primary systems, having been trained exclusively on the previous year's data, may have overfitted to patterns specific to that dataset, leading to better performance on the older test set but struggling on the newer one.

This implies that the primary models may have difficulty generalizing to entirely new data distributions, an important factor to consider in low-resource settings where the training data is limited and may not be representative of future data. It also underscores the importance of using diverse and varied datasets during training to help mitigate such biases and improve the overall robustness of the models.

## 7 Conclusion

This paper presented the systems and results of the Yes-MT team's participation in the WMT 2024 Low-Resource Indic Language Translation Shared Task. The experiments highlighted the potential of LLMs, especially when fine-tuned with techniques such as LoRA, in enhancing translation quality even under low-resource conditions. The contrastive submissions, which utilized fine-tuned LLMs, demonstrated significant improvements over the primary submissions that relied on training Transformers from scratch.

Our findings suggest that while training models from scratch can be challenging in low-resource settings due to data scarcity and generalization issues, fine-tuning pre-trained models can effectively bridge the gap, leveraging shared knowledge across languages to achieve better translation performance.

Future work could explore integrating monolingual data through back-translation or other data augmentation techniques, as well as further refin-

ing prompt engineering strategies to improve the structure and clarity of LLM outputs. Additionally, focusing on addressing potential biases in test data to help create more reliable translation systems.

## Acknowledgments

## References

Ali Araabi, Christof Monz, and Vlad Niculae. 2022. How effective is byte pair encoding for out-of-vocabulary words in neural machine translation? *arXiv preprint arXiv:2208.05225.*

Tom B Brown. 2020. Language models are few-shot learners. *arXiv preprint ArXiv:2005.14165.*

Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh M Khapra, and Pratyush Kumar. 2021. Indicbart: A pre-trained model for indic natural language generation. *arXiv preprint arXiv:2109.02903.*

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783.*

Jay Gala, Pranjal A Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, et al. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307.*

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685.*

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088.*

Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.

Shih-yang Liu, Zechun Liu, Xijie Huang, Pingcheng Dong, and Kwang-Ting Cheng. 2023. Llm-fp4: 4-bit floating-point quantized transformers. *arXiv preprint arXiv:2310.16836.*

Nick Mecklenburg, Yiyou Lin, Xiaoxiao Li, Daniel Holstein, Leonardo Nunes, Sara Malvar, Bruno Silva, Ranveer Chandra, Vijay Aski, Pavan Kumar Reddy Yannam, et al. 2024. Injecting new knowledge into large language models via supervised fine-tuning. *arXiv preprint arXiv:2404.00213.*

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Partha Pakray, Santanu Pal, Advaitha Vetagiri, Reddi Mohana Krishna, Arnab Kumar Maji, Sandeep Kumar Dash, Lenin Laitonjam, Lyngdoh Sarah, and Riyanka Manna. 2024. Findings of wmt 2024 shared task on low-resource indic languages translation. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the wmt 2023 shared task on low-resource indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694.

Maja Popović. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.

Matt Post. 2018. A call for clarity in reporting bleu scores. *arXiv preprint arXiv:1804.08771.*

Ashish Vaswani. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762.*

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2020. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934.*

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyan Luo. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. *arXiv preprint arXiv:2403.13372.*

# System Description of BV-SLP for Sindhi-English Machine Translation in MultiIndic22MT 2024 Shared Task

**Nisheeth Joshi[1#], Pragya Katyayan[2*], Palak Arora[1ŧ], Bharti Nathani[4**]**

[1]Speech and Language Processing Lab, Banasthali Vidyapith, Rajasthan, India
[2]School of Computer Science, University of Petroleum and Energy Studies, Uttrakhand, India
[#]nisheeth.joshi@rediffmail.com, [*]pragya.katyayan21@gmail.com, [ŧ]palak.arora.pa55@gmail.com,
[**]nbharti@banasthali.in

## Abstract

This paper presents our machine translation system that was developed for the WAT2024 MultiIndic MT shared task. We built our system for the Sindhi-English language pair. We developed two MT systems. The first system was our baseline system where Sindhi was translated into English. In the second system, we used Hindi as a pivot for the translation of text. In both the cases, we had identified the name entities and translated them into English as a preprocessing step. Once this was done, the standard NMT process was followed to train and generate MT outputs for the task. The systems were tested on the hidden dataset of the shared task

## 1 Introduction

This paper presents the system description of our neural machine translation system developed for the MultiIndic shared task organized at WMT 2024. We collected around two lac English-Hindi parallel corpus from Press Information Bureau's website[1] which had collection of news articles in English as well as in Hindi and then translated it into Sindhi (in Devanagari script). Thus, two NMT systems were trained on using this corpus. The first system was the baseline system which was trained using the Sindhi-English language pair. The second system had two NMT systems, Sindhi-Hindi and Hindi-English. This system used Hindi as the pivot language for translation.

## 2 System Overview

### 2.1 Preprocessing

Here, we did tokenization of text and also performed spelling correction. Then named entities from Sindhi text were extracted using the Bi-LSTM Sindhi POS tagger that was developed in-house (Nathani et al. 2023). The identified named entities were then classified into MUC-6 category (Grishman et al. 1996) through a rule-based approach. These tagged named entities were searched in a knowledge base which had translations of Sindhi/Hindi Organization and Location named entities in English. We extracted the named entities from the Sindhi/Hindi corpus using a rule-based NER system. Once the named entities were extracted, they were searched in a knowledge base that had translations of these named entities in English (Organization and Location names). If they were found then the same were replaced in the Sindhi/Hindi Corpus. In cases where the named entity translations were not present in the knowledge base, then they were transliterated and were replaced in Sindhi/Hindi corpus. This became our Named Entity Translation module which identified the named entities and accordingly translated/transliterated them into English (Sharma et al. 2023; Joshi & Katyayan 2023). The work of this module is shown in Figure 1.

### 2.2 Byte Pair Encoding

Here the source and the target corpus were divided into smaller units known as subwords. This task was performed to convert the words into smaller basic units which helped neural MT models in better handling of out of vocabulary (OOV) words.

### 2.3 Training of the Model

In training both systems we applied the same steps. For system 1 which was the baseline system, we had only one named entity translation module (Sindhi-English) while for the system 2 the named entity translation module performed Hindi-English translation/transliteration. The process followed

---

[1] https://pib.gov.in/

was; the POS tagging of Sindhi sentence was performed and NER was performed using a rule-based module. The identified named entities were translated as explained in the previous section. This produced an augmented corpus-based source sentence. For example, let us consider a Sindhi sentence, "निशीथ जोशी नई दिल्ली जे इंदिरा गांधी अंतर्राष्ट्रीय हवाई अड्डे खां जयपुर जो सफर करे रहियो आहे।" Here "निशीथ जोशी (Person)", "नई दिल्ली (location)", "जयपुर (location)", and "इंदिरा गांधी अंतर्राष्ट्रीय हवाई अड्डे (organization)" are named entities. Among these since "निशीथ जोशी" and "जयपुर" were not available in the knowledge base, so they were transliterated to "Nisheeth Joshi" and "Jaipur" respectively. The rest of the named entities had their categories in the knowledge base; thus, they were looked up in a sequential manner. "नई दिल्ली" was not found and was transliterated to "New Delhi", similarly "इंदिरा गांधी अंतर्राष्ट्रीय हवाई अड्डे" was translated to Indra Gandhi International Airport". The entire training corpus was augmented using this methodology. Figure 2 shows the working of the entire system.

The hyperparameters used in training both the systems are shown in table 1.

| Parameter | Value |
|---|---|
| No. of Encoding Layers | 6 |
| No. of Decoding Layers | 6 |
| **Early Stopping** | |
| metric | bleu |
| min_improvement | 0.2 |
| steps | 6 |
| Optimizer | Adam |
| beta_1 | 0.8 |
| beta_2 | 0.998 |
| learning_rate | 1.0 |
| droupout | 0.25 |
| **Regularization** | |
| type | l1_l2 |
| scale | 1e-4 |
| Minimum_learning_rate | 0.00001 |
| Max_steps | 1000000 |

Tabel 1: Hyperparameters Used in Training NMT Models

## 3   Evaluation

We participated in the shared task using the hidden corpus and submitted the outputs for both the systems viz baseline and pivot MT systems. The results of the same are shown in table 2.

| System | BLEU | chrF | chrF++ |
|---|---|---|---|
| System 1 | 19.4 | 44.6 | 43 |
| System 2 | 20 | 44.7 | 43.2 |

Tabel 2: Evaluation Results

The baseline system which translated Sindhi text into English had a BLEU score (Papineni et al. 2002) of 19.4, chrF score (Popović 2015) of 44.6 and chrF++ score (Popović 2017) of 43. From a human annotators perspective, this system produced fluent translations but in some cases lacked the desired quality. The second system which used Hindi as a pivot language (where Sindhi was translated into Hindi and then this Hindi translation was translated into English) produced slightly better results. Its BLEU score was 20, chrF score was 44.7 and chrF++ score was 43.2. This system generated translation which had improved adequacy and fluency scores.

## Acknowledgement

## References

Grishman, R., & Sundheim, B. M. 1996. *Design of the MUC-6 evaluation. In TIPSTER TEXT PROGRAM PHASE II*: Proceedings of a Workshop held at Vienna, Virginia, May 6-8, 1996 (pp. 413-422).

Joshi, N., & Katyayan, P. 2023. Improving English-Bharti Braille Machine Translation Through Proper Name Entity Translation. In ICIDSSD 2022: Proceedings of the 3rd International Conference on ICT for Digital, Smart, and Sustainable Development, ICIDSSD 2022, 24-25 March 2022,

New Delhi, India (p. 168). European Alliance for Innovation.

Nathani, B., Arora, P., Joshi, N., Katyayan, P., Rathore, S. S., & Dadlani, C. P. 2023. *Sindhi POS Tagger Using LSTM and Pre-Trained Word Embeddings*. In XVIII International Conference on Data Science and Intelligent Analysis of Information (pp. 37-45). Cham: Springer Nature Switzerland.

Papineni, K., Roukos, S., Ward, T., & Zhu, W. J. 2002. *Bleu: a method for automatic evaluation of machine translation*. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics (pp. 311-318).

Popović, M. 2015. *chrF: character n-gram F-score for automatic MT evaluation*. In Proceedings of the tenth workshop on statistical machine translation (pp. 392-395).

Popović, M. 2017. *chrF++: words helping character n-grams*. In Proceedings of the second conference on machine translation (pp. 612-618).

Sharma, R., Katyayan, P., & Joshi, N. 2023. *Improving the quality of neural machine translation through proper translation of name entities*. In 2023 6th International Conference on Information Systems and Computer Networks (ISCON) (pp. 1-4). IEEE.

Figure 1: Named Entity Translation Module

गुरु अर्जुनदेव, तेग़ बहादुरु, महात्मा गांधी, सरदारु भगतु सिंह, हेमू कालाणी, भगतु कंवरु ऐं हासारामु पमिनाणी शहीदनि जे शुमार में आहिनि।

**Tokenization**

**Sentence Normalization**

**POS Tagging**

**Name Entity Recognition**

**NE Translation**
Guru Arjundev, Teg Bahadur, Mahatma Gandhi, Sadguru, Bhagat Singh, Hemu Kalani, Bhattu Kanwar, Hasaram Paminani

**Intermediate Output**
Guru Arjundev, Teg Bahadur, Mahatma Gandhi, Sadguru, Bhagat Singh, Hemu Kalani, Bhattu Kanwar ऐं Hasaram Paminani शहीदनि जे शुमार में आहिनि।

**NMT Framework**

**Post-processing**

Guru Arjundev, Teg Bahadur, Mahatma Gandhi, Sadguru, Bhagat Singh, Hemu Kalani, Bhattu Kanwar and Hasaram Paminani are among the martyrs who have died for the cause of freedom.

Figure 2: BV-SLP NMT Approach

# WMT24 System Description for the MultiIndic22MT Shared Task on Manipuri Language

**Ningthoujam Justwant Singh**, **Kshetrimayum Boynao Singh**,

**Ningthoujam Avichandra Singh**, **Sanjita Phijam**, and **Thoudam Doren Singh**

Centre for Natural Language Processing (CNLP) & Dept. of CSE, NIT Silchar, India

{njustwant92,boynfrancis,avichandra0420,phijamsan.jk,thoudam.doren}@gmail.com

## Abstract

This paper presents a Transformer-based Neural Machine Translation (NMT) system developed by the Centre for Natural Language Processing and the Department of Computer Science and Engineering at the National Institute of Technology Silchar, India (NITS-CNLP) for the MultiIndic22MT 2024 Shared Task. The system focused on the English-Manipuri language pair for the WMT24 shared task. The proposed WMT system shows a BLEU score of 6.4, a chrF score of 28.6, and a chrF++ score of 26.6 on the public test set Indic-Conv dataset. Further, in the public test set Indic-Gen dataset, it achieved a BLEU score of 8.1, a chrF score of 32.1, and a chrF++ score of 29.4 on the English-to-Manipuri translation.

## 1 Introduction

The Centre for Natural Language Processing and the Department of Computer Science and Engineering at the National Institute of Technology Silchar, India (NITS-CNLP) participated in The MultiIndic22MT 2024 Shared Task (Dabre and Kunchukuttan, 2024) for English-Manipuri language pair in the WMT2024 shared task. The shared task involves developing Machine Translation (MT) for English and 22 Indic Languages (Assamese, Bengali, Bodo, Dogri, Konkani, Gujarati, Hindi, Kannada, Kashmiri (Arabic script), Maithili, Malayalam, Marathi, Manipuri (Meitei script), Nepali, Oriya, Punjabi, Sanskrit, Santali, Sindhi (Devanagari script), Tamil, Telugu, Urdu).

In recent years, there has been growing interest in developing effective machine translation systems for Manipuri (Singh et al., 2023a) (Singh and Singh, 2020) (Singh and Singh, 2022b) (Singh et al., 2023b), which is a language with a complex linguistic structure (Singh and Bandyopadhyay, 2010) and limited bitext. Various approaches have been explored to create models that can accurately translate between Manipuri and other languages(Singh and

Singh, 2022a). These efforts include the development of translation models that handle different scripts, such as Bengali and Meitei Mayek, and the integration of linguistic (Singh and Bandyopadhyay, 2005) (Singh and Bandyopadhyay, 2010) and that are essential for producing high-quality translations.

### 1.1 Brief Description of Manipuri language

The Manipuri language can be written in: Bengali and Meitei Mayek. It is one of the 22 official languages of India included in the $8^{th}$ schedule of the Indian constitution. Historically, computational linguistics research and translation efforts for Manipuri have predominantly focused on the Bengali script, due to its extensive availability of digital resources.

Most English-to-Manipuri translation models and linguistic resources have been developed using the Bengali script. Numerous projects have created bitext and bilingual dictionaries in this script, significantly advancing machine translation for Manipuri.

In contrast, the Meitei Mayek script, which holds cultural and historical significance for the Manipuri people, has not received similar attention. Although recent years have seen a revival of the Meitei Mayek script, highlighting the need for computational resources and tools to support its use in modern digital contexts, it still faces challenges due to the limited availability of textual data and digital resources.

Efforts to address this gap include digitizing ancient manuscripts and developing new textual resources in Meitei Mayek.

## 2 Our Approaches

### 2.1 Dataset and Preprocessing

The training dataset (Gala et al., 2023) provided by the WMT Shared Task 2024 consists of 42,740 bitext. After incorporating additional data from the

| Language | Sentence | Word |
|---|---|---|
| **English-Training** | 63506 | 1093014 |
| **Manipuri-Training** | 63506 | 894411 |
| **English-Validation** | 997 | 30772 |
| **Manipuri-Validation** | 997 | 31799 |
| **English-Testing**conv | 1502 | 14849 |
| **Manipuri-Testing** conv | 1502 | 12621 |
| **English-Testing** gen | 1023 | 25347 |
| **Manipuri-Testing** gen | 1023 | 23421 |

Table 1: This table presents the BLEU, chrF, and chrF++ scores for the English-to-Manipuri machine translation system.

Ministry of Electronics and Information Technology (MeitY), we ensure that the dataset is properly aligned to each other to confirm that it consists of bitext, along with removing duplicates and noise. As a result, we obtain a clean training dataset of 63,506 bitext. The validation dataset, also provided by the WMT Shared Task 2024, contains 997 bitext. For testing, we use the test set from the WMT Shared Task 2024, which includes the Indic-Conv and Indic-Gen datasets, comprising 1,502 and 1,023 bitext, respectively.

## 2.2 Hyperparameter

### 2.2.1 Sentencepiece Model

We train a model ($MT_{sp}$) system based on a basic Transformer architecture (Vaswani et al., 2017), utilizing the OpenNMT toolkit (Klein et al., 2017)[1]. In this model, we employ the SentencePiece (Kudo, 2018)[2] tokenization technique with a vocabulary size of 8,000 for both English and Manipuri. The model consists of 6 encoder and 6 decoder layers, each with 8 attention heads. The $MT_{sp}$ system is trained for 200,000 steps, with validation conducted every 5,000 steps, and model checkpoints saved at 5,000-step intervals.

It utilizes a bucket size of 262,144 and a batch size of 2048, along with 8,000 warmup steps. Optimization is performed using the Adam optimizer (Kingma and Ba, 2014). The ($MT_{sp}$) is trained with a feed-forward layer size of 2048, a hidden size of 512, and a label smoothing of 0.1.

### 2.2.2 Proposed Subword Model

Our proposed model ($WMT24_{proposed}$) is also a transformer model trained using the OpenNMT toolkit. For tokenization, we employ the Byte Pair

Encoding (BPE) method (Sennrich et al., 2016)[3] with the same vocab size 8000 for English and Manipuri. The proposed model shares the same hyperparameters as the ($MT_{sp}$), including training for 200,000 steps, with validation every 5,000 steps, and model checkpoints saved at 5,000-step intervals. It also uses the same bucket size of 262,144 and a batch size of 2048.

Both the ($MT_{sp}$) and ($WMT24_{proposed}$) models are configured with 8 attention heads, 6 encoder layers, 6 decoder layers, and a learning rate of 2, along with an attention dropout rate of 0.1. Optimization is performed using the Adam optimizer (Kingma and Ba, 2014), and the models share identical hyperparameters, including a feed-forward layer size of 2048, a hidden size of 512, and label smoothing of 0.1.

We train both the $WMT24_{proposed}$ and $MT_{sp}$ models using the complete set of 63,506 sentence pairs, which includes data from both the WMT Shared Task data and additional data provided by MeitY. We utilize the same validation sentences, and the testing data remains unchanged.

The performance of each model is evaluated using BLEU (Papineni et al., 2002), chrF (Popović, 2015), and chrF++ (Popović, 2017) metrics, utilizing the sacreBLEU tool (Post, 2018)[4] for score evaluation.

## 3 Results and Discussion

In this section, we discuss the experimental results and performance of the models. The reported BLEU, chrF, and chrF++ scores are calculated based on the de-tokenized text. The scores of the systems are given in Table 2. The English-to-Manipuri translation $WMT24_{proposed}$ model achieves a BLEU score of **6.4**, a chrF score of 28.6, and a chrF++ score of 26.6 on the Indic-Conv dataset. In contrast, the $MT_{sp}$ achieves a BLEU score of 5.1, a chrF score of 30.9, and a chrF++ score of 27.1 on the same dataset.

For the Indic-Gen dataset, the $WMT24_{proposed}$ achieves a BLEU score of **8.1**, a chrF score of 32.1, and a chrF++ score of 29.4, while the $MT_{sp}$ achieves a BLEU score of 6.8, a chrF score of 32.8, and a chrF++ score of 28.7. These results highlight the superior performance of the $WMT24_{proposed}$ compared to the $MT_{sp}$ model across all evaluation metrics.

---

[1]https://github.com/OpenNMT/OpenNMT
[2]https://github.com/google/sentencepiece

[3]https://github.com/rsennrich/subword-nmt
[4]https://github.com/mjpost/sacrebleu

| MT systems | Test Set | BLEU | chrF | chrF++ |
|---|---|---|---|---|
| WMT24 proposed | conv | **6.4** | **28.6** | **26.6** |
| MT$_{sp}$ | conv | 5.1 | 30.9 | 27.1 |
| WMT24 proposed | gen | **8.1** | **32.1** | **29.4** |
| MT$_{sp}$ | gen | 6.8 | 32.8 | 28.7 |

Table 2: This table presents the BLEU, chrF, and chrF++ scores for the English-to-Manipuri machine translation system.

## 3.1 Qualitative Analysis

In the table 3, sample 1, the word ꯑꯍꯤꯡ, meaning "tomorrow" is correctly translated in both models. The word "movie" has been translated to a more beautiful word in both translations as ꯍꯥꯎ ꯃꯤꯠꯍꯤ, which we call movies in the early period, while in the reference, it is translated as cinema, which is not so accurate. While "Mom" has been translated as ꯑꯃꯥ in the reference but it is transliterated in MT$_{sp}$ ꯃꯝ. In the second sample, the phrase "school and I" is accurately translated; the reference ꯑꯩ ꯁ꯭ꯀꯨꯜ is correctly represented in the output as ꯑꯩꯅ ꯁ꯭ꯀꯨꯜ, but the overall meaning of the sentence is not conveyed as the reference text is "not to go" while both translations have translated it as "go". In the third sample, the word "holiday" is translated properly, with the reference being ꯍꯣꯂꯤꯗꯦ, here both models show a better translation than the reference text. In the fourth sample, the phrase "14 April right" is accurately translated as ꯑꯦꯞꯔꯤꯜ ꯱꯴, but in the WMT24$_{proposed}$, the word o.t.p ꯑꯣ.ꯇꯤ.ꯄꯤ has been included, which changes the overall meaning. In the fifth sample, the name "Lelina" is correctly translated as ꯂꯦꯂꯤꯅ in the WMT24$_{proposed}$, but the meaning of the sentence cannot be conveyed as the phrase "thank you" has not been translated. The word "thank you" has been translated in the MT$_{sp}$. Still, the name "Lelina" is not translated. In the sixth sample, "Ambedkar Jayanti" is correctly translated as ꯑꯝꯕꯦꯗꯀꯔ ꯖꯌꯟꯇꯤ in the WMT24$_{proposed}$ model; however, the adequacy is hampered by the missing translation of "tomorrow" in the output, and the fluency is also affected by the ill-formed sentence structure. Meanwhile, in the MT$_{sp}$ model, the word "Jayanti" is missing. In sample 7, the word "municipal" ꯃꯨꯅꯤꯁꯤꯄꯦꯜ has been translated in both models, while MT$_{sp}$ performs better. Some keywords have been translated like the "senior citizen" ꯁꯤꯅꯤꯌꯔ ꯁꯤꯇꯤꯖꯦꯟ. In sample 8, the word "motorcycle" ꯃꯣꯇꯔꯁꯥꯏꯀꯜ is included in the WMT24$_{proposed}$, which is an extra word.

In the table 4, sample 1, the words "shoes," "clothes," "tie," "jewelry," "hairstyle," "make-up," "watch," "cosmetics," and "perfume" have been translated in both models. In the second sample, "dry" is translated as ꯑꯀꯥꯡ and "stone" as ꯅꯨꯡ; in both samples, the overall meaning is conveyed. In sample 3, "chilli powder" ꯑꯌꯤꯡ ꯍꯤꯗꯥꯛ has been translated correctly. In sample 4, the phrase "metro station" has been translated correctly in both models, but in this case, the MT$_{sp}$ model performs better. In the sample 5, the word "Xeres" has not been translated, and the overall meaning of the sentence cannot be conveyed. In the last sample, while the output contains some keywords from the reference, it fails to translate the overall meaning of the sentence.

Four native speakers assessed the adequacy and fluency of the translations. The overall output of the sample has been shown in the figure 1. This evaluation indicates the quality of the sample outputs, reflecting how fluent and adequate the translations are in conveying the intended meaning.



Figure 1: Adequacy and Fluency for the output samples

## 4 Conclusion

We develop and evaluate two Transformer-based machine translation (MT) systems tested on two different datasets (Indic-Conv and Indic-Gen) for translating English to Manipuri. One system (MT$_{sp}$) utilizes the OpenNMT toolkit with Senten-

| Result | Samples for Indic22-Conv test dataset |
|---|---|
| **Source 1:** | Mom, let's go for a movie tomorrow. |
| **Reference 1:** | ꯅꯦ, ꯑꯔꯥꯏ ꯂꯤꯡꯅ ꯑꯃꯔꯣꯄ ꯆꯠꯂꯁꯤ॥ |
| **WMT24<sub>proposed</sub> 1:** | ꯅꯦꯝ ꯑꯔꯦ ꯑꯔꯥꯏꯒꯤ ꯁꯤꯅ ꯅꯤꯡꯊꯧ ꯑꯃꯆꯤꯊꯥ ꯆꯠꯅꯥ ꯑꯣꯏ ॥ |
| **MT<sub>sp</sub> 1:** | ꯁꯦꯝ, ꯑꯔꯥꯏ ꯁꯤꯅ ꯅꯤꯡꯊꯧ ꯑꯃꯔ ꯆꯠꯂꯤ॥ |
| **Source 2:** | I don't have to go to school. |
| **Reference 2:** | ꯃꯥꯢ ꯂꯥꯢꯔꯤꯛ ꯆꯠꯂꯣꯢ ॥ |
| **WMT24<sub>proposed</sub> 2:** | ꯃꯥꯡ ꯂꯥꯢꯔꯤꯛ ꯁꯦꯝ ꯅꯦꯢꯇꯦ॥ |
| **MT<sub>sp</sub> 2:** | ꯃꯥꯢ ꯂꯥꯢꯔꯤꯛ ꯆꯠꯂꯤ ॥ |
| **Source 3:** | It is a holiday. |
| **Reference 3:** | ꯂꯩꯄꯥꯢ॥ |
| **WMT24<sub>proposed</sub> 3:** | ꯃꯁꯤ ꯂꯩꯄꯥ ꯑꯣꯏ॥ |
| **MT<sub>sp</sub> 3:** | ꯃꯁꯤ ꯂꯩꯄꯥ ꯂꯩꯔꯤꯏ॥ |
| **Source 4:** | Oh, tomorrow is the 14th of April right? |
| **Reference 4:** | ꯑꯣ, ꯑꯔꯥꯢꯗꯤ ꯑꯦꯄ꯭ꯔꯤꯜꯒꯤ ꯱꯴ ꯅꯤꯡꯊꯩꯔ ? |
| **WMT24<sub>proposed</sub> 4:** | ꯑꯣ.ꯑꯦꯜ.ꯑꯥꯔ. ꯑꯔꯦ ꯑꯦꯄ꯭ꯔꯤꯜ ꯱꯴ꯅꯤ॥ |
| **MT<sub>sp</sub> 4:** | ꯑꯣ, ꯑꯔꯥꯏꯒꯤ ꯑꯦꯄ꯭ꯔꯤꯜ ꯱꯴ ꯑꯔꯦꯅꯤ |
| **Source 5:** | Thank you, Lelina. |
| **Reference 5:** | ꯊꯥꯒꯠꯆꯔꯤ, ꯂꯦꯂꯤꯅꯥ॥ |
| **WMT24<sub>proposed</sub> 5:** | ꯇꯦꯡꯛ ꯌꯨ ꯂꯦꯂꯤꯅꯥ, ꯂꯦꯂꯤꯅꯥ ꯊꯣꯡꯂꯤꯟ॥ |
| **MT<sub>sp</sub> 5:** | ꯑꯄꯤꯛ ꯊꯥꯒꯠꯆꯔꯤ |
| **Source 6:** | It is Ambedkar Jayanti tomorrow! |
| **Reference 6:** | ꯑꯔꯥꯏ ꯑꯝꯕꯦꯗꯀꯔ ꯖꯌꯟꯇꯤ ꯅꯤ! |
| **WMT24<sub>proposed</sub> 6:** | ꯃꯁꯤ ꯑꯝꯕꯦꯗꯀꯔ ꯖꯌꯟꯇꯤ! |
| **MT<sub>sp</sub> 6:** | ꯃꯁꯤ ꯑꯔꯥꯏꯅꯤ ꯑꯝꯕꯦꯗꯀꯔ |
| **Source 7:** | Even the municipal corporation people also worked round the clock so that they can get the electricity back on time as there were kids and senior citizens present who were facing a lot of difficulties. |
| **Reference 7:** | ꯑꯡꯥꯡꯁꯤꯡ ꯑꯃꯁꯨꯡ ꯑꯍꯜ ꯑꯍꯣꯕꯁꯤꯡ ꯑꯗꯨꯒꯤ ꯃꯔꯝꯅ ꯏꯔꯩꯅꯡꯕꯗꯨ ꯃꯥꯏꯡꯆꯠ ꯁꯦꯝꯕ ꯃꯤꯁꯤꯡ ꯇꯧꯔꯕꯗꯤ ꯁꯤꯔꯔꯒ ꯑꯣꯏꯅ ꯁꯦꯝꯕ ꯃꯔꯝꯅ ꯁꯤꯖꯤꯜꯂꯕ ꯊꯕꯥ ꯂꯩ ॥ |
| **WMT24<sub>proposed</sub> 7:** | ꯃꯤꯖꯤꯜꯒꯤ ꯀꯣꯔꯄꯣꯔꯦꯁꯟ ꯇꯥꯢꯡ ꯁꯤꯗꯦꯢ ꯃꯍꯥꯢ ꯊꯥꯔꯩ ꯇꯥꯡꯒꯤ ꯁꯤꯄ ꯃꯃꯨꯛꯒꯤ ꯑꯣꯏꯅ ꯂꯥꯢ ꯂꯟꯆꯤꯔ ꯁꯥꯠꯏ ꯁꯍ ꯑꯣꯢꯂꯛ ꯁꯤꯖꯤꯟꯕ ꯑꯄꯨꯛ ꯃꯗꯨꯒꯤ ꯑꯃꯔꯩ ꯑꯃ ꯇꯣꯛꯏ॥ |
| **MT<sub>sp</sub> 7:** | ꯃꯤꯖꯤꯜꯒꯤ ꯀꯣꯔꯄꯣꯔꯦꯁꯟ ꯁꯤꯔꯦꯝ ꯑꯔꯩ ꯑꯃꯔ ꯑꯃ ꯁꯥꯡꯗꯥꯔ ꯑꯄꯨꯛ ꯑꯁꯥꯢꯔꯟꯁ ꯌꯥꯠꯕ ꯃꯗꯨꯗ ꯑꯄꯨꯛ ꯁꯥꯒꯜ ꯁꯍꯗꯤ, ꯃꯁꯤ ꯁꯍ ꯑꯣꯢꯂꯛ ꯑꯥ ꯁꯤꯝ ꯁꯤꯖꯤꯟ ꯑꯃꯒꯤ ꯂꯤꯅꯥꯔ ꯂꯥꯢꯡꯂꯕꯗ ꯅꯨꯡꯇꯤꯁꯤꯡ ꯁꯤꯝ ꯑꯁꯥꯢꯔꯟꯁ ꯂꯣꯢ॥ |
| **Source 8:** | There are a lot of organisations here which are catering help to the people, in terms of groceries, medical facilties and medicines and all the necessary items as and when it is needed. |
| **Reference 8:** | ꯁꯦꯝꯠ ꯃꯌꯨꯝꯗ, ꯃꯁꯤ ꯇꯣꯔꯥꯏꯕꯥꯒꯤ ꯁꯤꯖꯤꯛ ꯁꯦꯡ ꯑꯃꯁꯨꯡ ꯑꯄꯤꯛ ꯇꯦꯡꯕꯥꯡꯕ ꯑꯃꯁꯨꯡ ꯁꯤꯠ ꯊꯣꯢꯗꯣꯛꯁꯤꯡ ꯑꯀꯛꯅꯕ ꯑꯃꯁꯨꯡ ꯁꯤꯔꯔ ꯁꯦꯠ ꯁꯥꯕꯤꯛ ꯁꯩ ꯊꯧ ꯃꯌꯥꯢ ꯁꯤꯔꯠꯆꯔꯕꯗ ꯃꯊꯧꯕ ꯁꯤꯔꯦꯛ ꯃꯌꯨꯝꯗ ꯁꯥꯄꯤꯟ ꯇꯥꯅꯨꯛꯁ ꯑꯃꯔꯤ ꯑꯃ ꯗꯤ ॥ |
| **WMT24<sub>proposed</sub> 8:** | ꯃꯖ ꯑꯔꯦꯅ ꯁꯤꯃꯗꯨ ꯑꯃꯔꯤ ꯑꯃ ꯗꯤ, ꯃꯥꯗꯤ ꯑꯀꯛꯅꯕ ꯇꯦꯡꯕꯥꯡꯕꯒꯤ, ꯃꯦꯗꯤꯀꯦꯜꯇꯤꯒꯤ ꯂꯥꯢꯅꯕꯒꯤ, ꯁꯤꯠ–ꯇꯔꯥꯏꯒꯤꯗ, ꯇꯣꯔꯥꯢꯕꯥ–ꯇꯦꯡꯕꯥꯡꯕ ꯑꯃꯔꯤ ꯑꯀꯛꯅꯕꯒꯤ ꯁꯤꯃꯗꯨ ꯁꯩ ꯊꯧ ꯁꯤꯔꯠꯆꯔꯤ ॥ |
| **MT<sub>sp</sub> 8:** | ꯇꯦꯡꯕꯥꯡꯕꯒꯤ ꯁꯤꯔꯦꯛ ꯁꯤꯔꯠꯕꯥ ꯁꯤꯃꯌꯨꯝꯗ ꯗꯨꯡ ꯃꯌꯨꯝ ꯇꯥꯡꯢ, ꯃꯌꯨ ꯁꯣꯡꯕ ꯑꯀꯛꯅꯤ ꯁꯩ ꯊꯧ ꯁꯤꯖꯤꯛꯔ ꯅꯨꯡꯇꯤꯀꯥꯢ ꯑꯀꯛꯅꯤ ꯇꯣꯔꯥꯏꯕꯥꯒꯤ ꯁꯩ ꯊꯧ ꯁꯤꯔꯠꯁ॥ |

Table 3: Sample input and output of the English to Manipuri MT system on the Indic22-Conv test dataset.

| Result | Samples for Indic22-Gen test dataset |
|---|---|
| **Source 1:** | An appearance is a bunch of attributes related to the service person, like their shoes, clothes, tie, jewellery, hairstyle, make-up, watch, cosmetics, perfume, etc |
| **Reference 1:** | ꯑꯄꯤꯌꯔꯦꯟꯁ ꯑꯦꯜꯇꯥ ꯁꯔꯚꯤꯁ ꯄꯔꯁꯟꯒ ꯃꯔꯤ ꯂꯩꯅꯕ ꯁꯨꯔꯨꯛ ꯁꯝꯍꯥꯎꯅ, ꯁꯨꯖ, ꯄꯣꯠꯁꯤꯟ, ꯇꯥꯏ ꯁꯥꯖꯦꯠ, ꯃꯩꯛꯑꯞ, ꯁꯝꯁꯨ ꯑꯃꯁꯨꯡ ꯄꯔꯐꯤꯎꯝ ꯑꯁꯤꯅꯆꯤꯡꯕ ꯑꯣꯏꯕ ꯀꯔꯤ꯫ |
| **WMT24_proposed 1:** | ꯑ꯭ꯄꯤꯌꯔꯦꯟꯁ ꯁꯝꯍꯥꯎꯅ ꯁꯔꯚꯤꯁ ꯄꯔꯁꯟꯒ, ꯇꯥꯏ, ꯇ꯭ꯌꯥꯔ, ꯇꯨ, ꯖꯨꯜꯔꯤ, ꯍꯦꯌꯔꯁꯇꯥꯏꯜ, ꯁꯝꯁꯨ, ꯁꯔꯟꯒ, ꯑꯦꯇ, ꯃꯦꯛꯑꯞ, ꯄꯔꯐꯤꯎꯝ, ꯑꯃꯁꯨꯡ, ꯄꯣꯠꯁꯤꯟ ꯁꯔꯚꯤꯁ ꯄꯔꯁꯟꯒ꯫ |
| **MT_sp 1:** | ꯄꯔꯒ ꯁꯔꯚꯤꯁ ꯄꯔꯁꯟꯒ ꯃꯔꯤ ꯂꯩꯅꯕ ꯁꯝꯍꯥꯎꯅ, ꯁꯨꯔꯨꯛ ꯁꯝꯍꯥꯎꯅ, ꯇꯥꯏꯒꯤꯟ, ꯄꯣꯠꯁꯤꯟ, ꯖꯨꯜꯒꯤꯟ, ꯍꯦꯌꯔꯁꯇꯥꯏꯜ, ꯃꯩꯛꯑꯞ, ꯁꯝꯁꯨ, ꯄꯔꯐꯤꯎꯝ ꯑꯁꯤꯅꯆꯤꯡꯕ꯫ |
| **Source 2:** | Make this into powder with a dry grinder or in a stone pestle. |
| **Reference 2:** | ꯃꯁꯤ ꯑꯥꯡꯕ ꯒ꯭ꯔꯥꯏꯟꯗꯔ ꯅꯇ꯭ꯔꯒ ꯅꯨꯡ ꯄꯦꯁꯇꯜꯗ ꯃꯄꯥꯟ ꯑꯣꯏꯍꯟꯎ꯫ |
| **WMT24_proposed 2:** | ꯃꯁꯤ ꯁꯤꯡꯖ ꯁꯇꯣꯅ ꯄꯦꯁꯇꯜ ꯑꯣꯏꯅ ꯁꯨꯔꯨꯛ ꯁꯦꯝꯅꯕ ꯇꯧꯒꯗꯕꯅꯤ꯫ |
| **MT_sp 2:** | ꯃꯁꯤ ꯑꯥꯡꯕ ꯒꯔꯤ ꯑꯣꯏꯅ ꯁꯦꯝꯅꯕ ꯇꯧꯖꯧ ꯫ |
| **Source 3:** | The use of chilli powder in this region is done cautiously. |
| **Reference 3:** | ꯃꯐꯝ ꯑꯁꯤꯗ ꯃꯔꯤꯛꯄ ꯄꯥꯎꯗꯔ ꯁꯤꯖꯤꯟꯅꯕ ꯆꯦꯛꯁꯤꯟꯅ ꯇꯧꯏ꯫ |
| **WMT24_proposed 3:** | ꯇꯔꯦꯠ ꯑꯁꯤꯗ ꯁꯤꯇꯤꯕ ꯁꯦꯝꯒꯠ ꯇꯧꯒꯅꯤ ꯫ |
| **MT_sp Output 3:** | ꯇꯔꯦꯠ ꯑꯁꯤꯗ ꯁꯤꯇꯤ ꯁꯦꯝꯒꯠꯄ ꯃꯔꯤꯛꯄ ꯄꯥꯎꯗꯔ ꯁꯤꯖꯤꯟꯅꯤ ꯫ |
| **Source 4:** | The nearest Delhi Metro station is Arjan Garh, on the Yellow Line. |
| **Reference 4:** | ꯅꯛꯁꯤꯟꯕ ꯗꯦꯜꯍꯤ ꯃꯦꯇꯔꯣ ꯁ꯭ꯇꯦꯁꯟꯁꯤ ꯑꯔꯖꯟ ꯒꯥꯔ ꯑꯃꯁꯨꯡ ꯏꯂꯣ ꯂꯥꯏꯟꯗꯅꯤ꯫ |
| **WMT24_proposed 4:** | ꯅꯛꯁꯤꯟꯕ ꯗꯦꯜꯍꯤ ꯃꯦꯇꯔꯣ ꯌꯦꯜꯂꯣ ꯑꯁꯤ ꯑꯔꯖꯟ ꯒꯥꯔ ꯇꯦꯟꯒꯤꯟ ꯇꯦ ꯫ |
| **MT_sp 4:** | ꯅꯛꯁꯤꯟꯕ ꯗꯦꯜꯍꯤ ꯁ꯭ꯇꯦꯁꯟ ꯃꯦꯇꯔꯣ ꯌꯦꯜꯂꯣ ꯂꯥꯏꯟ ꯑꯔꯖꯟ ꯒꯥꯔꯒꯤ ꯃꯄꯥꯟꯗꯅꯤ꯫ |
| **Source 5:** | After him, came Xerxes II for a short while. |
| **Reference 5:** | ꯃꯍꯥꯛꯒꯤ ꯃꯇꯨꯡꯗ ꯃꯇꯝ ꯁꯥꯡꯗ ꯄꯔꯛꯇ ꯇꯔꯁꯤꯜ II ꯂꯥꯛꯏ꯫ |
| **WMT24_proposed 5:** | ꯃꯍꯥꯛꯒꯤ ꯃꯇꯨꯡꯗ ꯇꯔꯁꯤꯟ XII ꯑꯁꯤ ꯁꯥꯡꯔꯛꯇ ꯄꯔꯛꯅ ꯂꯥꯛꯏ꯫ |
| **MT_sp 5:** | ꯃꯍꯥꯛꯒꯤ ꯃꯇꯨꯡꯗ, ꯃꯍꯥꯛ ꯃꯇꯝ ꯁꯥꯡꯗ ꯄꯔꯛꯇ XI ꯂꯥꯛꯏ꯫ |
| **Source 6:** | In Karaikal liquor is cheaper than in the neighbouring Tamil Nadu, there are quite a few decent bars in Karaikal - the Niagra bar in the Nanda hotel, the Thunder bar in the Paris International, the City bar - a very famous one in the town, and The Sea Gulls Restaurant owned by the government of Pondicherry which is at the sea shore and is good to hang out in the evenings. |
| **Reference 6:** | ꯀꯔꯥꯏꯀꯥꯜꯗ ꯌꯨ ꯑꯁꯤ ꯌꯨꯝꯂꯣꯟꯕ ꯇꯦꯝꯁꯤꯟꯒꯤ ꯇꯝꯤꯜ ꯅꯥꯗꯨ, ꯀꯔꯥꯏꯀꯥꯜ ꯃꯤ ꯁꯦꯝꯒꯠ ꯁꯝꯨ ꯇꯣꯡꯕꯅ ꯄꯔꯛꯇ ꯌꯦꯝ ꯑꯃꯇ ꯇꯦ – ꯅꯟꯗ ꯍꯣꯇꯦꯜꯒꯤ ꯅ ꯇꯦꯔꯥꯕꯔ ꯕꯥꯔ– ꯄꯦꯔꯤꯁ ꯏꯟꯇꯔꯅꯦꯁꯅꯦꯜꯒꯤ ꯅ ꯁꯦꯝꯒꯠ ꯕꯥꯔ– ꯁꯤꯇꯤ ꯕꯥꯔ ꯑꯁꯤꯗ ꯃꯐꯝ ꯃꯌꯥꯝ ꯂꯧꯕ ꯕꯥꯔ– ꯑꯃꯁꯨꯡ ꯄꯟꯗꯤꯆꯦꯔꯤꯒꯤ ꯁꯔꯀꯥꯔꯒꯤ ꯃꯈꯥ ꯄꯣꯟꯕ ꯁꯤ ꯒꯜ ꯔꯦꯁ꯭ꯇꯣꯔꯦꯟꯠ ꯑꯁꯤ ꯄꯥꯎꯗ ꯃꯔꯤ ꯅꯨꯃꯤꯠ ꯂꯩꯁꯥꯕꯗ ꯃꯐꯝ ꯐꯖꯕ ꯑꯃꯅꯤ꯫ |
| **WMT24_proposed 6:** | ꯀꯔꯥꯏꯀꯥꯟ ꯑꯁꯤꯗ ꯌꯨ ꯌꯨꯝꯂꯣꯟꯕ ꯇꯦꯝꯁꯤꯟ ꯇꯦ ꯀꯔꯥꯏꯀꯥꯜ ꯁꯤ ꯇꯝꯤꯜ ꯅꯥꯗꯨ ꯅꯤꯡꯈꯥꯏ– ꯁꯨꯖꯠ ꯑꯁꯤ ꯀꯔꯥꯏꯀꯥꯜ ꯇꯝꯤꯜ, ꯅꯟꯗ ꯍꯣꯇꯦꯜꯒꯤ ꯅꯥꯒꯔ ꯕꯥꯔ ꯁꯦꯝꯒꯠ ꯁꯤ ꯕꯥꯔ ꯄꯦꯔꯤꯁ ꯇꯦꯔꯥꯕꯔ ꯑꯁꯤ ꯁꯦꯝꯒꯠ ꯑꯁꯤ ꯃꯌꯥꯃꯗ ꯂꯥꯏꯔꯦ ꯄꯟꯗꯤ ꯑꯃꯁꯨꯡ ꯇꯦ ꯑꯦꯞꯔꯤ ꯏꯟꯇꯔ ꯑꯁꯤ ꯃꯐꯝ ꯐꯖꯕ ꯁꯝꯔꯨꯛ ꯁꯝꯔꯨꯛ꯫ ꯐꯖꯟꯄꯣꯛꯇꯔꯦ ꯂ꯭ꯌꯣꯖ ꯑꯃꯁꯨꯡ ꯃꯗꯨꯗ ꯃꯐꯝ ꯃꯔꯤꯗꯅꯤ꯫ |
| **MT_sp Output 6:** | ꯀꯔꯥꯏꯀꯥꯟꯗ ꯌꯨ ꯌꯨꯝꯂꯣꯟꯕ ꯃꯔꯥꯡ ꯄꯥꯎꯗ ꯁꯝꯍꯥꯎ, ꯁꯦꯝꯒꯠ ꯇꯦꯝꯁꯤꯟ, ꯁꯦꯝꯒꯠ ꯇꯦꯝꯁꯤꯟ ꯄꯔꯤꯁꯒꯤ, ꯃꯁꯤ ꯇꯣꯡꯅꯤꯔꯦ ꯅꯣꯡꯒ ꯃꯗꯨ ꯁꯥꯡꯅ ꯄꯨꯡꯇꯥꯔꯨꯛ ꯂꯧꯕ – ꯁꯤꯇꯤꯅꯤꯟꯁꯤꯅꯍꯤ, ꯑꯃꯁꯨꯡ, ꯌꯨꯃꯂꯣꯟ, ꯏꯟꯇꯔꯒꯤ ꯄ꯭ꯔꯥꯏꯁ ꯁꯝꯍꯥꯎꯅ ꯃꯔꯤ ꯁꯨꯡ ꯃꯔꯤꯛ ꯃꯔꯤ ꯁꯨꯡꯅꯤ꯫ |

Table 4: Sample Input and Output of the English to Manipuri MT System on the Indic22-Gen test dataset.

cepiece tokenization for tokenization, while the proposed model (WMT24_proposed) employs the Open-NMT toolkit with Byte Pair Encoding (BPE). Both models are trained on a comprehensive dataset that includes data from WMT24 and MeitY.

The model is optimized with the Adam optimizer and is evaluated using BLEU, chrF, and chrF++ metrics. Additionally, the translations are assessed

for both adequacy and fluency. The models successfully convey the overall meaning of the source sentences, but they often lack fluency, producing disjointed or grammatically incorrect outputs.

Overall, the WMT24$_{proposed}$ produces translations that are more syntactically correct, contextually appropriate, and idiomatically fluent, while MT$_{sp}$ offers more direct, simpler translations that sometimes lose nuance or complex structure.

## Limitations

The proposed (WMT24$_{proposed}$) model translation conveys the main ideas of the reference sentence, despite certain errors and structural challenges. It captures some aspects of the overall meaning of the reference sentences. In the case of longer sentences, there is a large amount of adequacy. However, the fluency of these translations deteriorates as the length of the input sentences increases.

## Acknowledgements

## References

Raj Dabre and Anoop Kunchukuttan. 2024. Findings of wmt 2024's multiindic22mt shared task for machine translation of 22 indian languages. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. OpenNMT: Open-source toolkit for neural machine translation. In *Proceedings of ACL 2017, System Demonstrations*, pages 67–72, Vancouver, Canada. Association for Computational Linguistics.

T Kudo. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Kshetrimayum Boynao Singh, Avichandra Singh Ningthoujam, Loitongbam Sanayai Meetei, Sivaji Bandyopadhyay, and Thoudam Doren Singh. 2023a. NITS-CNLP low-resource neural machine translation systems of English-Manipuri language pair. In *Proceedings of the Eighth Conference on Machine Translation*, pages 967–971, Singapore. Association for Computational Linguistics.

Kshetrimayum Boynao Singh, Ningthoujam Avichandra Singh, Loitongbam Sanayai Meetei, Ningthoujam Justwant Singh, Thoudam Doren Singh, and Sivaji Bandyopadhyay. 2023b. A comparative study of transformer and transfer learning MT models for English-Manipuri. In *Proceedings of the 20th International Conference on Natural Language Processing (ICON)*, pages 791–796, Goa University, Goa, India. NLP Association of India (NLPAI).

Salam Michael Singh and Thoudam Doren Singh. 2020. Unsupervised neural machine translation for English and Manipuri. In *Proceedings of the 3rd Workshop on Technologies for MT of Low Resource Languages*, pages 69–78, Suzhou, China. Association for Computational Linguistics.

Salam Michael Singh and Thoudam Doren Singh. 2022a. An empirical study of low-resource neural machine translation of manipuri in multilingual settings. *Neural Comput. Appl.*, 34(17):14823–14844.

Salam Michael Singh and Thoudam Doren Singh. 2022b. Low resource machine translation of english–manipuri: A semi-supervised approach. *Expert Systems with Applications*, 209:118187.

Thoudam Doren Singh and Sivaji Bandyopadhyay. 2005. Manipuri morphological analyzer. In *In the Proceedings of the Platinum Jubilee International Conference of LSI, Hyderabad, India*.

Thoudam Doren Singh and Sivaji Bandyopadhyay. 2010. Manipuri-English bidirectional statistical machine translation systems using morphology and dependency relations. In *Proceedings of the 4th Workshop on Syntax and Structure in Statistical Translation*, pages 83–91, Beijing, China. Coling 2010 Organizing Committee.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

# NLIP_Lab-IITH Multilingual MT System for WAT24 MT Shared Task

**Maharaj Brahma**     **Pramit Sahoo**     **Maunendra Sankar Desarkar**

Natural Language and Information Processing Lab (NLIP)

Indian Institute of Technology Hyderabad

Hyderabad, India

{cs23resch01004,ai23mtech14004}@iith.ac.in, maunendra@cse.iith.ac.in

## Abstract

This paper describes NLIP Lab's multilingual machine translation system for the WAT24 shared task on multilingual Indic MT task for 22 scheduled languages belonging to 4 language families. We explore pre-training for Indic languages using alignment agreement objectives. We utilize bi-lingual dictionaries to substitute words from source sentences. Furthermore, we fine-tuned language direction-specific multilingual translation models using small and high-quality seed data. Our primary submission is a 243M parameters multilingual translation model covering 22 Indic languages. In the IN22-Gen benchmark, we achieved an average chrF++ score of 46.80 and 18.19 BLEU score for the En-Indic direction. In the Indic-En direction, we achieved an average chrF++ score of 56.34 and 30.82 BLEU score. In the In22-Conv benchmark, we achieved an average chrF++ score of 43.43 and BLEU score of 16.58 in the En-Indic direction, and in the Indic-En direction, we achieved an average of 52.44 and 29.77 for chrF++ and BLEU respectively. Our model[1] is competitive with IndicTransv1 (474M parameter model).

## 1 Introduction

Multilingual Neural Machine Translation (MNMT) has shown remarkable success in building translation systems for world languages in a single model (Johnson et al., 2017). These successes have led researchers to increase the model capacity catering to hundreds of world languages (Fan et al., 2020), (NLLB Team et al., 2022). It also led to multilingual translation models for particular languages under particular geographical groups such as Indic (Ramesh et al., 2022; Gala et al., 2023), African (Nekoto et al., 2020). Indic languages are interesting, with diverse languages belonging to various language families and written scripts.

This paper describes our system submission for the WAT 24 MultiIndic22MT task (Dabre and Kunchukuttan, 2024), which includes 22 scheduled Indian languages belonging to 4 language families across 12 written scripts. We participated in the constrained translation task. We explore an alignment agreement-based pre-training objective. Specifically, we substitute words from source sentences for equivalent words in a random language. The pre-training data consists of a sentence pair from the original data and code-switched augmented sentences. Our primary submission is a fine-tuned transformer-based multilingual model with 243M parameters. Experimental results show that our system achieves an average chrF++ score of 46.80 for the En-Indic direction in the IN22-Gen benchmark. We achieved an average chrF++ score of 46.80 and 18.19 BLEU score for the En-Indic direction. In the Indic-En direction, we achieved an average chrF++ score of 56.34 and 30.82 BLEU score. Compared with the IndicTransv2 model for Indic-Indic translation, our system lags most minor for pan_Guru-snd_Deva with 0.3 chrF++ scores. Due to computational constraints, we train our model on a reduced corpus.

## 2 Dataset

### 2.1 Pre-training data

In this section, we describe the dataset used for pre-training. We use the official Bharat Parallel Collection Corpus (BPCC) (Gala et al., 2023) but reduce the corpus size due to computational limitations. We also exclude sentences from the comparable directory. For languages with over 10 million parallel sentences, we reduce the no. of sentences by half. The corpus statistics are shown in Table 1. To handle skew data distribution and have good representation for low-resource languages, we use heuristic-based temperature sampling (Arivazhagan et al., 2019; Conneau et al., 2020) for data

---

[1]Our code and models are available at https://github.com/maharajbrahma/WAT2024-MultiIndicMT

sampling with temperature sampling (T = 5) shown in Figure 1. We utilize small, high-quality data from BPCC, namely ILCI, Massive, NLLB Seed, Daily, and Wiki, for direction-specific fine-tuning.

## 2.2 Alignment Augmentation

For alignment augmentation we English-centric bilingual dictionaries from MUSE[2] and GATITOS[3]. We use top 4000 words in dictionaries, replaced with a probability of 30% from the bi-lingual dictionary. We consider only replacing words in the languages that have dictionaries.

## 3 Methodology

Our pre-training approach is inspired by aligning embeddings (Lin et al., 2020; Yang et al., 2020) through substituting words from a bi-lingual dictionary. We pre-trained a universal model that covers En-Indic and Indic-En. We named this model **"IndicRASP"**. IndicRASP is fine-tuned into a language direction-specific model called **"IndicRASP-Seed"** using small and high-quality seed data.

**IndicRASP (IR)**: IndicRASP is pre-trained on data from 22 Indic languages sourced from BPCC. During pre-training, we randomly substitute English words for corresponding Indic language words, resulting in code-switched augmented sentences. The alignment augmentation technique helps to bring semantically similar embeddings closer together. We get 56M sentences after alignment augmentation. We combined training sentences from the original En-Indic and Indic-En[4] and obtained 282M sentences for pre-training.

**IndicRASP-Seed (IR Seed)**: To further enhance the performance of IndicRASP, we fine-tuned the model to be language-direction specific. We consider high-quality seed data from the BPCC corpus: ILCI, NLLB Seed, Massive, Daily, and Wiki. We sampled a total of 2.26M sentences and fine-tuned IndicRASP for both En-Indic and Indic-En directions.

## 4 Experiments

**Setting:** We use the standard sequence-to-sequence Transformer big model as our architecture for pre-training. It uses 6 encoder and 6 decoder layers, with an embedding size of 1024. The embeddings between the encoder and decoder are shared, with a feed-forward network size of 4096 and 16 attention heads.

**Training:** We pre-train the model with the Adam optimizer (Kingma and Ba, 2014) with $\beta_1$ set to 0.9 and $\beta_2$ set to 0.98. We set the warmup initial learning rate to 1e-07 and the learning rate to 5e-4, with a warmup step of 4000. We train the models with a dropout rate of 0.1 and a label smoothing rate of 0.1. During fine-tuning, we consider a learning rate of 3e-5 and a dropout rate of 0.2. All experiments are conducted on 8 NVIDIA A100 GPUs.

**Baseline Models:** We consider two baselines:

1. **IndicTransv1 (Ramesh et al., 2022)**: Indic-Transv1 (IT1) is a multilingual transformer (Vaswani et al., 2017) translation model for 11 Indic languages trained. It is a 474M parameter trained on the 49.7M sentence pair on the Samanantar dataset.

2. **IndicTransv2 (Gala et al., 2023)**: Indic-Transv2 (IT2) is a 1B parameter model trained on the BPCC corpus for 22 Scheduled Indian languages.

**Language-Direction Specific Models:** For our primary submission, we fine-tune IndicRASP with direction-specific small seed data for En-Indic and Indic-En. For the Indic-Indic model, we fine-tune the IndicRASP-Seed (En-Indic direction) on the Indic-Indic corpus extracted from the BPCC corpus.

**Evaluation:** We use the dev set of BPCC IN-Gen as our validation and evaluate our model on the test set of BPCC IN-Gen and IN-Conv. We report our results on lexical-based automatic metrics BLEU (Papineni et al., 2002), and chrF++ (Popović, 2017). We use the sacreBLEU library for evaluation, with a chrF word order of 2.

## 5 Results

We list the results of our model on the IN22-Gen in Table 3, 4 for chrF++ and BLEU, respectively. Similarly, Table 5 and 6 results for chrF++ and BLEU in IN22-Conv. Table 2 shows the performance of our primary submission on a hidden test set. Our findings described for IN22-Gen are:

- IndicRASP achieves an average chrF++ score of 45.50, and IndicRASP-Seed achieves 46.80 with an improvement of (+1.30) for the En-Indic direction. Similarly, IndicRASP

---

| Language | Script | # of sentences (M) | Language | Script | # of sentences (M) |
|----------|--------|-------------------|----------|--------|-------------------|
| Assamese | Bengali | 1.42 | Manipuri | Metei | 0.04 |
| Bodo | Devanagari | 0.12 | Manipuri | Bengali | 0.37 |
| Bengali | Bengali | 16.39 | Marathi | Devanagari | 9.37 |
| Dogri | Devanagari | 0.02 | Nepali | Devanagari | 1.68 |
| Konkani | Devanagari | 0.10 | Odia | Oriya | 5.80 |
| Gujarati | Gujarati | 10.12 | Punjabi | Gurmuki | 9.75 |
| Hindi | Devanagari | 19.24 | Sanskrit | Devanagari | 0.28 |
| Kannada | Kannada | 11.60 | Santali | Olck | 0.02 |
| Kashmiri | Devanagari | 0.20 | Sindhi | Devanagari | 0.01 |
| Kashmiri | Arabic | 0.15 | Tamil | Tamil | 10.18 |
| Maithili | Devanagari | 0.09 | Telugu | Telugu | 11.54 |
| Malayalam | Malayalam | 11.69 | Urdu | Arabic | 2.99 |

Table 1: Statistics of the dataset. Total of 113.65 million bi-texts.



Figure 1: Number of sentences in each language and the sampled distribution with the T=5

| Language pair | BLEU | chrF | chrF++ |
|---------------|------|------|--------|
| asm_Beng-eng_Latn | 19.9 | 50.4 | 47.8 |
| ben_Beng-eng_Latn | 22.1 | 50.1 | 48.0 |
| brx_Deva-eng_Latn | 17.8 | 47.6 | 45.3 |
| guj_Gujr-eng_Latn | 16.8 | 45.4 | 43.2 |
| hin_Deva-eng_Latn | 23.1 | 50.5 | 48.5 |
| kas_Arab-eng_Latn | 12.4 | 38.5 | 36.5 |
| mal_Mlym-eng_Latn | 20.3 | 48.3 | 46.2 |
| npi_Deva-eng_Latn | 18.0 | 46.7 | 44.5 |
| san_Deva-eng_Latn | 9.3 | 34.7 | 32.6 |
| sat_Olck-eng_Latn | 11.0 | 36.3 | 34.0 |
| snd_Deva-eng_Latn | 21.2 | 47.1 | 45.5 |
| tel_Telu-eng_Latn | 13.8 | 40.6 | 38.4 |
| urd_Arab-eng_Latn | 20.3 | 45.6 | 43.9 |

Table 2: Indic-En scores results on hidden test set

achieves an average BLEU score of 16.82 and 18.19 for the IndicRASP-Seed. It suggests that fine-tuning small, high-quality language directions improves the alignment augmented IndicRASP model. We can observe similar results for Indic-En.

- By comparing IT1 and IndicRASP-Seed, we find that IndicRASP-Seed has a chrF++ improvement of +1.30 for En-Indic; however, in the Indic-En direction, IndicRASP-Seed is lagging behind by 1.63.

- By comparing IT2 and IndicRASP-Seed, we find that IndicRASP-Seed lags behind by 1.74 chrF++ scores for En-Indic direction. In the Indic-En direction, the IndicRASP-Seed lags behind significantly by a 7.13 chrF++ score from IT2.

- For En-Indic languages highlighted in bold in Table 3, namely Manipuri, Oriya, and Santali, IndicRASP-Seed performs better than Indic-Transv2 with chrF++ score difference of 0.5, 2.3, and 6.5 respectively.

- We observe that our setup performs better in the En-Indic direction than in Indic-En. This is possibly due to the reduction of the dataset.

In Table 7, we show the performance of IndicRASP-Seed for Indic-Indic direction in the IN22-Gen and IN22-Conv datasets. We observe that the IT2 is better than the IndicRASP-Seed in all language pairs, particularly for mal_Mlym-hin_Deva, IndicRASP-Seed lags highest behind by a 5.6 chrF++ score, and

| Language | En − Indic | | | | Indic − En | | | |
|---|---|---|---|---|---|---|---|---|
| | IT1 | IT2 | IR | IR Seed | IT1 | IT2 | IR | IR Seed |
| asm_Beng | 35.9 | 47.1 | 43.0 | 44.8 | 56.1 | 66.5 | 57.2 | 58.4 |
| ben_Beng | 48.6 | 51.8 | 47.3 | 48.5 | 58.4 | 64.5 | 55.6 | 56.9 |
| brx_Deva | – | 47.8 | 46.3 | 46.5 | – | 61.8 | 53.5 | 54.5 |
| doi_Deva | – | 57.9 | 58.1 | 58.0 | – | 72.7 | 64.0 | 64.5 |
| gom_Deva | – | 45.2 | 42.1 | 42.9 | – | 58.7 | 50.7 | 51.6 |
| guj_Gujr | 47.2 | 53.4 | 47.6 | 48.8 | 60.3 | 66.9 | 57.7 | 59.2 |
| hin_Deva | 53.3 | 56.6 | 52.4 | 54.8 | 60.7 | 65.0 | 57.6 | 58.7 |
| kan_Knda | 46.7 | 50.9 | 46.3 | 47.9 | 58.8 | 65.1 | 55.4 | 56.5 |
| kas_Arab | – | 40.2 | 37.6 | 39.5 | – | 60.5 | 52.6 | 53.8 |
| mai_Deva | – | 48.7 | 46.7 | 47.3 | – | 66.4 | 58.5 | 59.5 |
| mal_Mlym | 45.7 | 50.8 | 45.6 | 47.4 | 56.9 | 64.5 | 54.2 | 56.0 |
| mni_Mtei | – | 44.5 | 44.5 | 45.0 | – | 60.3 | 51.9 | 52.7 |
| mar_Deva | 44.3 | 50.9 | 44.2 | 46.7 | 57.7 | 65.1 | 55.8 | 57.3 |
| npi_Deva | – | 49.0 | 44.8 | 47.8 | – | 69.4 | 60.6 | 62.1 |
| ory_Orya | 40.3 | 43.8 | 43.0 | 46.1 | 60.0 | 67.6 | 57.6 | 59.3 |
| pan_Guru | 48.0 | 50.7 | 48.3 | 47.9 | 57.2 | 63.0 | 54.5 | 56.1 |
| san_Deva | – | 38.6 | 34.6 | 36.2 | – | 56.0 | 45.9 | 46.9 |
| sat_Olck | – | 33.4 | 39.6 | 39.9 | – | 47.7 | 47.2 | 48.2 |
| snd_Deva | – | 36.5 | 34.2 | 35.2 | – | 57.0 | 51.3 | 52.6 |
| tam_Taml | 45.5 | 49.6 | 45.4 | 46.4 | 53.9 | 59.7 | 51.3 | 53.2 |
| tel_Telu | 46.5 | 52.5 | 47.2 | 48.8 | 57.7 | 64.9 | 55.7 | 56.8 |
| urd_Arab | – | 68.0 | 63.1 | 62.4 | – | 73.1 | 63.5 | 64.7 |
| Avg. | 45.64 | 48.54 | 45.50 | 46.80 | 57.97 | 63.47 | 55.10 | 56.34 |

Table 3: chrF++ (↑) scores on IN22-Gen

| Language | En − Indic | | | | Indic − En | | | |
|---|---|---|---|---|---|---|---|---|
| | IT1 | IT2 | IR | IR Seed | IT1 | IT2 | IR | IR Seed |
| asm_Beng | 9.9 | 19.3 | 15.0 | 17.8 | 32.5 | 42.5 | 30.7 | 31.6 |
| ben_Beng | 18.1 | 20.7 | 15.6 | 17.2 | 33.4 | 40.9 | 29.2 | 30.3 |
| brx_Deva | – | 17.0 | 15.9 | 16.2 | – | 39.0 | 27.2 | 28.2 |
| doi_Deva | – | 33.8 | 33.7 | 33.4 | – | 53.7 | 41.2 | 41.7 |
| gom_Deva | – | 18.7 | 14.7 | 16.4 | – | 34.0 | 23.8 | 24.9 |
| guj_Gujr | 17.9 | 25.6 | 18.2 | 19.6 | 36.3 | 43.5 | 31.3 | 32.6 |
| hin_Deva | 28.3 | 33.5 | 27.0 | 28.0 | 36.1 | 40.4 | 29.8 | 30.6 |
| kan_Knda | 13.4 | 17.7 | 13.0 | 15.6 | 34.8 | 40.5 | 29.0 | 30.0 |
| kas_Arab | – | 14.4 | 12.4 | 13.4 | – | 38.6 | 28.3 | 29.5 |
| mai_Deva | – | 19.2 | 17.0 | 17.8 | – | 43.2 | 32.8 | 33.8 |
| mal_Mlym | 13.9 | 16.4 | 12.0 | 13.1 | 31.4 | 41.0 | 28.2 | 30.1 |
| mni_Mtei | – | 17.4 | 17.5 | 18.2 | – | 39.0 | 27.7 | 28.9 |
| mar_Deva | 13.9 | 21.4 | 13.8 | 17.5 | 33.5 | 41.9 | 29.8 | 31.1 |
| npi_Deva | – | 16.8 | 12.6 | 15.6 | – | 48.2 | 35.7 | 38.0 |
| ory_Orya | 10.2 | 14.4 | 12.3 | 17.4 | – | 45.1 | 31.4 | 32.6 |
| pan_Guru | 23.5 | 25.8 | 23.7 | 22.6 | 33.5 | 41.1 | 29.5 | 30.9 |
| san_Deva | – | 10.9 | 8.4 | 9.1 | – | 31.9 | 20.6 | 21.8 |
| sat_Olck | – | 5.5 | 8.7 | 8.8 | – | 25.1 | 23.1 | 24.3 |
| snd_Deva | – | 13.9 | 10.1 | 11.1 | – | 33.4 | 25.8 | 27.0 |
| tam_Taml | 11.9 | 14.7 | 11.3 | 11.7 | 28.9 | 36.1 | 25.6 | 27.1 |
| tel_Telu | 15.5 | 19.7 | 15.3 | 16.2 | 33.5 | 42.5 | 30.5 | 31.5 |
| urd_Arab | – | 49.4 | 41.8 | 43.4 | – | 53.8 | 40.1 | 41.6 |
| Avg. | 16.0 | 20.28 | 16.82 | 18.19 | 30.0 | 40.7 | 29.60 | 30.82 |

Table 4: BLEU (↑) scores on IN22-Gen

| Language | En − Indic | | | | Indic − En | | | |
|---|---|---|---|---|---|---|---|---|
| | IT1 | IT2 | IR | IR Seed | IT1 | IT2 | IR | IR Seed |
| asm_Beng | 36.4 | 46.8 | 40.9 | 44.9 | 52.5 | 62.9 | 52.7 | 57.7 |
| ben_Beng | 47.5 | 49.7 | 45.1 | 47.6 | 55.2 | 58.4 | 51.7 | 55.3 |
| brx_Deva | – | 45.3 | 43.8 | 44.2 | – | 56.3 | 50.1 | 50.9 |
| doi_Deva | – | 53.9 | 55.4 | 55.2 | – | 65.0 | 59.1 | 59.9 |
| gom_Deva | – | 42.5 | 39.8 | 39.9 | – | 51.7 | 46.6 | 47.3 |
| guj_Gujr | 49.1 | 53.1 | 46.9 | 48.5 | 56.9 | 62.0 | 54.7 | 58.1 |
| hin_Deva | 48.6 | 49.6 | 48.0 | 48.2 | 57.4 | 60.1 | 54.8 | 56.7 |
| kan_Knda | 32.6 | 33.8 | 31.7 | 32.3 | 44.0 | 47.5 | 40.4 | 43.9 |
| kas_Arab | – | 35.6 | 28.7 | 34.3 | – | 52.6 | 45.9 | 47.6 |
| mai_Deva | – | 44.3 | 39.8 | 43.0 | – | 57.8 | 52.3 | 52.9 |
| mal_Mlym | 43.8 | 45.7 | 41.7 | 42.9 | 50.6 | 54.3 | 47.2 | 50.7 |
| mni_Mtei | – | 40.2 | 40.8 | 41.1 | – | 52.5 | 48.5 | 49.1 |
| mar_Deva | 43.7 | 48.6 | 42.2 | 44.7 | 54.2 | 58.5 | 50.9 | 55.2 |
| npi_Deva | – | 51.5 | 44.4 | 49.9 | – | 63.0 | 56.0 | 59.1 |
| ory_Orya | 38.9 | 40.2 | 39.1 | 41.4 | 55.6 | 60.3 | 52.4 | 56.6 |
| pan_Guru | 54.0 | 57.8 | 53.1 | 54.1 | 58.1 | 62.7 | 54.8 | 58.5 |
| san_Deva | – | 35.5 | 29.3 | 33.5 | – | 48.3 | 40.2 | 42.6 |
| sat_Olck | – | 34.6 | 41.7 | 41.7 | – | 43.5 | 46.4 | 47.4 |
| snd_Deva | – | 30.3 | 31.8 | 33.2 | – | 49.6 | 49.5 | 50.1 |
| tam_Taml | 37.7 | 39.1 | 37.4 | 38.3 | 44.1 | 45.8 | 40.8 | 43.6 |
| tel_Telu | 42.5 | 45.5 | 40.8 | 42.4 | 48.5 | 52.9 | 45.8 | 49.3 |
| urd_Arab | – | 61.6 | 54.6 | 53.9 | – | 65.5 | 57.4 | 61.2 |
| Avg. | 43.16 | 44.78 | 41.66 | 43.43 | 52.46 | 53.22 | 49.92 | 52.44 |

Table 5: chrF++ (↑) scores on IN22-Conv

| Language | En − Indic | | | | Indic − En | | | |
|---|---|---|---|---|---|---|---|---|
| | IT1 | IT2 | IR | IR Seed | IT1 | IT2 | IR | IR Seed |
| asm_Beng | 11.6 | 19.7 | 15.3 | 18.5 | 31.3 | 43.8 | 31.8 | 36.7 |
| ben_Beng | 20.1 | 21.3 | 17.5 | 19.1 | 32.9 | 36.4 | 29.0 | 32.2 |
| brx_Deva | – | 15.4 | 13.6 | 14.7 | – | 35.5 | 26.8 | 27.9 |
| doi_Deva | – | 32.4 | 34.1 | 34.4 | – | 45.6 | 36.8 | 38.1 |
| gom_Deva | – | 14.2 | 11.3 | 11.2 | – | 29.9 | 23.2 | 23.7 |
| guj_Gujr | 23.2 | 27.2 | 20.9 | 22.3 | 34.7 | 41.1 | 32.0 | 35.4 |
| hin_Deva | 28.4 | 30.1 | 27.4 | 27.5 | 35.5 | 39.3 | 32.5 | 34.0 |
| kan_Knda | 6.1 | 6.7 | 5.1 | 5.8 | 21.1 | 24.9 | 17.8 | 19.8 |
| kas_Arab | – | 11.3 | 6.5 | 9.4 | – | 31.8 | 23.1 | 25.2 |
| mai_Deva | – | 18.9 | 15.3 | 18.0 | – | 36.6 | 28.7 | 29.3 |
| mal_Mlym | 11.1 | 11.3 | 9.1 | 9.4 | 27.6 | 31.6 | 23.8 | 27.4 |
| mni_Mtei | – | 14.2 | 14.6 | 15.2 | – | 31.9 | 26.1 | 26.9 |
| mar_Deva | 15.5 | 19.4 | 14.7 | 16.2 | 32.2 | 36.7 | 28.5 | 32.6 |
| npi_Deva | – | 21.2 | 14.3 | 19.4 | – | 42.4 | 33.5 | 36.9 |
| ory_Orya | 11.3 | 12.3 | 11.7 | 13.9 | 33.6 | 38.8 | 30.4 | 34.1 |
| pan_Guru | 32.0 | 35.7 | 30.8 | 31.5 | 36.8 | 43.0 | 33.2 | 37.0 |
| san_Deva | – | 6.3 | 3.9 | 5.5 | – | 26.1 | 17.8 | 19.5 |
| sat_Olck | – | 6.6 | 10.9 | 10.6 | – | 23.1 | 23.7 | 25.0 |
| snd_Deva | – | 7.4 | 8.3 | 9.2 | – | 27.5 | 26.5 | 27.2 |
| tam_Taml | 7.7 | 7.6 | 7.2 | 7.2 | 20.8 | 22.7 | 18.0 | 19.7 |
| tel_Telu | 12 | 14.1 | 10.9 | 11.2 | 26.3 | 31.0 | 23.6 | 26.3 |
| urd_Arab | – | 43.7 | 33.5 | 34.6 | – | 45.9 | 35.7 | 40.0 |
| Avg. | 16.27 | 18.14 | 15.31 | 16.58 | 30.25 | 33.36 | 27.39 | 29.77 |

Table 6: BLEU (↑) scores on IN22-Conv

pan_Guru-snd_Deva lags behind by a 0.3 chrF++ score.

| Language pair | IT2 | IR Seed |
|---|---|---|
| **IN22-Gen** | | |
| ben_Beng-hin_Deva | 48.7 | 44.0 (-4.7) |
| hin_Deva-ben_Beng | 45.7 | 41.3 (-4.4) |
| hin_Deva-mal_Mlym | 44.4 | 39.2 (-5.2) |
| mal_Mlym-hin_Deva | 48.0 | 42.4 (-5.6) |
| pan_Guru-snd_Deva | 30.8 | 30.5 (-0.3) |
| snd_Deva-pan_Guru | 41.1 | 37.5 (-3.6) |
| tam_Taml-tel_Telu | 43.5 | 38.3 (-5.2) |
| tel_Telu-tam_Taml | 45.4 | 41.5 (-3.9) |
| **IN22-Conv** | | |
| ben_Beng-hin_Deva | 44.3 | 40.8 (-3.5) |
| hin_Deva-ben_Beng | 44.0 | 39.2 (-4.8) |
| hin_Deva-mal_Mlym | 40.9 | 36.8 (-4.1) |
| mal_Mlym-hin_Deva | 40.8 | 37.6 (-3.2) |
| pan_Guru-snd_Deva | 29.4 | 28.5 (-0.9) |
| snd_Deva-pan_Guru | 43.8 | 40.8 (-3.0) |
| tam_Taml-tel_Telu | 37.4 | 32.5 (-4.9) |
| tel_Telu-tam_Taml | 36.6 | 33.5 (-3.1) |

Table 7: Indic-Indic chrF++ (↑) scores results on IN22-Gen and IN22-Conv dataset

## 6  Conclusion

This paper presents our system for the WAT24 shared task on the MultiIndic22MT 2024 Shared Task. We focus on a universal model using pretrain-ing Indic languages with alignment augmentation and further obtaining direction-specific models using finetuning on small and high-quality seed data. We submit a competitive 243M parameter model covering 22 Indic languages that achieves a comparable performance with a 474M parameter model covering 11 languages.

## Limitations

The present study particularly focuses on pre-training objectives on a parallel corpus. However, techniques such as utilizing monolingual corpus (Pan et al., 2021) along with alignment objective remain unexplored. Also, large language models can be potentially leveraged to generate datasets for low-resource Indic languages. Further, we restricted the alignment augmentation of substitute words from source sentences (English words). However, words from target sentences can also be substituted can explored.

## Acknowledgements

We thank the reviewer for the valuable feedback. We are also grateful to the Department of Computer Science and Engineering, Indian Institute of Technology Hyderabad, for providing the necessary computing resources to conduct the experiments.

## References

Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George F. Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen, and Yonghui Wu. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges. *CoRR*, abs/1907.05019.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Raj Dabre and Anoop Kunchukuttan. 2024. Findings of wmt 2024's multiindic22mt shared task for machine translation of 22 indian languages. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation. *arXiv preprint*.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google's multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.

Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980.

Zehui Lin, Xiao Pan, Mingxuan Wang, Xipeng Qiu, Jiangtao Feng, Hao Zhou, and Lei Li. 2020. Pretraining multilingual neural machine translation by leveraging alignment information. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2649–2663, Online. Association for Computational Linguistics.

Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohungbe, Solomon Oluwole Akinola, Shamsuddeen Muhammad, Salomon Kabongo Kabenamualu, Salomey Osei, Freshia Sackey, Rubungo Andre Niyongabo, Ricky Macharm, Perez Ogayo, Orevaoghene Ahia, Musie Meressa Berhe, Mofetoluwa Adeyemi, Masabata Mokgesi-Selinga, Lawrence Okegbemi, Laura Martinus, Kolawole Tajudeen, Kevin Degila, Kelechi Ogueji, Kathleen Siminyu, Julia Kreutzer, Jason Webster, Jamiil Toure Ali, Jade Abbott, Iroro Orife, Ignatius Ezeani, Idris Abdulkadir Dangana, Herman Kamper, Hady Elsahar, Goodness Duru, Ghollah Kioko, Murhabazi Espoir, Elan van Biljon, Daniel Whitenack, Christopher Onyefuluchi, Chris Chinenye Emezue, Bonaventure F. P. Dossou, Blessing Sibanda, Blessing Bassey, Ayodele Olabiyi, Arshath Ramkilowan, Alp Öktem, Adewale Akinfaderin, and Abdallah Bashir. 2020. Participatory research for low-resourced machine translation: A case study in African languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2144–2160, Online. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Xiao Pan, Mingxuan Wang, Liwei Wu, and Lei Li. 2021. Contrastive learning for many-to-many multilingual neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 244–258, Online. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia,

Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Gowtham Ramesh, Sumanth Doddapaneni, Aravinth Bheemaraj, Mayank Jobanputra, Raghavan Ak, Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Divyanshu Kakwani, Navneet Kumar, et al. 2022. Samanantar: The largest publicly available parallel corpora collection for 11 indic languages. *Transactions of the Association for Computational Linguistics*, 10:145–162.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. advances in neural information processing systems. *Advances in neural information processing systems*, 30(2017).

Zhen Yang, Bojie Hu, Ambyera Han, Shen Huang, and Qi Ju. 2020. CSP:code-switching pre-training for neural machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2624–2636, Online. Association for Computational Linguistics.

# DCU ADAPT at WMT24: English to Low-resource Multi-Modal Translation Task

**Sami Ul Haq, Rudali Huidrom, Sheila Castilho**
ADAPT Centre, Dublin City University, Dublin, Ireland
{sami.haq, rudali.huidrom, sheila.castilho}@adaptcentre.ie

## Abstract

This paper presents the system description of "DCU_NMT's" submission to the WMT-WAT24 English-to-Low-Resource Multimodal Translation Task. We participated in the English-to-Hindi track, developing both text-only and multimodal neural machine translation (NMT) systems. The text-only systems were trained from scratch on constrained data and augmented with back-translated data. For the multimodal approach, we implemented a context-aware transformer model that integrates visual features as additional contextual information. Specifically, image descriptions generated by an image captioning model were encoded using BERT and concatenated with the textual input.

The results indicate that our multimodal system, trained solely on limited data, showed improvements over the text-only baseline in both the challenge and evaluation sets, suggesting the potential benefits of incorporating visual information.

## 1 Introduction

The increasing prominence of multimodal content in the machine translation (MT) community highlights its potential to improve translation quality by incorporating visual context, which is otherwise inaccessible through textual information alone. This approach has significant implications for commercial applications, including the translation of image captions in online news articles and the translation of product descriptions in e-commerce platforms (Belz et al., 2017; Calixto et al., 2017; Lala et al., 2017; Zhou et al., 2018). By integrating visual information, multimodal MT systems can achieve more accurate and contextually appropriate translations.

Despite MT achieving near-human performance for many high-resource languages, significant challenges remain, particularly for low-resource languages (Popel et al., 2020; Costa-jussà et al., 2022).

In recent years, the integration of additional modalities, such as images, into MT systems has gained prominence as a critical area of research (Sulubacak et al., 2020; Parida et al., 2021b,a). This multimodal approach seeks to address the limitations of traditional text-only MT by incorporating supplementary contextual information, thereby improving translation accuracy and expanding the applicability of MT across a broader spectrum of languages and specialised domains.

The WMT-WAT 2024 Shared Task[1] introduces the "English to Lowres Multi-Modal Translation Task," utilizing the Hindi, Bengali, Malayalam, and Hausa Visual Genome datasets. Participants are given an image, a specific rectangular region within it, and a short English caption describing the region. The task is to translate the caption into one of the target languages: Hindi, Bengali, Malayalam, or Hausa.

In this system description paper, we explain our approach for the tasks in which we participated in English (EN) to Hindi (HI) (i) Text only and (ii) Multimodal translation. We released the code and data produced during research through GitHub[2].

## 2 Dataset

We use the data sets provided by the organizers for the relevant tasks. The Visual Genome datasets for Hindi, Bengali, Malayalam, and Hausa include 29,000 training examples, 1,000 examples for development, and 1,600 examples for evaluation. These datasets are based on a shared set of images, with some variations due to independent sanity checks conducted for each language. For evaluation, the WMT-WAT 2024 Multimodal Shared Task utilises 1,600 examples from the evaluation set and 1,400 examples from the challenge

---

[1] https://www2.statmt.org/wmt24/multimodallowresmt-task.html
[2] https://github.com/sami-haq99/DCU_NMT_WMT-WAT24

set. In this submission, we denote evaluation set as "EV" and challenge set as "CH" respectively. The statistics of the dataset are shown in Table 1. Due to time constrained, We only trained our systems for English-Hindi language pair.

| Set | Sentences | Tokens | |
| | | English | Hindi |
| --- | --- | --- | --- |
| Train | 28930 | 143164 | 145448 |
| D-Test | 998 | 4922 | 4978 |
| E-Test | 1595 | 7853 | 7852 |
| C-Test | 1400 | 8186 | 8639 |

Table 1: Statistics of our data used in the English→Hindi Multimodal translation task.

## 3   Experimental Details

In this section, we present our experimental details for the tasks we participated in.

### 3.1   Text-only translation

For the EN-HI text-only translation task, we have two submissions: one restricted and the other using additional monolingual data.

Back-translation enables the effective use of monolingual data to improve the MT system, especially in a low-resource context (Sennrich et al., 2016; Ul Haq et al., 2020) where model struggles to learn reliable alignments from limited parallel data. For our experiments, we used backtranslated data generated from Flickr8k image captioning data set to enrich text-only data (Parida et al., 2022). For our text-only baseline, we trained the sentence-level transformer model from scratch using all training data until convergence.

### 3.2   Multimodal translation

For the EN-HI multimodal translation system, we employ a context-aware model, an extension of the Transformer architecture designed to incorporate additional contextual information during translation. Unlike traditional neural machine translation (NMT) models that translate sentences independently, context-aware NMT relaxes this assumption by conditioning the translation not only on the current source sentence but also on auxiliary information from within or outside the document. Given that the HVG data set is limited to the caption translation of specific image regions, we hypothesize that providing the model with additional context, such as a comprehensive description of the entire image, could enhance the accuracy of the generated translations. To take advantage of visual features, we extracted image captions from HVG image dataset and used them as additional context for translation. Additionally, we used pretrained BERT as additional encoder to encode and aggregate contextual features (Wu et al., 2022).

We used BLIP[3], an image caption model, to generate a description of the HVG image dataset. As HVG contains short descriptions of specified regions of images in English and Hindi, we generate captions of entire image to be fed as additional information to multimodal context-aware model. Since our context-aware model expects context during the training and evaluation stage, we generated captions for the entire HVG dataset, including the evaluation (EV) and Challenge (CH) test sets. The overview of our multimodal translation system is depicted in Figure 1.

Two step training strategy is followed, we first train a strong sentence-level transformer model using all the training data until convergence, then the context-aware model is initialized from best checkpoint and fine-tuned on context-aware data. We select the best model on the validation data. Contextual features are encoded using pre-trained *bert-base* model released under transformers package [4].The model incorporates two special tokens: [CLS], which is added at the beginning of a sentence, and [SEP], which is employed to separate different sequences. The context is concatenated with sentences as follows:

$x_{ctx} =$ *[CLS] surfer on a surfboard riding a wave in the ocean [SEP] man surfing in ocean [SEP]*

Several techniques exist for context integration (Castilho et al., 2020; Wu et al., 2022; Haq et al., 2022), we used $1 - fixed - sequence$ on the source and target side as context. In this approach, a single previous sentence or external sequence is considered context for current sentence being translated. After that Bert encoded features are extracted as defined in equation 1. Although the context-aware multi-encoder models are exposed to additional contextual information, the translation is still performed at sentence level.

$$C = BERT(x_{ctx}) \tag{1}$$

---

[3] https://huggingface.co/Salesforce/blip-image-captioning-large

[4] https://github.com/huggingface/transformers

Figure 1: Overview of multimodal translation system.

Our translation models are based on transformer architecture with 6 encoder/decoder blocks, 512 embedding input, and 1024 FFN layer dimension size. Dropout rate is 0.3 for all tasks. We use the Adam optimizer and $5 \times 10^{-4}$ learning rate schedule with 4000 warmup steps. Model training was conducted on two GPUs, with a batch size of 6000 tokens per GPU. Our Transformer implementation is based on the Fairseq (Ott et al., 2019) toolkit.

## 4 Results

Our results for EN-HI text-only and multimodal translation are presented in Table 2.

| Modality | System | BLEU | |
| | | EV | CH |
|---|---|---|---|
| text-only | Transformer | 40.20 | 29.20 |
| text-only | Transformer$_{bt}$ | **42.70** | **35.90** |
| multimodal | Context-aware$_{src\_tgt}$ | 40.60 | 28.60 |
| multimodal | Context-aware$_{src}$ | 40.60 | **30.30** |

Table 2: WMT_WAT2024 Automatic evaluation results for EN→HI on Evaluation (EV) and Challenge (CH) test sets. "Transformer$_{bt}$" denotes NMT model trained with back-translated data. For multimodal task, "$src\_tgt$" represents context-aware model with visual contextual features used on both encoder and decoder side while $src$ indicates context used only on the encoder side.

For text-only translation, the baseline system (Transformer) obtains BLEU scores of 40.20 on the evaluation set (EV) and 29.20 on the challenge test (CH). In contrast, the Transformer$_{bt}$ (Transformer with back-translated data) system demonstrates improved results, with BLEU scores of 42.70 for EV and 35.90 for CH. This improvement suggests that back-translation enhances translation quality by incorporating additional synthetic data, which is particularly advantageous for the challenge set (CH).

In multimodal translation, the context-aware$_{src\_tgt}$ approach achieves BLEU scores of 40.60 for EV and 28.60 for CH. Compared with a text-only restricted baseline, the EV score slightly exceeds that of the Transformer (40.20), the CH score is lower, indicating that while the multimodal context benefits the evaluation set, it does not consistently improve performance on the challenge set. Conversely, the context-aware$_{src}$ (source only context) method achieves BLEU scores of 40.60 for EV and 30.30 for CH, showing a modest improvement for the challenge set compared to the $src\_tgt$ and text-only methods (except Transformer$_{bt}$).

## 5 Discussion

The baseline Transformer model achieves BLEU scores of 40.20 for the evaluation set (EV) and 29.20 for the challenge test (CH). The Transformer$_{bt}$ model shows marked improvement, with BLEU scores of 42.70 for EV and 35.90 for CH, highlighting back-translation's effectiveness in enhancing performance, particularly in challenging scenarios.

In multimodal translation, the context-aware$_{src\_tgt}$ method, which utilises visual context on both the encoder and decoder sides, scores 40.60 for EV and 28.60 for CH. It slightly outperforms the baseline Transformer on EV but underperforms on CH, suggesting that while visual context can help in simpler cases, it may complicate results in more difficult scenarios.

The context-aware $src$ only approach, using visual context only with the source text, achieves BLEU scores of 40.60 for EV and 30.30 for CH. It shows modest improvement over $src\_tgt$ for CH but does not surpass the Transformer + Back-translation in overall performance. This is obvious because multimodal translation systems are trained on constrained resources while Transfomer$_{bt}$ use 8k additional synthetic parallel sentences for training.

These findings underscore the value of back-translation in improving text-only translation, especially for more challenging tasks. For multimodal translation, while visual context can be beneficial, its effectiveness varies with context integration choice. The results suggest that different methods may be better suited to different types of translation challenge, indicating a need for further research to optimize the use of visual context.

# 6 Conclusion

Our results have showed that the Transformer with back-translated data consistently outperforms the text-only and multimodal systems in both evaluation tasks, demonstrating a significant benefit of back-translation, particularly for challenging scenarios. Our multimodal systems, despite not utilizing back-translated data, still outperformed the text-only baseline, highlighting the potential of visual context in improving translation accuracy. However, multimodal systems employing visual context on both the encoder and decoder sides do not exhibit a clear advantage over the text-only model or other multimodal approaches. Notably, the multimodal method shows diminished effectiveness for the challenge test (CH), suggesting that while additional visual context may enhance performance in certain cases, it can also introduce complexities that potentially undermine translation accuracy. These findings highlight the need for further investigation into optimizing the integration of visual context to improve translation outcomes across varying task difficulties.

# Acknowledgements

# References

Anja Belz, Erkut Erdem, Katerina Pastra, and Krystian Mikolajczyk. 2017. Proceedings of the sixth workshop on vision and language. In *Proceedings of the Sixth Workshop on Vision and Language*.

Iacer Calixto, Daniel Stein, Evgeny Matusov, Pintu Lohar, Sheila Castilho, and Andy Way. 2017. Using Images to Improve Machine-Translating E-Commerce Product Listings. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 637–643, Valencia, Spain. Association for Computational Linguistics.

Sheila Castilho, Maja Popović, and Andy Way. 2020. On context span needed for machine translation evaluation.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Sami Ul Haq, Sadaf Abdul Rauf, Arslan Shaukat, and Muhammad Hassan Arif. 2022. Context-aware neural machine translation using selected context. In *2022 19th International Bhurban Conference on Applied Sciences and Technology (IBCAST)*, pages 349–352. IEEE.

Chiraag Lala, Pranava Madhyastha, JK Wang, and Lucia Specia. 2017. Unraveling the contribution of image captioning and neural machine translation for multimodal machine translation. In *The Prague Bulletin of Mathematical Linguistics*, volume 108-1, pages 197–208. De Gruyter Open.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.

Shantipriya Parida, Subhadarshi Panda, Satya Prakash Biswal, Ketan Kotwal, Arghyadeep Sen, Satya Ranjan Dash, and Petr Motlicek. 2021a. Multimodal neural machine translation system for english to bengali. In *Proceedings of the First Workshop on Multimodal Machine Translation for Low Resource Languages (MMTLRL 2021)*, pages 31–39.

Shantipriya Parida, Subhadarshi Panda, Stig-Arne Grönroos, Mark Granroth-Wilding, and Mika Koistinen. 2022. Silo NLP's participation at WAT2022. In *Proceedings of the 9th Workshop on Asian Translation*, pages 99–105, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.

Shantipriya Parida, Subhadarshi Panda, Ketan Kotwal, Amulya Ratna Dash, Satya Ranjan Dash, Yashvardhan Sharma, Petr Motlicek, and Ondřej Bojar. 2021b. Nlphut's participation at wat2021. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 146–154.

Martin Popel, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature communications*, 11(1):1–15.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Umut Sulubacak, Ozan Caglayan, Stig-Arne Grönroos, Aku Rouhe, Desmond Elliott, Lucia Specia, and Jörg Tiedemann. 2020. Multimodal machine translation through visuals and speech. *Machine Translation*, 34:97–147.

Sami Ul Haq, Sadaf Abdul Rauf, Arsalan Shaukat, and Abdullah Saeed. 2020. Document level NMT of low-resource languages with backtranslation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 442–446, Online. Association for Computational Linguistics.

Xueqing Wu, Yingce Xia, Jinhua Zhu, Lijun Wu, Shufang Xie, and Tao Qin. 2022. A study of bert for context-aware neural machine translation. *Machine Learning*, pages 1–19.

Mingyang Zhou, Runxiang Cheng, Yong Jae Lee, and Zhou Yu. 2018. A visual attention grounding neural model for multimodal machine translation. *arXiv preprint arXiv:1808.08266*.

# English-to-Low-Resource Translation: A Multimodal Approach for Hindi, Malayalam, Bengali, and Hausa

**Ali Hatami**[1,4]  and  **Shubhanker Banerjee**[2,4]  and  **Mihael Arcan**[3]  and
**Bharathi Raja Chakravarthi**[1,4]  and  **Paul Buitelaar**[1,4]  and  **John Philip McCrae**[1,2,4]

[1]Insight SFI Research Centre for Data Analytics, Data Science Institute, University of Galway
[2]ADAPT Centre, Data Science Institute, University of Galway
[3]Lua Health, Galway, Ireland
[4]School of Computer Science, University of Galway
ali.hatami@insight-centre.org

## Abstract

Multimodal machine translation leverages multiple data modalities to enhance translation quality, particularly for low-resourced languages. This paper uses a multimodal model that integrates visual information with textual data to improve translation accuracy from English to Hindi, Malayalam, Bengali, and Hausa. This approach employs a gated fusion mechanism to effectively combine the outputs of textual and visual encoders, enabling more nuanced translations that consider both language and contextual visual cues. The model's performance was evaluated against the text-only machine translation model based on BLEU, ChrF2 and TER. Experimental results demonstrate that the multimodal approach consistently outperforms the text-only baseline, highlighting the potential of integrating visual information in low-resourced language translation tasks.

## 1 Introduction

In recent years, neural network-based translation models have been widely used in translation tasks, demonstrating remarkable performance in terms of fluency and precision compared to previous generations of machine translation systems (Cho et al., 2014). The Transformer model, in particular, has shown significant improvements in machine translation tasks. A crucial component of the Transformer model is the cross-attention mechanism, which enhances the model's ability to capture semantic dependencies by combining self-attention—allowing source words to interact with one another—with attention mechanisms that involve target words (Vaswani et al., 2017).

Despite the broader context focus in text-only translation models, understanding the input text remains a challenge. In natural language, lexical ambiguity (Rios Gonzales et al., 2017) occurs when a single word has multiple meanings or interpretations, complicating text comprehension. For example, in the domain of finance and economics,

the word "bank" almost always refers to a financial institution rather than the side of a river.

Multimodal Machine Translation (MMT), a sub-area of NMT, has been introduced to utilise visual information from other modalities, such as images, to translate an aligned sentence in a source language into a target language. Recent studies (Yao and Wan, 2020; Zhao et al., 2022; Wang and Xiong, 2021) demonstrate the potential of leveraging multimodal information, alongside textual content, to enhance translation quality. Visual cues, as an additional source of information, can provide valuable insights that complement textual information, enabling MMT models to better understand and produce more accurate and contextually appropriate translations. The concept behind MMT is to integrate visual information to help disambiguate input words, detect the correct scenes in the source language, and select the appropriate translation in the target language (Hatami et al., 2022). MMT is particularly beneficial when dealing with low-resource languages where there is not sufficient parallel data to train the model.

This paper aims to explore the benefit of using visual information in translating English into four different low-resource languages, Hindi, Malayalam, Bengali and Hausa. We used a gated fusion approach to integrate textual and visual information in the encoder and generate the text in the target language on the decoder side. In the baseline, we train the model on the input text without considering the aligned image. For the multimodal model, we trained four different models for each language. We explain our methodology in Section 3, our experimental setup in Section 4, results in Section 5, and we conclude our findings in Section 6.

## 2 Related Work

There are various approaches proposed to integrate visual information with text-only translation models. These approaches typically utilise a visual

attention mechanism in either the decoder or encoder to capture the relationships between words in a sentence and image features. The common method involves extracting visual information by employing Convolutional Neural Networks (CNN) and then integrating this information with textual features.

Regarding visual features, existing studies on MMT employ two types of visual features: global and local visual features. Global features represent the entire image as a single vector without attention to the spatial layout of the image. On the other hand, local features describe an image as a sequence of equally sized patches (Calixto et al., 2017). Local features are extracted from multiple points in the image and are more robust to clutter than global features (Lisin et al., 2005). CNNs can be used to extract both global and local features from the image (Zheng et al., 2019).

Global image features are used in the encoder in addition to word sequences (Huang et al., 2016). Alternatively, they can be used to initialise the hidden parameters of the encoder and decoder of a RNN (Calixto and Liu, 2017). Element-wise multiplication was used to initialise the hidden states of the encoder/decoder in the attention-based model (Caglayan et al., 2017). Visual attention mechanism was employed to link visual and corresponding text semantically (Zhou et al., 2018).

Several approaches have been proposed to improve the quality of the visual modality in Multimodal Machine Translation (MMT). For instance, a multimodal Transformer-based self-attention mechanism was introduced to encode relevant information in images (Yao and Wan, 2020). A graph-based multimodal fusion encoder was developed to capture various relationships between modalities (Yin et al., 2020). Additionally, a translate-and-refine mechanism was implemented using images in a second-stage decoder to refine a text-only Neural Machine Translation (NMT) model, particularly for handling ambiguous words. A latent variable model was also employed to extract the multimodal relationships between image and text modalities (Calixto et al., 2019).

Recent methods aim to reduce noise in visual information and select visual features relevant to the text. For example, object-level visual modelling has been used to mask irrelevant objects and specific words in the source text to enhance visual feature learning (Wang and Xiong, 2021).

Object detection in the image encoder has been employed to extract visual features from object regions within an image, which are then applied to a doubly-attentive decoder model (Zhao et al., 2022).

In this paper, we adopt the gated fusion MMT model (Wu et al., 2021), which integrates visual and textual representations through a gate mechanism. This gated fusion mechanism allows the model to adjust the amount of visual information that contributes to the translation process.

## 3 Methodology

The objective of our experiments is to evaluate the impact of visual features on translation quality in low-resource languages. Following Wu et al. (2021), we conduct experiments to assess both the text-only Transformer and the gated fusion multimodal Transformer (gated fusion MMT) using the shared task data for Hindi, Bengali, Malayalam, and Hausa. In this section, we provide descriptions of the model architectures mentioned above.

### 3.1 Text-only Machine Translation

For the text-only translation model, we use the training and development sets for Hindi, Bengali, Malayalam, and Hausa to train the Transformer-based model. This model serves as our baseline for evaluating the multimodal model. The text-only Transformer architecture was introduced by Vaswani et al. (2017). It consists of an encoder-decoder structure, where both the encoder and decoder are composed of stacked layers of self-attention, and feed-forward neural networks.

First, we tokenize the sentences into subwords in the training, development, and test sets. We then train four translation models on the tokenized sentences for these language pairs. Tokenization helps the model better learn the language and handle out-of-vocabulary words, especially in low-resource languages. During the inference step, we translate the tokenized test sentences from English into the four low-resource languages.

### 3.2 Multimodal Machine Translation

For the multimodal model, we use the gated fusion approach (Wu et al., 2021) to fuse both textual and visual information. Gated fusion MMT incorporates visual information into the translation process in a controlled and interpretable manner using a gating mechanism. The textual component is similar to the text-only model, with tokenized sentences

| Dataset | Hindi | Bengali | Malayalam | Hausa |
|---|---|---|---|---|
| Training Set | 28,932 | 28,930 | 29,000 | 28,930 |
| Development Set (D-Test) | 998 | 998 | 1,000 | 998 |
| Evaluation Set (E-Test) | 1,595 | 1,595 | 1,600 | 1,595 |
| Challenge Test Set (C-Test) | 1,400 | 1,400 | 1,400 | 1,400 |
| **Total** | 32,925 | 32,923 | 33,000 | 32,923 |

Table 1: Number of sentences of Visual Genome dataset for Hindi, Bengali, Malayalam and Hausa.

fed into the model. On the visual side, each sentence is paired with an image, and for each image, we have the coordinates of the rectangular region corresponding to the part of the image that relates to the sentence (see Figure 1).

For each language, we trained two models: one that considers the entire image and another that considers only the specific rectangular region. We use the pre-trained ResNet-101 CNN (He et al., 2016) to extract visual features from the images. In this study, we extract visual representations from both the whole image and the designated rectangular region, which is aligned with the text caption. The motivation for using the partial image (rather than the full image) is that objects outside the rectangular region may be irrelevant to the text caption and could potentially degrade translation model performance (Hatami et al., 2023).

Both the textual and visual representations are fed into the gated fusion model, allowing it to be trained based on both modalities. We then use these multimodal models to translate test sentences that are aligned with images. More detailed information about the multimodal models can be found in Section 4.2.2.

## 4 Experimental Setup

### 4.1 Dataset

The Hindi Visual Genome (HVG) (Parida et al., 2019), Bengali Visual Genome (BVG) (Sen et al., 2022), Malayalam Visual Genome (MVG) (Parida et al., 2019), and Hausa Visual Genome (HaVG) (Abdulmumin et al., 2022) datasets are multimodal datasets designed for English-to-Hindi, English-to-Bengali, English-to-Malayalam, and English-to-Hausa machine translation, respectively (Figure 1). These datasets, based on the original Visual Genome dataset, contain real-world images annotated with region-specific captions. The captions have been translated into the respective languages through a combination of automated translation and manual post-editing by native speakers to ensure

contextual accuracy.

The MVG, HVG, BVG, and HaVG datasets are divided into training, development, evaluation, and challenge test sets, as outlined in Table 1.

**Training Set**: The training sets for Malayalam, Hindi, Bengali, and Hausa contain 29,000, 28,932, 28,930, and 28,930 image-caption pairs, respectively. Each pair consists of an image, a selected region in the image, and its corresponding English and Malayalam/Hindi/Bengali/Hausa captions. The captions have been manually refined to align with the visual context of the images.

**Development Set (D-Test)**: The development sets contain 1,000 image-caption pairs in the Malayalam dataset and 998 pairs in the Hindi, Bengali, and Hausa datasets. These sets are used to validate and fine-tune model performance during the training process.

**Evaluation Set (E-Test)**: The evaluation sets include 1,600 image-caption pairs in the Malayalam dataset and 1,595 pairs in the Hindi, Bengali, and Hausa datasets. These sets are used for evaluating model performance on unseen data, providing a benchmark for generalization capabilities.

**Challenge Test Set (C-Test)**: The challenge test sets for all four languages consist of 1,400 image-caption pairs. These sets are designed to focus on ambiguous English words that require visual context to resolve their meaning in Malayalam, Hindi, Bengali, or Hausa. The ambiguous words were identified based on embedding similarity, and the corresponding images help disambiguate their meaning, providing a robust test for multimodal translation systems (Hatami et al., 2024).

### 4.2 Machine Translation Models

#### 4.2.1 Text-only Translation Model

A text-only Transformer model serves as the baseline in our experiment, utilizing only the textual captions of images for translation. The model is trained using the OpenNMT toolkit (Klein et al., 2018) on the Visual Genome dataset for English-

Figure 1: Examples from the Visual Genome dataset show English caption of the rectangular region (solid red line) with translation in Hindi, Bengali, Malayalam and Hausa.

to-Hindi, Bengali, Malayalam, and Hausa translations. It comprises a 6-layer Transformer architecture with attention mechanisms in both the encoder and decoder stages, trained for 50k steps.

The encoder processes a sequence of tokens (words or subword units) and generates context-aware representations for each token. The decoder generates the output sequence (e.g., translated text) by leveraging the encoded representations from the encoder along with the previously generated tokens. It employs multi-head self-attention and feed-forward layers, incorporating additional attention mechanisms to effectively focus on the encoded input. The core innovation of the Transformer is the **self-attention mechanism**, which computes attention scores across all tokens in the sequence, creating weighted representations that capture contextual relationships between tokens.

Since the Transformer model does not inherently process sequences in a fixed order, as recurrent neural networks (RNNs) do, it uses **positional encodings** to inject information about the position of tokens in the sequence. These positional encodings are added to the input embeddings, enabling the model to differentiate between tokens based on their positions within the sequence. To enhance its ability to capture different types of relationships between tokens, the Transformer employs **multi-head attention**. This involves splitting the self-attention process into multiple parallel attention heads, each learning a different set of attention weights. The outputs from all heads are then con-catenated and linearly transformed to provide a richer, more comprehensive representation of the input sequence.

SentencePiece (Kudo and Richardson, 2018) is employed to segment words into subword units, offering a language-independent approach to tokenization without requiring pre-processing steps, thereby enhancing the model's adaptability and versatility in handling raw text.

### 4.2.2 Multimodal Machine Translation

In the MMT model, we adopt the gated fusion MMT model (Wu et al., 2021), which fuses visual and text representations by employing a gate mechanism. Gated fusion is a mechanism used to integrate visual information from images with textual information from source sentences during the translation process. The main idea behind gated fusion is to control the amount of visual information that is blended into the textual representation using a gating matrix.

The source sentence $x$ is fed into a vanilla Transformer encoder to obtain a textual representation $H_{text}$ of dimension $T{\times}d$ [1]. The image $z$ is processed using a pre-trained ResNet-101 CNN (He et al., 2016), which has been trained on the ImageNet dataset (Russakovsky et al., 2014), to extract a 2048-dimensional average-pooled visual representation, denoted as $Embed_{image}(z)$. The visual representation $Embed_{image}(z)$ is projected to the

---

[1]T is the number of tokens (words) in the input sentence, and d is the dimensionality of the representation

| English → Hindi | BLEU ↑ | ChrF2 ↑ | TER ↓ |
|---|---|---|---|
| Text-only MT | 38.26 | 58.65 | 42.54 |
| Multimodal MT (entire image) | **39.65**\* | **59.34**\* | **41.92**\* |
| Multimodal MT (partial image) | 38.64 | 58.84 | 42.62 |

| English → Bengali | BLEU ↑ | ChrF2 ↑ | TER ↓ |
|---|---|---|---|
| Text-only MT | 39.85 | 64.32 | 39.24 |
| Multimodal MT (entire image) | **41.92**\* | **65.96**\* | **38.37**\* |
| Multimodal MT (partial image) | 39.45 | 64.75 | 39.65 |

| English → Malayalam | BLEU ↑ | ChrF2 ↑ | TER ↓ |
|---|---|---|---|
| Text-only MT | 28.94 | 58.74 | 54.87 |
| Multimodal MT (entire image) | **32.34**\* | **61.15**\* | **53.94**\* |
| Multimodal MT (partial image) | 28.76 | 58.63 | 54.58 |

| English → Hausa | BLEU ↑ | ChrF2 ↑ | TER ↓ |
|---|---|---|---|
| Text-only MT | 39.86 | 61.21 | 47.59 |
| Multimodal MT (entire image) | **41.25**\* | **62.94**\* | **46.48**\* |
| Multimodal MT (partial image) | 38.31 | 60.87 | 47.62 |

Table 2: BLEU, ChrF2 and TER scores for text-only and multimodal models for English to Hindi, Bengali, Malayalam and Hausa on the test set (\* represents a statistically significant result compared to the baseline text-only model at a significance level of $p < 0.05$).

same dimension as $H_{text}$ using a weight matrix $W_z$, denoted as:

$$\text{Embed}_{\text{image}}(z) = W_z \text{ResNet}_{\text{pool}}(z)$$

where $W_z$ is a learned projection matrix.

To determine the amount of visual information to fuse with the textual representation, a gating matrix $\Lambda$ of dimension $T \times d$ is generated ($[0, 1]^{T \times d}$). This matrix is computed using a sigmoid function applied to both the projected visual representation and the textual representation:

$$\Lambda = \sigma\left(W_\Lambda \text{Embed}_{\text{image}}(z) + U_\Lambda H_{\text{text}}\right)$$

where $W_\Lambda$ and $U_\Lambda$ are learned parameters, and $\sigma$ is the sigmoid function. The gating matrix $\Lambda$ makes the fusion process interpretable, as it controls how much visual context is used in translation. A larger value in $\Lambda$ indicates that the model is relying more on the visual context, while a smaller value indicates a stronger reliance on the textual representation alone.

The final representation $H$ that combines both textual and visual information is given by:

$$H = H_{\text{text}} + \Lambda \text{Embed}_{\text{image}}(z)$$

This fused representation $H$ is then passed into the Transformer decoder for generating the target translation.

### 4.3 Evaluation Metrics

We use three evaluation metrics: BLEU (Papineni et al., 2002), ChrF2 (Popović, 2015), and TER (Snover et al., 2006). BLEU assesses the precision of translation by comparing candidate translations to reference translations based on *n-grams*. ChrF2 evaluates the similarity between character *n-grams* in machine-generated and reference translations, particularly beneficial for languages with complex writing systems. TER quantifies the number of edits needed to align machine translations with human-generated references. We conduct statistical significance testing using the *sacreBLEU*[2] toolbox.

---

[2]https://github.com/mjpost/sacrebleu

## 5 Results and Discussion

In this section, we present the results of our experiments, where we trained our models on the Visual Genome dataset and evaluated the translation quality using the BLEU, ChrF2, and TER metrics. We compare the translation quality of our proposed models with text-only baseline models, where the text-only NMT model was trained solely on text captions without images, across test sets for four languages.The MMT models were trained on both text captions and original images with entire images and just considering the coordinates of a part of the image related to the caption (partial image).

The results in Table 2 demonstrate the performance of both text-only and multimodal models across four language pairs: English to Hindi, Bengali, Malayalam, and Hausa. For English to Hindi, the MMT model that utilizes the entire image outperforms the text-only model, achieving a BLEU score of 39.65, ChrF2 score of 59.34, and TER score of 41.92. These improvements are statistically significant over the text-only MT model at $p < 0.05$, highlighting the benefit of incorporating visual context into the translation process. Similar trends are observed for English to Bengali, where the entire image-based MMT achieves a BLEU score of 41.92, a ChrF2 score of 65.96, and a TER score of 38.37, all of which are significantly better than the text-only model.

For English to Malayalam, the entire image-based multimodal model also shows clear advantages, with a BLEU score of 32.34, ChrF2 of 61.15, and TER of 53.94, outperforming the text-only model on all metrics. Finally, in the case of English to Hausa, the entire image-based multimodal MT model again demonstrates superior performance, achieving a BLEU score of 41.25, ChrF2 of 62.94, and TER of 46.48, compared to the text-only model. Across all language pairs, the partial image-based multimodal models do not consistently outperform the text-only models, suggesting that complete visual context is necessary for achieving the best translation quality.

## 6 Conclusion

This paper demonstrates the significant advantages of employing a multimodal machine translation approach that integrates visual information with textual data, especially in the case of low-resourced languages like Hindi, Malayalam, Bengali, and Hausa. The results indicate that the gated fusion MMT model enhances translation accuracy and provides a more nuanced understanding of context, leading to improved performance over traditional text-only models. By leveraging visual context, we can address the challenges faced in translating low-resourced languages, highlighting the importance of incorporating diverse data modalities to enrich the translation process.

## Acknowledgements

## References

Idris Abdulmumin, Satya Ranjan Dash, Musa Abdullahi Dawud, Shantipriya Parida, Shamsuddeen Muhammad, Ibrahim Sa'id Ahmad, Subhadarshi Panda, Ondřej Bojar, Bashir Shehu Galadanci, and Bello Shehu Bello. 2022. Hausa visual genome: A dataset for multi-modal English to Hausa machine translation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6471–6479, Marseille, France. European Language Resources Association.

Ozan Caglayan, Walid Aransa, Adrien Bardet, Mercedes García-Martínez, Fethi Bougares, Loïc Barrault, Marc Masana, Luis Herranz, and Joost van de Weijer. 2017. LIUM-CVC submissions for WMT17 multimodal translation task. In *Proceedings of the Second Conference on Machine Translation*, pages 432–439, Copenhagen, Denmark. Association for Computational Linguistics.

Iacer Calixto and Qun Liu. 2017. Incorporating global visual features into attention-based neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 992–1003, Copenhagen, Denmark. Association for Computational Linguistics.

Iacer Calixto, Qun Liu, and Nick Campbell. 2017. Doubly-attentive decoder for multi-modal neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1913–1924, Vancouver, Canada. Association for Computational Linguistics.

Iacer Calixto, Miguel Rios, and Wilker Aziz. 2019. Latent variable model for multi-modal translation. In

*Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6392–6405, Florence, Italy. Association for Computational Linguistics.

Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 103–111, Doha, Qatar. Association for Computational Linguistics.

Ali Hatami, Mihael Arcan, and Paul Buitelaar. 2024. Enhancing translation quality by leveraging semantic diversity in multimodal machine translation. In *Proceedings of the 16th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 154–166, Chicago, USA. Association for Machine Translation in the Americas.

Ali Hatami, Paul Buitelaar, and Mihael Arcan. 2022. Analysing the correlation between lexical ambiguity and translation quality in a multimodal setting using WordNet. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop*, pages 89–95, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.

Ali Hatami, Paul Buitelaar, and Mihael Arcan. 2023. A filtering approach to object region detection in multimodal machine translation. In *Proceedings of Machine Translation Summit XIX, Vol. 1: Research Track*, pages 393–405, Macau SAR, China. Asia-Pacific Association for Machine Translation.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Po-Yao Huang, Frederick Liu, Sz-Rung Shiang, Jean Oh, and Chris Dyer. 2016. Attention-based multimodal neural machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 639–645, Berlin, Germany. Association for Computational Linguistics.

Guillaume Klein, Yoon Kim, Yuntian Deng, Vincent Nguyen, Jean Senellart, and Alexander Rush. 2018. OpenNMT: Neural machine translation toolkit. In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 177–184, Boston, MA. Association for Machine Translation in the Americas.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

D.A. Lisin, M.A. Mattar, M.B. Blaschko, E.G. Learned-Miller, and M.C. Benfield. 2005. Combining local and global image features for object class recognition. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05) - Workshops*, pages 47–47.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Shantipriya Parida, Ondřej Bojar, and Satya Ranjan Dash. 2019. Hindi Visual Genome: A Dataset for Multi-modal English to Hindi Machine Translation. *Computación y Sistemas*, 23(4):1499–1505. Presented at CICLing 2019, La Rochelle, France.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Annette Rios Gonzales, Laura Mascarell, and Rico Sennrich. 2017. Improving word sense disambiguation in neural machine translation with sense embeddings. In *Proceedings of the Second Conference on Machine Translation*, pages 11–19, Copenhagen, Denmark. Association for Computational Linguistics.

Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. 2014. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575.

Arghyadeep Sen, Shantipriya Parida, Ketan Kotwal, Subhadarshi Panda, Ondřej Bojar, and Satya Ranjan Dash. 2022. Bengali Visual Genome: A Multimodal Dataset for Machine Translation and Image Captioning. In *Intelligent Data Engineering and Analytics*, pages 63–70. Springer.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Dexin Wang and Deyi Xiong. 2021. Efficient object-level visual context modeling for multimodal machine translation: Masking irrelevant objects helps grounding. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 2720–2728. AAAI Press.

Zhiyong Wu, Lingpeng Kong, Wei Bi, Xiang Li, and Ben Kao. 2021. Good for misconceived reasons: An empirical revisiting on the need for visual context in multimodal machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 6153–6166. Association for Computational Linguistics.

Shaowei Yao and Xiaojun Wan. 2020. Multimodal transformer for multimodal machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4346–4350, Online. Association for Computational Linguistics.

Yongjing Yin, Fandong Meng, Jinsong Su, Chulun Zhou, Zhengyuan Yang, Jie Zhou, and Jiebo Luo. 2020. A novel graph-based multi-modal fusion encoder for neural machine translation. In *Annual Meeting of the Association for Computational Linguistics*.

Yuting Zhao, Mamoru Komachi, Tomoyuki Kajiwara, and Chenhui Chu. 2022. Region-attentive multimodal neural machine translation. *Neurocomputing*, 476:1–13.

Yufeng Zheng, Jun Huang, Tianwen Chen, Yang Ou, and Wu Zhou. 2019. CNN classification based on global and local features. In *Real-Time Image Processing and Deep Learning 2019*, volume 10996, page 109960G. International Society for Optics and Photonics, SPIE.

Mingyang Zhou, Runxiang Cheng, Yong Jae Lee, and Zhou Yu. 2018. A visual attention grounding neural model for multimodal machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3643–3653, Brussels, Belgium. Association for Computational Linguistics.

# OdiaGenAI's Participation in WAT2024 English-to-Low Resource Multimodal Translation Task

Shantipriya Parida[1], Shashikanta Sahoo[2], Sambit Sekhar[3],
Upendra Kumar Jena[4], Sushovan Jena[5], Kusum Lata[6]

[1]Silo AI, Finland; [2]Government College of Engineering Kalahandi, India;
[3]Odia Generative AI, India; [4]Creanovation Technologies Pvt. Ltd., India;
[5]IIT Mandi, India; [6]Sharda University, India
correspondence: shantipriya.parida@silo.ai

## Abstract

This paper covers the system description of the team "ODIAGEN's" submission to the 11th Workshop on Asian Translation (WAT 2024). We participated in the English-to-LowRes Multimodal Translation Task, in two of the tasks, i.e. Text-only Translation and Multi-modal Translation. For Text-only Translation, we trained the `Mistral-7B` model for English to Multi-lingual (Hindi, Bengali, Malayalam, Hausa). For Multi-modal Translation (using both image and text), we trained the `PaliGemma-3B` model for English to Hindi translation.

## 1 Introduction

Machine translation (MT) is a well-established area within Natural Language Processing (NLP), focusing on the development of software that can automatically translate text or speech between languages. While substantial progress has been made in achieving human-level translation for high-resource languages, significant challenges persist for low-resource languages (Popel et al., 2020) (Parida et al., 2023). Recent research has also investigated how to effectively incorporate other modalities, such as images, into the translation process.

Since 2013, the WAT (Workshop on Asian Translation) has been an open evaluation campaign centered on Asian languages (Nakazawa et al., 2021). The multimodal translation tasks in WAT2024 involve image caption translation, where the input includes a descriptive caption in the source language paired with the image it describes, and the output is a caption in the target language. This multimodal input leverages image context to clarify source words with multiple meanings.

The evaluation of these translation tasks is conducted using established metrics such as Bilingual Evaluation Understudy (BLEU) (Papineni et al., 2002) and Rank-based Intuitive Bilingual Evaluation Score (RIBES) (Isozaki et al., 2010). In this system description paper, we (team "ODIAGEN") outline our approach to the tasks and sub-tasks in which we participated.

- Task 1: English→Hindi (EN-HI) Multimodal Translation
  - EN-HI text-only translation
  - EN-HI multimodal translation
- Task 2: English→Malayalam (EN-ML) Text-only Translation
- Task 3: English→Bengali (EN-BN) Text-only Translation
- Task 4: English→Hausa (EN-HA) Text-only Translation

## 2 Datasets

We used only Hindi (Parida et al., 2019), Bengali (Sen et al., 2022), Malayala, and Hausa (Abdulmumin et al., 2022) Visual Genome datasets specified by the organizer for text-only and multi-modal translation without any additional synthetic data.

### 2.1 Pre-processing

#### 2.1.1 For Text-only

A few Hindi samples were excluded due to identical Hindi and English text in the Hindi dataset, and one Malayalam sample was removed for similar reasons. Formatting issues in the Hindi dataset were corrected, and duplicate samples were excluded from all language datasets. Image metadata (image_id, X, Y, Width, Height) was excluded from the text-to-text translation task. The final dataset sentence/sample count is provided in Table 1.

All four different language datasets were combined to make a common translation

| Language | Sentence Count |
|----------|----------------|
| Hindi | 28,927 |
| Bengali | 28,927 |
| Malayalam | 28,922 |
| Hausa | 28,927 |

Table 1: Training Dataset Sentence Count

dataset with a single task of translating from English to instructed Target Language like Hindi, Bengali, Malayalam, and Hausa.



Figure 1: Instruction set in different language

### 2.1.2 For Multimodal involving both Text and Image

The multimodal dataset comprises both text and images. The text portions of the dataset (train and test sets) are organized in simple tab-delimited plain text files. Each text file contains seven columns as follows:

- Column 1: imageid,
- Column 2: X,
- Column 3: Y,
- Column 4: Width,
- Column 5: Height,
- Column 6: English text,
- Column 7: Hindi Text.

The X, Y, Width, and Height columns define the rectangular region in the image described by the caption.

The `Mistral-7B` model (Beyer et al., 2024) requires data in the format $[xmin, ymin, xmax, ymax]$. We interpreted the provided X and Y coordinates as the center coordinates of the rectangular region and calculated $[xmin, ymin, xmax, ymax]$ as the coordinates of the bottom-left and top-right corners of the rectangular box.

### 2.2 Instruction Dataset

#### 2.2.1 For Text-only

Alpaca prompt format was used to prepare instruction data sets for text-to-text translation for all languages. Sample prompt format is given below.

```
""" Below is an instruction that describes a
translation task, paired with an input that
provides context in Source Language. Write
a response that appropriately completes
translation to desired Target Language.


### Instruction:
{}

### Input:
{}

### Response:
{}"""
```

A raw training sample data for Hindi translation after prompt formatting is shown below, similar method was used for other language translations.

Below is an instruction that describes a translation task, paired with an input that provides context in the Source Language. Write a response that appropriately completes translation to desired Target Language.
Instruction: Translate to Hindi
Input: it is an indoor scene
Response: यह एक इनडोर दृश्य है

#### 2.2.2 For Multi-modal involving both Image and Text

We passed the prompts in a CSV file with fields 'image id', and 'message'. The prompt in the "message" field is in the below format:
'message': ['content': 'describe the image in Hindi <loc ymin><loc xmin><loc ymax><loc xmax>', 'role': 'user','content': English text, 'role':'assistant']

### 2.3 Tokenization

Both model `unsloth/mistral-7b-v0.3` and tokenizer were used from unsloth library, tok-

| Set | Sentences | Tokens | | | | |
|---|---|---|---|---|---|---|
| | | English | Hindi | Malayalam | Bengali | Hausa |
| Train | 28,927 | 143,164 | 145,448 | 107,126 | 113,978 | 113,978 |
| D-Test | 998 | 4,922 | 4,978 | 3,619 | 3,936 | 3,936 |
| E-Test | 1,595 | 7,853 | 7,852 | 5,689 | 6,408 | 6,408 |
| C-Test | 1,400 | 8,186 | 8,639 | 6,044 | 6,657 | 6,657 |

Table 2: Statistics of our data used in the English→Hindi, English→Malayalam, English→Bengali, and English→Hausa tasks: the number of sentences and tokens in text-text translation.

| Set | Images | English Words | Hindi Words |
|---|---|---|---|
| Train | 28,927 | 143,164 | 145,448 |
| D-Test | 998 | 4,922 | 4,978 |
| E-Test | 1,595 | 7,853 | 7,852 |
| C-Test | 1,400 | 8,186 | 8,639 |

Table 3: Statistics of our data used in the English→Hindi multi-modal translation.

enizer is based on SentencePiece with Byte-Pair Encoding (BPE). This is the standard approach for tokenization in many modern transformer-based language models, including those similar to Mistral.



Figure 2: Instruction set in English-Hindi for multimodal translation

## 3 Experimental Details

This section describes the complete pipeline used to produce the translation systems for the WAT English-to-Low Resource Multimodal shared task submission.

### 3.1 EN-HI, EN-ML, EN-BN, EN-HA Text-only Translation

For EN–HI, EN-BN, EN-ML, and EN–HA text-only (E-Test and C-Test) translation, the study fine-tunes the pre-trained Mistral-7B model (Jiang et al., 2023), which has been fine-tuned utilizing only HVG, BVG, MVG, and HaVG Datasets; aiming to develop a high-quality machine translation system.

The Mistral-7B model is a cutting-edge language model that has been fine-tuned specifically for developing high-quality machine

translation systems. Leveraging its 7 billion parameters, Mistral-7B (Jiang et al., 2023) excels in capturing linguistic nuances and context, making it exceptionally adept at translating between languages with high accuracy. The fine-tuning process involves training the model on extensive and diverse datasets, allowing it to understand and generate translations that are not only precise but also contextually relevant.

### 3.2 EN-HI Multimodal Translation

This section discusses the multimodal translation pipeline for EN-HI. For EN-HI multimodal (E-Test and C-Test) translation, we used the object tags extracted from the HVG dataset images for image features and concatenated them with the text. The PaliGemma-3B model (Beyer et al., 2024) is finetuned on the Hindi-Visual-Genome dataset for English to Hindi Translation when a specific location is given in the input prompt as explained in Section 2.2.2. We used the script from LLaMa Factory (Zheng et al., 2024) with our configuration to fine-tune this model. During fine-tuning, we froze the vision tower and adjusted the parameters in the language model and projector layer. The hyperparameters are shown in Table 4.

## 4 Results

### 4.1 Text-only Translation

We present the official automatic evaluation results of our models for all the tasks we participated in Table 2, along with sample outputs in Table 3. After the fine-tuning process, these

| Hyperparameter | Value |
|---|---|
| Train Batch Size | 2 |
| Eval Batch Size | 8 |
| Learning Rate | $3 \times 10^{-6}$ |
| Epochs | 10 |
| Warm-up Steps | 50 |
| LR Scheduler | Cosine |
| Gradient Accumulation Steps | 8 |
| Optimizer | "Adam" |

Table 4: Training Hyperparameters.

models were used to generate translations for two distinct sets in each language: the evaluation set and the challenge set. The translation quality was assessed using the BLEU (Bilingual Evaluation Understudy) score and the RIBES (Ranking by Incremental Bilingual Evaluation System) score.

The English-to-Hindi model achieved a BLEU score of 41.60 on the evaluation set and 44.10 on the challenge set. Similarly, it attained a RIBES score of 0.82115 on the evaluation set and 0.8154 on the challenge set. These results underscore the model's robust performance and its ability to manage more complex or less typical translation tasks.

In the case of the English-to-Bengali model, a BLEU score of 43.70 was achieved on the evaluation set, with a slightly lower score of 35.60 on the challenge set. Similarly, it attained a RIBES score of 0.78975 on the evaluation set and 0.73534 on the challenge set. This indicates a robust overall performance and a commendable capability to handle nuanced translations specific to the Bengali language.

For the English-to-Malayalam model, the system achieved a BLEU score of 33.10 on the evaluation set and 18.10 on the challenge set. Similarly, it attained a RIBES score of 0.66837 on the evaluation set and 0.50594 on the challenge set. Despite a slightly lower score on the challenge set, the model still demonstrates a respectable performance in translating English to Malayalam.

Lastly, for the English-to-Hausa model, the system achieved a BLEU score of 49.80 on the evaluation set and 24.40 on the challenge set. Similarly, it attained a RIBES score of 0.81289 on the evaluation set and 0.66363 on the challenge set. This indicates a robust overall performance and a commendable capability

to handle nuanced translations specific to the Hausa language.

## 4.2 Multi-modal Translation Involving both Image and Text

Contrary to our expectations, the PaliGemma-3B model showed very poor results on the mentioned dataset and we tried to investigate the factors behind it. By qualitative analysis, we figured out that the location coordinates that we normalized during pre-processing may not be the right approach required for PaliGemma-3B. We found that the normalized [xmin, ymin, xmax, ymax] coordinates provided in the input prompt did not perfectly align with the model-generated captions. Instead, they pointed to a neighboring location in the image with a significant overlap. However, this mismatch in location led to a very poor BLEU score for the predicted captions.

## 5 Availability

The text-to-text and multimodal datasets, as well as the models, are freely available for research and non-commercial use under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 License via Hugging Face.

We have also released our experimental code on GitHub.[1]

### 5.1 EN-HI/ML/BN/HA Text-only Translation

Dataset: https://huggingface.co/datasets/OdiaGenAIdata/wat24_text_to_text_translation

Model: https://huggingface.co/OdiaGenAI-LLM/wat_mistral_7b_translate

### 5.2 EN-HI Multimodal Translation

Dataset: https://huggingface.co/datasets/sahoosk/Hindi-visual-genome_Train

Model: https://huggingface.co/sam2ai/odia-paligemma-2b-9900-v1.1

## 6 Conclusion

In this system description paper, we presented our approach for four tasks in WAT2024: (a) English→Hindi text-only and multimodal translation, (b) English→Malayalam text-only

---

[1] https://github.com/shantipriyap/ODIAGEN_WAT2024

| Translation Model | Translation Type | BLEU Score (Evaluation Set) | BLEU Score (Challenge Set) |
|---|---|---|---|
| English to Hindi | Text-to-Text<br>Multimodal | 41.60<br>0.50 | 44.10<br>- |
| English to Bengali | Text-to-Text | 43.70 | 35.60 |
| English to Malayalam | Text-to-Text | 33.10 | 18.10 |
| English to Hausa | Text-to-Text | 49.80 | 24.40 |

Table 5: Comparison of BLEU Scores for Different Translation Models and Types

| Translation Model | Translation Type | RIBES Score (Evaluation Set) | RIBES Score (Challenge Set) |
|---|---|---|---|
| English to Hindi | Text-to-Text<br>Multimodal | 0.8212<br>0.1512 | 0.8155<br>- |
| English to Bengali | Text-to-Text | 0.7898 | 0.7353 |
| English to Malayalam | Text-to-Text | 0.6684 | 0.5059 |
| English to Hausa | Text-to-Text | 0.8129 | 0.6636 |

Table 6: Comparison of RIBES Scores for Different Translation Models and Types

| | MALAYALAM | HINDI | BENGALI | HAUSA |
|---|---|---|---|---|
| English-Sentence-1 | silver car is parked | fine thin red hair | A stop light | A stop light |
| Target-Original | സിൽവർ കാർ പാർക്ക് ചെയ്തു | सूक्ष्म पतले लाल बाल | একটি স্টপ লাইট | Hasken tasha |
| Target-Translated | വെള്ളി കാർ പാർക്ക് ചെയ്തിരിക്കുന്നു | ठीक पतले लाल बाल | একটি স্টপ আলো | Hasken tasha |
| Gloss | Silver car has been parked | Correct thin red hair | A stop light | A stop light |
| Remarks (Comparison) | Translated version is more formal | Original version is better; "Fine" mistranslated by our model. | Original version is more colloquial | Both are identical |
| | | | | |
| English-Sentence-2 | eye of the pumpkin | the cross is black | This is a person | three zebras in the wild |
| Target-Original | മത്തങ്ങയുടെ കണ്ണ് | क्रॉस काला है | এটি একজন ব্যক্তি | alfadarai uku a cikin daji |
| Target-Translated | പമ്പ്കിന്റെ കണ്ണ് | क्रॉस काला है | এটি একজন ব্যক্তি | alfadarai uku a cikin daji |
| Gloss | Pumpkin's eyes | The cross is black | This is a person | Three zebras in the wild |
| Remarks (Comparison) | Model doesn't translate "pumpkin," which is colloquial | Both are identical | Both are identical | Both are identical |
| | | | | |
| English-Sentence-3 | pen on the paper | date and time of photo | the bird is black | a girl is standing. |
| Target-Original | പേപ്പറിൽ പേന | फोटो की तारीख और समय | পাখিটি কালো | yarinya tana tsaye |
| Target-Translated | പേപ്പറിൽ പേന | फोटो की तारीख और समय | পাখিটি কালো | yarinya tana tsaye |
| Gloss | Pen on the paper | Date and time of photo | The bird is black | A girl is standing |
| Remarks (Comparison) | Both are identical | Both are identical | Both are identical | Both are identical |

Table 7: Comparison between original translations and our model's translations for English-Malayalam, English-Hindi, and English-Bengali language pairs.

translation, (c) English→Bengali text-only translation, and (d) English→Hausa text-only translation. The results for the multimodal English→Hindi translation, which involves both image and text, were suboptimal due to improper normalization of the location coordinates for the `PaliGemma-3B` model. As a result, the model was unable to accurately map the provided coordinates in the prompt to the original image. We utilized the `PaliGemma-3B` model with a resolution of 448, which performed well in the translation tasks but failed to generate results relevant to the precise coordinates. Due to limitations in time and computing resources, addressing this issue has been deferred to future work. The code has been released on GitHub for use by other researchers.

## Acknowledgments

## References

Idris Abdulmumin, Satya Ranjan Dash, Musa Abdullahi Dawud, Shantipriya Parida, Shamsuddeen Hassan Muhammad, Ibrahim Sa'id Ahmad, Subhadarshi Panda, Ondřej Bojar, Bashir Shehu Galadanci, and Bello Shehu Bello. 2022. Hausa visual genome: A dataset for multi-modal english to hausa machine translation. arXiv preprint arXiv:2205.01133.

Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. 2024. Paligemma: A versatile 3b vlm for transfer. arXiv preprint arXiv:2407.07726.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In Proceedings of the 2010 conference on empirical methods in natural language processing, pages 944–952.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. arXiv preprint arXiv:2310.06825.

Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Shohei Higashiyama, Hideya Mino, Isao

| Image | Prompt | Predicted Answer (Hindi) | Gloss |
|---|---|---|---|
|  | describe the image in Hindi <loc166><loc177><loc263><loc298>. | एक महिला एक पुस्तक की दुकान से गुजर रही है। | A woman is passing by a bookstore. |
|  | describe the image in Hindi <loc6><loc120><loc31><loc136>. | एक बड़ा सफेद भवन जिसके शीर्ष पर एक घंटाघर है। | A large white building with a clock tower at the top. |
|  | describe the image in Hindi <loc95><loc284><loc214><loc300>. | एक सफेद भोजन कक्ष में एक कांच की टेबल सेट होती है। | A glass table set in a white dining room. |

Table 8: Comparison of user prompts, predicted answers in Hindi, and their English translations with corresponding images.

Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, et al. 2021. Overview of the 8th workshop on asian translation. In Proceedings of the 8th Workshop on Asian Translation (WAT2021), pages 1–45.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pages 311–318.

Shantipriya Parida, Ondřej Bojar, and Satya Ranjan Dash. 2019. Hindi visual genome: A dataset for multi-modal english to hindi machine translation. Computación y Sistemas, 23(4):1499–1505.

Shantipriya Parida, Alakananda Tripathy, Satya Ranjan Dash, and Shashikanta Sahoo. 2023. Mdolc: Multi dialect odia song lyric corpus.

Martin Popel, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. Nature communications, 11(1):1–15.

Arghyadeep Sen, Shantipriya Parida, Ketan Kotwal, Subhadarshi Panda, Ondřej Bojar, and Satya Ranjan Dash. 2022. Bengali visual genome: A multimodal dataset for machine translation and image captioning. In Intelligent Data Engineering and Analytics: Proceedings of the 9th International Conference on Frontiers in Intelligent Computing: Theory and Applications (FICTA 2021), pages 63–70. Springer.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyan Luo. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. arXiv preprint arXiv:2403.13372.

# Arewa NLP's Participation at WMT24

**Mahmoud Said Ahmad[1], Auwal Abubakar Khalid[2], Lukman Jibril Aliyu[3],**
**Babangida Sani[4], Mariya Sunusi Abdullahi[5]**

[1]*Federal University of Technology Babura (FUTB),* [2]*Bayero University Kano (BUK),* [3]*Arewa Data Science Academy,*
[4]*Arewa Data Science Academy,* [5]*Bayero University Kano (BUK)*
*msahmad.cs@futb.edu.ng, aka2000078.mcs@buk.edu.ng, lukman.j.aliyu@gmail.com,*
*sunusimariya@gmail.com, bsani480@gmail.com*

## Abstract

This paper presents the work of our team, "ArewaNLP," for the WMT 2024 shared task. The paper describes the system submitted to the Ninth Conference on Machine Translation (WMT24). We participated in the English-Hausa text-only translation task. We fine-tuned the OPUS-MT-en-ha transformer model and our submission achieved competitive results in this task. We achieve a BLUE score of 27.76, 40.31 and 15.85 on the Development Test, Evaluation Test and Challenge Test respectively.

## 1 Introduction

Machine translation (MT) is widely regarded as one of the most successful applications of natural language processing (NLP). It has seen significant advancements, particularly in the accuracy of its results. While MT has achieved near-human performance for several language pairs, it still faces challenges when dealing with low-resource languages or when incorporating other modalities (such as images.(Parida et al., 2021).

In the broader field of machine learning and deep learning, multimodal processing involves training models using a combination of different information sources such as images, audio, text, or video. By incorporating multimodal data, models can learn features from various subsets of these sources (depending on the data modality), leading to improved prediction accuracy. Multimodal machine translation leverages information from multiple modalities, with the expectation that these additional modalities will offer valuable alternative perspectives on the input data. Despite machine translation's near-human performance for several language pairs, it still faces difficulties in translating low-resource languages and effectively utilizing other modalities. (Sen et al., 2022).

WMT is a workshop on Machine Translation. WMT24 features the English-to-Low-Resource Multimodal Translation Shared Task, which involves Bengali, Hausa, Hindi, and Malayalam datasets from the Visual Genome project. These datasets include both text and images, providing a rich resource for research in English-to-[Hindi, Bengali, Malayalam, Hausa] Machine Translation and Multimodal studies.(Parida et al., 2024; Scientist, 2024).

In this system description paper, we outline our approach to the English-Hausa text-only translation task.

## 2 Dataset

We utilized the Hausa Visual Genome (HaVG) dataset (Abdulmumin et al., 2022) provided by the organizers. This dataset comprises 32,923 images with corresponding descriptions, divided into training, development, test, and challenge-test sets. The training set includes 28,930 English and Hausa sentence pairs, while the development set contains 998 sentences, the evaluation test set has 1,595 sentences, and the challenge test set consists of 1,400 sentences. A summary of the sentence statistics is provided in Table 1.

## 3 Experimental Details

The experimental setup involved fine-tuning a pre-trained sequence-to-sequence language model, specifically the OPUS-MT-en-ha model, which was pre-trained on English-Hausa data. Fine-tuning was performed using PyTorch and Hugging Face Transformers. For the English-Hausa text-only translation task, we fine-tuned the OPUS-MT-en-ha model[1], a translation model pre-trained on English-Hausa data by the Language Technology Research Group at the University of Helsinki[2] .

---

[1]https://huggingface.co/Helsinki-NLP/opus-mt-ha-en
[2]https://github.com/Helsinki-NLP

| Set | Sentences | Tokens | |
| --- | --- | --- | --- |
| | | English | Hausa |
| Training set | 28,930 | 147,219 | 144,864 |
| Development test | 998 | 5,068 | 4,978 |
| Evaluation test | 1,595 | 8,079 | 7,952 |
| Challenge test | 1,400 | 8,411 | 9,514 |
| Total | 32,923 | - | - |

Table 1: Statistics of data used in the English-Hausa text-only translation: the number of sentences and tokens.

## 3.1 Preprocessing

The Hausa Visual Genome dataset was prepared to train the translation model. The preprocessing phase involved preparing the Hausa Visual Genome (HaVG) dataset for training the translation model, The data was loaded using 'pandas' and converted into Hugging Face 'Dataset' objects for both English and Hausa texts. We employed the 'Helsinki-NLP/opus-mt-en-ha' tokenizer to tokenize the text, truncating or padding sequences to a maximum length of 128 tokens. The tokenized data was then formatted for PyTorch, including input IDs, attention masks, and labels, to ready it for training.

## 3.2 Model Fine-Tuning

Model fine-tuning is a crucial step in which the pre-trained model is adapted to the specific task of English-Hausa translation. We fine-tuned a pre-trained sequence-to-sequence language model using PyTorch and Hugging Face Transformers. The model was trained for 3 epochs with an AdamW[3] optimizer and a linear learning rate scheduler. Training was conducted on a GPU in batches of 8, with evaluation performed after each epoch. Upon completion, the fine-tuned model and tokenizer were saved. Fine-tuning not only enhanced the model's translation accuracy but also allowed it to perform well on different test sets, although it faced challenges with more difficult content as seen in the Challenge Test results.

This methodology enabled the model to achieve competitive BLEU scores on the various test sets, demonstrating its effectiveness in translating between English and Hausa, albeit with some room for improvement in handling more complex or less familiar content

[3]https://keras.io/api/optimizers/adamw

## 4 Results

Table 4 presents the results of automatic evaluation of our model.

**Development Test (D-Test BLEU: 27.76):** The model scored 27.76 on the Development Test set. This is a solid result, indicating that the model produces translations that are reasonably accurate, though there's some room for improvement. This test set is typically used during the model's development phase to fine-tune its performance.

**Evaluation Test (E-Test BLEU: 40.31):** On the Evaluation Test set, the model achieved a BLEU score of 40.31, which is quite a bit higher than on the Development Test set. This suggests that the model is particularly good at translating the kinds of sentences found in this set, perhaps because they are similar to what the model has seen during training.

**Challenge Test (C-Test BLEU: 15.85):** The model scored 15.85 on the Challenge Test set, which is significantly lower than the other two scores. This suggests that the Challenge Test set contains more difficult or unfamiliar content, making it harder for the model to produce accurate translations.

**Zero-shot vs. Finetuned Scenarios**

The zero-shot evaluation BLEU scores (table 3) are very low compared to the fine-tuned results (table 4). This demonstrates that without prior exposure or training on this specific data, the model struggles to perform accurate translations. These low BLEU scores suggest that the model's ability to generalize to completely unseen data (zero-shot scenario) is limited.

The significant difference between fine-tuned and zero-shot BLEU scores across all sets illustrates the importance of HaVG data. Fine-tuning has allowed the model to learn the translation patterns within the datasets, leading to far superior performance compared to the zero-shot setting.

**English-Hausa Translation Examples**

Table 2 presents sample English sentences alongside their Hausa translations, sourced from the challenge test set. Some examples are straightforward, where the model successfully translated simple, clear sentence structures. However, other examples are more challenging, showcasing the model's ability to handle complex or ambiguous translations. For instance, in examples 7 and 8, the word "cross" appears, which can refer to either a cruciform symbol or the act of crossing a street. The model accurately interpreted the context in both cases, delivering correct translations for each meaning. These more difficult examples illustrate the differences between the Dev, Eval, and Challenge sets, with the Challenge set specifically designed to test the model's performance by including context-dependent and nuanced sentences. The model's ability to navigate these complexities demonstrates its overall effectiveness.

| S/N | English | Hausa Translation |
|-----|---------|-------------------|
| 1 | A second pizza in a pan. | Pizza na biyu a cikin kwanon suya. |
| 2 | A girl on the tennis court is preparing to hit the ball. | Wata yarinya a filin wasan tanis tana shirin buga kwallon. |
| 3 | Knife block sitting on counter with knives in it. | Sandar wuka zaune akan kan tebur tare da wukake a ciki. |
| 4 | The players' socks are blue. | Yan wasan safa sune shui. |
| 5 | Balconies on the second story of the buildings. | Baranda akan bene na biyu na gine-ginen. |
| 6 | Beige stairway going to second level. | Matakala na beige zuwa bene na biyu. |
| 7 | The woman is waiting to cross the street. | Matar tana jira ta tsallaka titi. |
| 8 | A black cross on a vertical stabilizer. | Gicciye mai baar fata akan mai tsaye tsaye. |
| 9 | Man cross country skiing. | Mutum ya tsallaka kan asa a lokacin tsere. |

Table 2: Sample of English to Hausa translations generated by our model.

| D-Test BLEU | E-Test BLEU | C-Test BLEU |
|-------------|-------------|-------------|
| 1.87 | 1.95 | 2.56 |

Table 3: Results of text-only translation task: Zero-shot

| D-Test BLEU | E-Test BLEU | C-Test BLEU |
|-------------|-------------|-------------|
| 27.76 | 40.31 | 15.85 |

Table 4: Results of text-only translation task: Fine-tuned model

## 5 Conclusion

This paper describes our system for English-to-Hausa text-only translation. The system performs well on more standard test sets (especially the Evaluation Test) but struggles with more challenging or unusual content, as seen in the Challenge Test results. This indicates that while the system is effective in many scenarios, it may need further training to handle more complex translation tasks. We plan to extend our work to include English-Hausa multimodal translation and image captioning tasks in the future.

## Ethics Statement

In our work on the English-to-Hausa text-only translation task, we adhered to the highest standards of ethical research and data use. The datasets employed, including the Hausa Visual Genome dataset, were provided under appropriate licenses, and we ensured that all data used was handled in accordance with the terms specified by the providers. Our research also followed guidelines for responsible AI development, including fairness, transparency, and privacy considerations. We took particular care to avoid biases in our models that could negatively impact the communities whose languages we are working with. Additionally, we acknowledge the potential risks of deploying machine translation systems in sensitive contexts and emphasize the importance of human oversight in such applications.

Technology Research Group at the University of Helsinki for the OPUS-MT-en-ha model. Computational resources provided by our institutions were crucial for this research.

# References

Idris Abdulmumin, Satya Ranjan Dash, Musa Abdullahi Dawud, Shantipriya Parida, Shamsuddeen Muhammad, Ibrahim Sa'id Ahmad, Subhadarshi Panda, Ondřej Bojar, Bashir Shehu Galadanci, and Bello Shehu Bello. 2022. Hausa visual genome: A dataset for multi-modal English to Hausa machine translation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6471–6479, Marseille, France. European Language Resources Association.

Shantipriya Parida, Ondřej Bojar, Idris Abdulmumin, and Shamsuddeen Hassan Muhammad. 2024. Wat2024 english-to-lowres multi-modal translation task. https://ufal.mff.cuni.cz/wat2024-multimodal. Accessed: 2024-08-29.

Shantipriya Parida, Subhadarshi Panda, Ketan Kotwal, Amulya Ratna Dash, Satya Ranjan Dash, Yashvardhan Sharma, Petr Motlicek, and Ondřej Bojar. 2021. NLPHut's participation at WAT2021. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 146–154, Online. Association for Computational Linguistics.

ML Scientist. 2024. Join the wmt2024: English-to-lowres multi-modal translation task! https://mlscientist.com/join-the-wmt2024-english-to-lowres-multi-modal-translation-task/. Accessed: 2024-08-29.

Arghyadeep Sen, Shantipriya Parida, Ketan Kotwal, Subhadarshi Panda, Ondřej Bojar, and Satya Ranjan Dash. 2022. Bengali visual genome: A multimodal dataset for machine translation and image captioning. In *Intelligent Data Engineering and Analytics*, pages 63–70, Singapore. Springer Nature Singapore.

# Multimodal Machine Translation for Low-Resource Indic Languages: A Chain-of-Thought Approach Using Large Language Models

**Pawan Kumar Rajpoot**[*]
pawan.rajpoot2411@gmail.com

**Nagraj N Bhat**[*]
nagbhat25@gmail.com

**Ashish Shrivastava**[*]
ashish3586@gmail.com

## Abstract

This paper presents the approach and results of team v036 in the English-to-Low-Resource Multi-Modal Translation Task at the Ninth Conference on Machine Translation (WMT24). Our team tackled the challenge of translating English source text to low-resource Indic languages, specifically Hindi, Malayalam, and Bengali, while leveraging visual context provided alongside the text data. We used InternVL2 for extracting the image context along with Knowledge Distillation from bigger LLMs to train Small Language Model on the tranlsation task. During current shared task phase, we submitted best models (for this task), and overall we got rank 3 on Hindi, Bengali, and Malyalam datasets. We also open source our models on huggingface.[1]

## 1 Introduction

With the recent advances in text generative AI Achiam et al. (2023); Dubey et al. (2024); Yang et al. (2024) and Diffusion based Dhariwal and Nichol (2021); Nichol and Dhariwal (2021); Saharia et al. (2022); Ramesh et al. (2022) models, multimodal approaches have gained significant traction. The concept of a model to understand both text and visual contexts provides a unique advantage for these models to understand the real world. On the other end, Machine Translation has been one of the most important task in NLP world. Since its origin, the MT task has undergone large shifts from rule based Nirenburg (1989); Chen et al. (2007) to complex Neural network based approaches and recently Transformer Vaswani (2017); Yin and Read (2020); Xu et al. (2024) based approaches. In the recent days with the advancements in the field of NLP, Multimodal Machine Translation (MMT) has evolved as an important research field, wherein the Model utilizes both vision and text information to achieve the translation task. This would better equip the model with additional context information and thus reducing the issues due to polysemy or missing text context. MMT finds its application in various fields like Media, Retail, Automobile etc. In this work we explore the problem of English to Lowres Multimodal Translation for Hindi, Bengali and Malyalam languages. The task requires translating a short English caption of the rectangular region to one of these languages, given the image context. There are multiple approaches possible which can be largely classified into :

- Text-only translation (Source image not used)

- Image captioning (English source text not used)

- Multi-modal translation (uses both the image and the text)

We strongly feel that Multi-modal translation approach would best solve the problem due to more context information. In this paper, we propose a novel unconstrained approach to solve the Lowres MMT task for Hindi, Bengali and Malyalam languages. Our solution tries to merge the best of both text and language contexts. In particular, our key contributions are:

- Fusing Multimodal image context with improved language understanding : We provide a concise yet effective approach to combine context information from vision to text description

- Advanced Chain of Thought reasoning for language translation: Our approach to step by step reasoning ustilizing the COT, gives a whole new perspective to enhance the ability of the model to comprehend better.

---

[1]https://huggingface.co/team-v036
[*]Authors contributed equally to this work

- **Custom finetuning** : Our approach of custom finetuning on target languages on training samples, equips the model to better perform on the MMT task.

## 2 Data

Visual Genome introduced by Krishna et al. (2017) is a rich dataset to enable the modelling of complex cognitive interactions and relations between objects in an image. Based on this dataset, Parida et al. (2019) introduce the Hindi Visual Genome dataset, which is a multi-modal dataset consisting of text and images suitable for English-Hindi multi-modal machine translation task. They select short English segments (captions) from Visual Genome along with the associated images and automatically translate them to Hindi with a careful manual post-editing( Parida et al. (2019) ) The dataset takes into account ambiguous English words based on the embedding. similarity and manual selection of certain cases where image helps to resolve the ambiguity( Parida et al. (2019) ). Hence this is a perfect dataset suited for the task. Similarly Sen et al. (2022) propose the Bengali Visual Genome Dataset which is manually labelled on HVG samples and Parida et al. (2019) curated the malyalam Visual Genome Dataset.

All three (Hindi, Bengali and Malyalam) dataset consists of 29k training samples, 1k dev set, 1.6k evaluation set and 1.4k challenge set.

The evaluation of the models were performed with BLEU metrics (Papineni et al. (2002)) on challenge and evaluation set independently. Along with these a manual labeller evaluation is also performed, subject to availability.

## 3 Related Work

MMT has gained increasing attention in recent years as a way to leverage visual information to improve translation quality. Several shared tasks and datasets have been introduced to advance research in this area, with a particular focus on low-resource languages. The Workshop on Asian Translation (WAT) has played a key role in promoting MMT research for Asian languages. Parida et al. (2019) introduced the first Hindi Visual Genome task at WAT 2019, using the Hindi Visual Genome 1.0 dataset (Parida et al. (2019)). This dataset contains English image captions paired with Hindi translations and associated image regions. The task evaluates systems on their ability to translate from English to Hindi while incorporating visual context. Subsequent iterations of WAT expanded the Hindi Visual Genome used an updated Hindi Visual Genome 1.1 dataset and introduced new evaluation tracks, including Hindi image captioning. The latest WAT 2021 (Nakazawa et al. (2021)) further refined the Hindi task and introduced a new English-Malayalam MMT task using the Malayalam Visual Genome dataset (Parida and Bojar, 2021). This represented the first multimodal translation dataset for Malayalam. For the Hindi task, recent approaches have focused on leveraging object tags extracted from images (Gupta et al. (2021)) and region-specific captioning (Parida et al. (2021)) to enhance translation quality. The introduction of the Malayalam task provides an opportunity to evaluate MMT techniques on a new low-resource language. While Hindi and Malayalam have been addressed in shared tasks, Bengali has seen less attention for MMT despite being widely spoken. The creation of a Bengali Visual Genome dataset, following the model of Hindi and Malayalam, would fill an important gap and enable MMT research for another major South Asian language. Overall, the development of these language-specific visual genome datasets has been crucial for advancing MMT for low-resource Indian languages. They provide much-needed benchmarks and drive innovation in incorporating visual context for translation. Expanding to additional languages like Bengali represents an important direction for broadening the scope of MMT research in the Indian context.

## 4 Approach

Our overall approach follow a three step process as seen in Figure 1.

### 4.1 Stage 1: Fusing Multimodal image context with improved language understanding

In this stage, we first extract context from cropped visual data using a powerful open-source Multimodal Large Language Model (MLLM)- InternVl2-8B (Chen et al. (2023, 2024)). This model demonstrates powerful capabilities in handling complex multimodal data and achieves state of art numbers on many open VQA tasks. Figure 2 shows a sample image and its description. We feed the output of segement description as an input into a Rapid Automatic Keyword Extraction (RAKE) algorithm (Rose et al. (2012)) which is an efficient keyword extraction algorithm. The top extracted key phrases

Figure 1: Overall Approach.

are selected and used as hastags to provide context to the source English text. This way we condense the full description into short and concise information for the next stage. This ensures that the further step do not completely rely on the image context also but rather use the original text but still use the relevant information from the image descriptors. This step is common for both train and evaluation process.

## 4.2 Stage 2: Advanced Chain of Thought reasoning for language translation

Chain-of-thought (CoT) (Wei et al. (2022)) prompting enables complex reasoning capabilities through intermediate reasoning steps. It is shown that the models ability is substantially improved by making them produce step by step reasoning. We employ this ability of Large Language model to solve the task in a more understandable and reasonable approach by decomposing the problem into multiple sub-problems. We use State of the Art LLaMa 405B (Dubey et al. (2024)) model to generate the CoTs for the training data. The model is provided with the English caption text that needs to be translated along, the hashtags generated in previous step, and the target language caption. The model is then asked to generate the step by step reasoning for converting the source caption text to target language text along with the condensed context information provided. Following prompt (Table 1) template is used to get the CoT from the bigger model.

Table 1: Prompt for CoT reasoning generation from bigger model (LLAMA 3.1 405B)

TASK:
ASSUMING YOU ARE A ENGLISH - HINDI translation expert, For given context of an image related to a original sentence, English sentence and translation of the sentence in hindi. Give reason on why this translation is the correct translation....ASsume that you secretly know the answer......DO NOT TRY TO FIX Translation...reason for whatever is given only.....reason SHOULD be proper Chain of Thought format in properly divided steps for the answer......give maximum 5 steps which are most important ones .......
Context: {RAKE HASTAGS}
English Sentence: {SOURCE TEXT}
Hindi Sentence: {TARGET TEXT}

Table 2: Prompt for SFT LLAMA 3.1 8B

TASK:
ASSUME YOU ARE AN ENGLISH-HINDI translation expert. Given an image description in English, image context and reasoning/CoT in English, translate the image description in Hindi. Use the image context to solve ambiguity if required. Note: DO NOT USE the image context in translations, just use them for disambiguation.
IMAGE DESCRIPTION:
{SOURCE TEXT}
IMAGE CONTEXT:
{RAKE HASHTAGS}
REASONING:
{GENERATED CoT}
RESPONSE:
{TARGET TEXT}

## 4.3 Stage 3: Custom fine-tuning

In stage 3 we train a smaller model to perform the task of translation, we finetune a LLaMA 3.1-8B-Instruct model on training samples using data from previous stages. The model is trained by providing the English caption along with hashtag contexts, CoT reasoning and the final answer. A sample prompt is shown in Table 2. We use LORA finetuning with rank=64 and alpha=128. We use following template for the training data so that the CoT step is more aligned to as what humans think, that is first source, then CoT and finally the target text.

During **inference**, we provide the finetuned model with source English caption and context and ask it to come up with the Reasoning and the answer. We then use a post processor script to filter out the final answer from the model output.

These 3 fundamental steps are performed for all 3 languages and we curate one PEFT model for each language.

Figure 2: Sample Data with reference the image segment, its corresponding source and target text along with the key phrase extraction from Internvl2 descriptions

## 5 Experimental setup

In our experimental setup, we fine-tuned the LLaMA 3.1-8B-Instruct model for the translation task using Quantized Low-Rank Adaptation (QLoRA) (Hu et al. (2021); Dettmers et al. (2024)). The LoRA configuration was carefully selected to balance performance and computational efficiency. We set the rank (r) to 64 and the alpha parameter to 128, with a lora_alpha value of 0.05. Notably, we applied LoRA to all target modules in the model architecture, ensuring a comprehensive adaptation across the entire network.

For the optimization process, we employed a learning rate (lr) of 0.003, coupled with a cosine learning rate scheduler. This scheduling strategy allowed for dynamic adjustments to lr, potentially aiding in convergence and generalization. The model was trained for two epochs, striking a balance between sufficient learning and computational constraints.

The chosen LoRA hyperparameters strike a balance between model capacity and computational efficiency, with the rank of 64 providing sufficient expressiveness for the adaptation.

By leveraging Quantized LoRA and carefully selected training parameters, we aim to achieve high-quality translation performance while minimizing computational resources and training time. We used A100 40GB VRAM and 84GB RAM single node machine to fine tune our models.

## 6 Results

The results show that our Multimodal approach of using multistage image description extraction clubbed with CoT is an effective approach to solve this task leveraging the knowledge of Large language models. Table 3 shows the results for all 3 Indic languages on Evaluation and Challenge set.

Our numbers are very close to SOTA numbers. The SOTA (baseline) approach is based on a fine-tuning of NLLB model on captions of Object tags of original along with synthetic images using DETR model. However, we do not use any additional image set in our process

Table 3: Results (BLEU Scores) on languages comparing to SOTA .

| Language | Evaluation Set | Challenge Set |
| --- | --- | --- |
| Hindi | 0.446/**0.45** | 0.432/**0.534** |
| Bengali | 0.441/**0.506** | 0.339/**0.487** |
| Malyalam | 0.427/**0.519** | 0.333/**0.422** |

The data analysis of the final output revealed a set of cases where the output is technically correct, yet contains variations in tokens compared to the gold set. This suggests that human evaluation could potentially yield higher accuracy, and relying solely on the BLEU score for this task may not fully capture the quality of the output.

## Limitations

Given that our approach heavily depends on multiple stages involving large language models, it may not be ideally suited for environments with limited resources. The complexity and computational demands of such models could pose challenges in settings where processing power, memory, or bandwidth are constrained. Additionally, this approach leverages the inherent knowledge embedded within the LLMs being used. The effectiveness of the method is closely tied to the pre-existing information and understanding that these models have acquired during training, which may be influenced by the data used for its training.

## Ethics Statement

Our work proposes an innovative approach to addressing the challenge of translating low-resource English to Indic languages - Hindi, Bengali, and Malayalam. In conducting our research, we have carefully considered the ethical implications of data usage. As a result, we have chosen to exclusively rely on the data provided by the Task administrators for our experiments, refraining from incorporating any additional external data sources. This ensures that our approach remains transparent and aligns with the ethical standards expected in this field. However, while using this approach for real world application, data privacy and consent should be given careful considerations.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Yu Chen, Andreas Eisele, Christian Federmann, Eva Hasler, Michael Jellinghaus, and Silke Theison. 2007. Multi-engine machine translation with an open-source smt decoder. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 193–196.

Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.

Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Kshitij Gupta, Devansh Gautam, and Radhika Mamidi. 2021. Volta at semeval-2021 task 6: Towards detecting persuasive texts and images using textual and multimodal ensemble. *arXiv preprint arXiv:2106.00240*.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73.

Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Shohei Higashiyama, Hideya Mino, Isao Goto, Win Pa Pa, Anoop Kunchukuttan, Shantipriya Parida, et al. 2021. Overview of the 8th workshop on asian translation. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 1–45.

Alexander Quinn Nichol and Prafulla Dhariwal. 2021. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR.

Sergei Nirenburg. 1989. Knowledge-based machine translation. *Machine Translation*, 4(1):5–24.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Shantipriya Parida, Ondřej Bojar, and Satya Ranjan Dash. 2019. Hindi visual genome: A dataset for multi-modal english to hindi machine translation. *Computación y Sistemas*, 23(4):1499–1505.

Shantipriya Parida, Subhadarshi Panda, Ketan Kotwal, Amulya Ratna Dash, Satya Ranjan Dash, Yashvardhan Sharma, Petr Motlicek, and Ondřej Bojar. 2021. Nlphut's participation at wat2021. In *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*, pages 146–154.

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3.

Stuart J Rose, Vernon L Crow, Nicholas O Cramer, et al. 2012. Rapid automatic keyword extraction for information retrieval and analysis. Technical report, Pacific Northwest National Laboratory (PNNL), Richland, WA (United States).

Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu

Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494.

Arghyadeep Sen, Shantipriya Parida, Ketan Kotwal, Subhadarshi Panda, Ondřej Bojar, and Satya Ranjan Dash. 2022. Bengali visual genome: A multimodal dataset for machine translation and image captioning. In *Intelligent Data Engineering and Analytics: Proceedings of the 9th International Conference on Frontiers in Intelligent Computing: Theory and Applications (FICTA 2021)*, pages 63–70. Springer.

A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *arXiv preprint arXiv:2401.08417*.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Kayo Yin and Jesse Read. 2020. Attention is all you sign: sign language translation with transformers. In *Sign Language Recognition, Translation and Production (SLRTP) Workshop-Extended Abstracts*, volume 4.

## A  Appendix

A sample example of COT generation is shown below:



```
Example 1
Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

###TASK:
ASSUME YOU ARE AN ENGLISH-BENGALI translation expert.
Given an image description in English, image context and resoning in English, translate the image description in Bengali. Use the image context to improve the translation, if required.
Note: DO NOT TRY to correct the image description or image context in English.

###IMAGE DESCRIPTION: couples walking in the rain

###IMAGE CONTEXT: #typical rainy day #also dressed casually #white tank top

###REASONING:
Step 1: "Couples" refers to two people in a romantic relationship.
Step 2: In Bengali, "দম্পতি" (dampati) specifically means a married couple.
Step 3: The context shows a casual, relaxed atmosphere, implying a romantic couple.
Step 4: "Walking" is accurately translated to "হাঁটছেন" (hãtchhen), which is the correct verb form.
Step 5: "In the rain" is correctly translated to "বৃষ্টিতে" (brstite), which sets the scene for the action.

###RESPONSE: দম্পতিরা বৃষ্টিতে হাঁটছেন
```

Figure 3: Sample COT training data

A few data analysis samples where we note that the translation is mostly valid but the gold may have different set of words.

source : gold religious cross on top of golden ball
gold: स्वर्ण गेंद के शीर्ष पर स्वर्ण धार्मिक क्रॉस
ours: सोने की गेंद के ऊपर सोने का धातु

source : front springs of motocross motorcycle
gold: मोटोक्रॉस मोटरसाइकिल के सामने स्प्रिंग्स
ours: मोटोक्रॉस मोटरसाइकिल के आगे के स्प्रिंग

source: A display of different phone models
gold: विभिन्न फोन मॉडलों का प्रदर्शन
ours: विभिन्न मोबाइल मॉडलों का एक प्रदर्शन

source: rapid in fast water
gold: तेजी से पानी में
ours: तेज़ पानी में तेज़

source: tail fine view of a red plane
gold: एक लाल विमान की पूंछ का दृश्य
ours: लाल विमान का टेल फाइन व्यू

source : white springs boxes and gears below train engine
gold: सफेद स्प्रिंग्स बक्से और ट्रेन इंजन के नीचे गियर
ours: ट्रेन इंजन के नीचे सफेद स्प्रिंग्स के डिब्बे और गियर हैं।

# Chitranuvad: Adapting Multi-Lingual LLMs for Multimodal Translation

**Shaharukh Khan, Ayush Tarun, Ali Faraz, Palash Kamble, Vivek Dahiya,**
**Praveen Pokala \*, Ashish Kulkarni \*, Chandra Khatri \*,**
**Abhinav Ravi \* and Shubham Agarwal \***

Krutrim AI, Bangalore, India
`* Senior Contributors`
Contact: {shaharukh.khan, abhinav.ravi, shubham.agarwal1}@olakrutrim.com

## Abstract

In this work, we provide the system description of our submission as part of the *English-to-Lowres Multimodal Translation Task* at the Workshop on Asian Translation (WAT2024). We introduce Chitranuvad, a multimodal model that effectively integrates Multilingual LLM and a vision module for Multimodal Translation. Our method uses a ViT image encoder to extract visual representations as visual token embeddings which are projected to the LLM space by an adapter layer and generates translation in an autoregressive fashion. We participated in all the three tracks (Image Captioning, Text-only and Multimodal translation tasks) for Indic languages (ie. English translation to Hindi, Bengali and Malyalam) and achieved SOTA results for Hindi in all of them on the Challenge set while remaining competitive for the other languages in the shared task.

## 1 Introduction

Recently, there has been an increased interest in Multimodal Machine Translation (MMT) task (Calixto and Liu, 2017; Delbrouck and Dupont, 2017; Elliott and Kádár, 2017; Yao and Wan, 2020) which involves translation between language pairs, incorporating other modalities (like images) as an auxiliary information. The visual cues act as 'symbol grounding' (Fodor, 1975; Harnad, 1990, 2003, 2005), helping to resolve ambiguities in language (Rainie et al., 2012; Hu et al., 2014; Specia et al., 2016; van Miltenburg et al., 2019; Caglayan et al., 2020) by learning to connect language and perception (Mooney, 2008; Bisk et al., 2020). For example, in order to correctly translate the word *court* in Figure 1, the model has to infer from the image that the statement is about tennis court and not the court as government institution.

Prior works mostly focused on translation from English to European languages (Elliott et al., 2016; Specia et al., 2016) while the Indic languages remain largely unexplored, with an exception of the

MMT shared task at the Workshop on Asian Translation (WAT) (Nakazawa et al., 2019, 2020, 2021, 2022, 2023).

The *English-to-Lowres Multimodal Translation Task* at WAT-2024 targets the MMT task for three Indic medium-to-low-resource languages (Hindi, Bengali, Malayalam) and a low-resource African language Hausa. To assess the importance of the image modality, the task has been decoupled into three tracks: **1).** *Text-only translation* where the source image is not used, **2).** *Image Captioning* where English source text is not used and **3).** *Multimodal translation* which uses both the image and the text. We participated in all the three tracks for Indic languages only (Hindi, Bengali, Malayalam) under a non-constrained and proprietary multilingual and multimodal Large Language Model (LLM): *Chitranuvad*[1].

In this paper, we provide a description of our multimodal LLM where we leverage a multilingual LLM backbone Krutrim (Team, 2024b), coupled with a visual image encoder. Our contributions could thus be summarized as follows:

- We introduce Chitranuvad, a Large Multimodal model, adapted for multi-lingual translation, which leverages images and language modalities to provide an image grounded translation of the English sentence in the target Indic languages.
- We showcase the effectiveness of task specific finetuning on the Visual Genome translation datasets and achieve SOTA performance.
- We evaluate Chitranuvad and prior baselines on the English-to-Lowres Multimodal Translation Task and demonstrate the ability of our model to perform grounded translation, using different training strategies and ablations.

The rest of the paper is organized as: Section 2 presents related research on multimodal machine

---

[1] Chitranuvad literally means Image Translate in Hindi

Figure 1: Multimodal Machine Translation task as part of English-to-lowres track where the source sentence is translated to multiple Indic languages (Hindi, Bengali, Malayalam) grounded in the image. Meaning of words like "court" and "right" in the translations can vary significantly depending on the visual context.

translation while Section 3 explains our Chitranuvad model recipe in detail. We present the datasets used in Section 4, followed by experimental findings in Section 5 and conclusion in Section 6.

## 2   Related Work

Early Neural Machine Translation (NMT) and Image captioning systems (Show, 2015; Gao et al., 2018) were based on Recurrent Neural Networks (RNNs) and their variants (Cho et al., 2014; Sutskever et al., 2014; Cho, 2014; Hochreiter and Schmidhuber, 1997), often incorporating attention mechanisms (Bahdanau et al., 2014). The seminal work of transformers (Vaswani, 2017) paved the way for the development of high-quality image captioning (Chen et al., 2021) as well as translation systems (Lewis, 2019), even for low-resource languages (Dabre et al., 2021; Gala et al., 2023a). Multimodal Machine Translation (MMT) systems witnessed a similar shift in their approrach (Caglayan et al., 2016; Yao and Wan, 2020; Guo et al., 2023). Prior submissions to the MMT task at Workshop on Asian Translation (Gain et al., 2021; Gupta et al., 2021; Parida et al., 2022; Dash et al., 2023; Shahid et al., 2023) also fall in this category.

The next generation of Multimodal LLMs (Lu et al., 2024a; Laurençon et al., 2024; Tong et al., 2024; Xue et al., 2024) can handle a variety of complex tasks, including machine translation and captioning, by utilizing cutting-edge architectures

as an unified general purpose agent. These models often rely on pre-trained LLMs, with an exception of few, which train the models from scratch (Team, 2024a; Lu et al., 2024b). Most of these Vision Language Models (VLMs) follow the architecture of (Liu et al., 2023a) where a CLIP (Radford et al., 2021) or a similar encoder is used to encode the image and projected into LLM's representation space using an adapter layer. Notably, Wang et al. (2023) offers a departure from conventional architectures by using distinct matrices and Feed Forward Networks for image modalities. Recent developments replace the image encoder with SigLIP (Zhai et al., 2023a) and the single-layer MLP projector with attention-based pooling (Laurençon et al., 2024).

Advanced backbone LLMs (Brown et al., 2020; Touvron et al., 2023; Achiam et al., 2023; Team et al., 2023; Jiang et al., 2024; Team et al., 2024) however have a primary focus for English and European languages. There have been relatively few LLMs for Indic languages, such as Airavata (Gala et al., 2024), Navarsa (Labs, 2023), Kannada LLaMA, Tamil LLaMA (Balachandran, 2023), Odia LLaMA (Kohli et al., 2023), to name a few. However, most of these LLMs are an extension and finetuned version of LLaMA/Gemma for Indic languages, which don't fully capture the nuances of the language. This could be attributed to the fact that Indic languages are under-represented in Common Crawl (which majorly forms the train-

Figure 2: Chitranuvad model architecture with the three stage training pipeline described in Section 3.

ing corpus of LLMs), despite India constituting 18% of the global population. Hindi, for example, does not show-up in the top 20 languages despite being the 3rd most spoken (Buck et al., 2014; Penedo et al., 2023). Closed-source models such as Krutrim (Team, 2024b) and Sutra (Bendale et al., 2024) represent exceptions, as they are trained from scratch. Currently, PALO (Maaz et al., 2024) is a multimodal LLM that supports only Hindi and Bengali. However, to the best of our knowledge, there are no other open-source multimodal LLMs trained specifically for low-resource Indic languages. In contrast, we developed a multilingual multimodal system that supports 10 Indic languages.

## 3  Model and Training Recipes

Figure 2 provides an overview of our architecture and the multi-stage training pipeline. Our *Chitranuvad* model architecture borrows heavily from LLaVA-like models (Liu et al., 2023a, 2024), where we use pre-trained Krutrim LLM (Team, 2024b) instead, as the autoregressive multi-lingual LLM backbone. Our Krutrim LLM is trained across 10 languages and natively supports all the 3 Indic languages (Hindi, Bengali, Malayalam) used as part of the shared task.

For the multimodal training, we first encode images through a vision encoder. Next, the modality projection (adapter/connector) layer projects the vision embeddings into the LLM embedding

space, creating a sequence of visual tokens. The multi-lingual LLM then generates the response conditioned on these visual embedding tokens. The Krutrim LLM model supports a context length of 4096 tokens, out of which 576 tokens are used for the image representation, obtained after the modality projector layer. For the projection layer, we experiment with both single layer projection (Liu et al., 2023b) as well as a two-layer MLP vision-language connector with non-linearity (Liu et al., 2023a). We also experiment with pre-trained CLIP ViT-L/14@336px (Radford et al., 2021) as well as SigLIP-SO400M (Zhai et al., 2023b) for the vision encoder. Similar to the LLaVA model, we generate multi-turn conversational data for instruction tuning our model, which we expand upon in Section 4. We train our model in multiple stages:

**Stage 1: Pre-Training (PT) for Feature Alignment.** In this stage, we do the pre-training with image-text pairs, where the projector layer is trained while the vision encoder and LLM is kept frozen. Here, each sample is treated as a single-turn conversational instruction tuning data.

**Stage 2: Instruction Tuning.** Similar to LLaVA models (Liu et al., 2023b,a), we also keep the vision encoder frozen during the second stage of training. However, here we also update the LLM weights apart from tuning the modality projection layer. This stage aims to build a general purpose Multimodal agent (chatbot) which can follow com-

| Split | #Instances | English | Hindi | Bengali | Malayalam |
|-------|-----------|---------|-------|---------|-----------|
| Train | 28930 | 5.09 | 5.13 | 4.07 | 3.86 |
| Valid | 998 | 5.08 | 5.04 | 4.06 | 3.75 |
| Test | 1595 | 5.07 | 4.95 | 4.14 | 3.76 |
| Challenge | 1400 | 6.04 | 6.35 | 4.92 | 4.48 |

Table 1: Total number of instances and average number of tokens for the text in English and splits of different Visual Genome datasets in other languages.

---

**Multimodal Translation**:
*Human:* You are given an image and coordinates of a bounding box as: x1={x1}, y1={y1}, x2={x1+x2}, y2={y1+y2}. Using the context of the objects or items available in the bounding box translate the following sentence from English into {lang} language. You are also provided labels of the objects in the image as: {labels}. English sentence is: {sentence}.
*System:* {translation}.
**Text only translation**:
*Human:* Translate the following sentence from English into {lang} language. English sentence is: {sentence}.
*System:* {translation}.
**Image captioning**: *Human:* You are given an image and coordinates of a bounding box as: x1={x1}, y1={y1}, x2={x1+x2}, y2={y1+y2}. You are also provided labels of the objects in the image as: {labels}. Provide a short caption of the object in {lang} language.
*System:* {caption}.

---

Table 2: Different prompt templates for creating task specific fine-tuning data, used in Stage 3 training.

plex instructions across multiple-turns of the conversation. We focus on developing a specialized multimodal translation system in the Stage 3.

**Stage 3: Task-specific Fine-Tuning.** We follow a similar recipe to that of Stage 2 for the (Machine Translation) task-specific fine-tuning and update weights for the projection layer and the LLM while keeping the vision encoder frozen. Here, we experiment with both LoRA style training (Hu et al., 2021; Houlsby et al., 2019) as well as full parameter fine-tuning on the shared task translation data.

## 4 Dataset

In this section, we describe the data resources utilized throughout our experiments.

**Stage 1:** In our initial experiments, we use the LLaVA-Pretrain-LCS-558K data for pre-training our model in Stage 1. However, recent works (Tong et al., 2024) showed that more adapter data is beneficial for the model, such as the 1.2M ShareGPT4V-PT (Chen et al., 2023) image-captioning dataset, which we use in Stage 1 training. We also trans-

lated this data in the 10 Indic languages that our LLM natively supports, using an in-house text Machine Translation system. We sample translations across different languages (including English) in an equal ratio and ensure that PT data limits to 1.2M data points in our final data mix.

**Stage 2:** For the second stage instruction tuning, eliciting visual reasoning abilities, we experiment with both LLaVA-Instruct-150K (Liu et al., 2023b) and LLaVA-1.5-665K (Liu et al., 2024) where we find continued improvements with the 665K version. Similar to pre-training data, we also translated the LLaVA-1.5-665K into multiple languages. Recently released Cauldron dataset (Laurençon et al., 2024) is a collection of 50 academic Vision-language tasks. In our final submission, we also include the translated versions and the original English language based Cauldron apart from the proprietary multi-modal dataset in the training mix. It must be noted that the English only Visual Genome might be a part of this stage's training data through various academic datasets, though not for the translation task.

**Stage 3:** For the Stage 3, we work with the aligned multi-lingual Visual Genome (Krishna et al., 2017) datasets, i.e. Hindi (Parida et al., 2019), Bengali (Sen et al., 2022) and Malayalam (Parida and Bojar, 2021), bundled as part of the shared task. Each row in the dataset consists of the following fields: *i).* MS COCO (Lin et al., 2014) image id *ii).* English utterance *iii).* Translated utterance in Hindi/ Bengali/ Malayalam *iv).* Bounding box of the area in the image that the utterance is based on. While there is also a track for Hausa language (Abdulmumin et al., 2022), we don't include this in our training data. Table 1 provides the statistics of the different versions of the dataset, which we transform into instruct tuning format, similar to Stage 1 and 2 data (see Table 2). To increase the efficacy of our model, we enrich the dataset with the labels of different objects in the image (object tags), similar to (Gupta et al., 2021). We use SOTA (state-of-the-art) YOLOv8 (Varghese and Sambath, 2024) for object detection compared to the prior works, which relied on Faster R-CNN models (Wu et al., 2019; Girshick, 2015). We also calculate the Intersection-over-union (IoU) for the detected and the dataset provided bounding boxes to get the most relevant object tag. However, we found a decreased performance against including the labels of all the detected objects.

| Submission | Hi-Ch | | Hi-Test | | Bn-Ch | | Bn-Test | | Ml-Ch | | Ml-Test | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ |
| SILO NLP | 29.6 | 0.73 | 36.2 | 0.79 | 22.6 | 0.61 | 41.0 | 0.77 | 14.6 | 0.39 | 30.8 | 0.60 |
| NLP Voices | 41.8 | 0.81 | 43.1 | 0.82 | 32.9 | 0.71 | 39.8 | 0.75 | 19.6 | 0.54 | 30.6 | 0.64 |
| Volta | 51.7 | 0.86 | 44.1 | 0.82 | - | - | - | - | - | - | - | - |
| ODIAGEN | 53.6 | 0.86 | 44.6 | 0.83 | **47.8** | **0.82** | **49.2** | **0.8** | 39.7 | 0.75 | 46.6 | 0.75 |
| Ours (leaderboard) | **54.1** | **0.86** | 43.3 | 0.81 | 44.2 | 0.79 | 45.1 | 0.77 | 34.0 | 0.65 | 37.8 | 0.63 |
| Ours† | **55.3** | **0.87** | **44.7** | **0.83** | <u>46.7</u> | <u>0.81</u> | <u>48.1</u> | <u>0.79</u> | **40.6** | **0.75** | **51.7** | **0.88** |

Table 3: English-to-lowres leaderboard scores for Text-only task for Indic languages (Team 007). In the following tables, †denotes the results after submission deadline using the IndicTrans2 evaluation scripts, all the other results are reported using the shared task dashboard.

| Method | Hi-Ch | | Hi-Test | | Bn-Ch | | Bn-Test | | Ml-Ch | | Ml-Test | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ |
| ODIANLP | 0 | 0.04 | 0.8 | 0.06 | - | - | - | - | - | - | - | - |
| NLPHUT | - | - | - | - | - | - | - | - | 0.9 | 0.02 | 0.9 | 0.05 |
| Ours (leaderboard) | **1.3** | **0.13** | **2.8** | **0.18** | **0.4** | **0.04** | **1.8** | **0.11** | 0.3 | 0.04 | 0.9 | **0.06** |

Table 4: English-to-lowres leaderboard scores for Image captioning track. Ours is the only multi-lingual model which can handle all the 3 Indic languages for image captioning.

## 5 Experimental Results and Discussion

This section details our experimental setup and presents the results of our comparative studies.

### 5.1 Implementation

We use HuggingFace Transformers (Wolf et al., 2019) based on PyTorch (Paszke et al., 2019) for our experiments. We consider PALO (Maaz et al., 2024) as a multi-lingual multi-modal baseline and use the code provided with the repository[2]. The shared task provides a leaderboard based on the automatics metrics of BLEU (Papineni et al., 2002) and RIBES (Isozaki et al., 2010). For reporting BLEU, we used the evaluation scripts[3] provided with (Gala et al., 2023b) and the official repository for RIBES[4]. Similar to previous works (Gupta et al., 2021), we also report the results after tokenizing the outputs using indic-tokenizer[5]. Our Stage 1 and Stage 2 tuning follow similar hyperparameters as the LLaVA model (Liu et al., 2023b) unless specified otherwise. For Stage 3 fine-tuning, we conducted multiple experiments for hyperparameter search of learning rate (1e-3, 1e-4, and 1e-5); as well as multiple epochs (1, 2, 3, and 5). We observed rapid over fitting after only one epoch while a learning rate of 1e-4 yielded the highest overall performance. All our further experiments are reported based on this configuration.

[2] https://github.com/mbzuai-oryx/PALO
[3] https://github.com/AI4Bharat/IndicTrans2
[4] https://github.com/nttcslab-nlp/RIBES
[5] https://github.com/ltrc/indic-tokenizer

### 5.2 Results for different tracks

Table 3, 4 and 5 present the results for text-only, image captioning and the Multimodal translation task respectively. For the text-only task, our Chitranuvad model was trained with image data till Stage 2. In Stage 3, we only finetune with text only translations. During inference, we prompt the model with text only translations and dont provide images. Our model achieves SOTA on Hindi and Malayalam Challenge and Test sets while being competitive on the Bengali dataset (see Table 3). We were the only submission which could do image captioning in all the 3 languages (see Table 4). For the MMT task, we achieved SOTA on Hindi Challenge and Malayalam test set while being competitive on the other languages. We also provide cherry-picked system outputs of our best Multimodal LLM in Table 3. From our manual inspection, we saw that our generated translations are better than the ground truth. For example, in the last snippet, our model correctly translates the word 'downhill', which the gold translation fails to capture.

### 5.3 0-shot on the Shared task data

We evaluate the efficacy of our model after Stage 2 as the 0-shot setting, where we don't fine-tune specifically for the shared task translation data. In our preliminary experiments, we only use the English versions of the datasets mentioned in Section 4 for both Adapter tuning (Stage 1) and Instruction tuning (Stage 2). Exceptionally, our Krutrim LLM still retained multi-lingual capabilities, ev-

| Submission | Hi-Ch | | Hi-Test | | Bn-Ch | | Bn-Test | | Ml-Ch | | Ml-Test | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ |
| IIT-P | 37.5 | 0.79 | 42.5 | 0.81 | - | - | - | - | - | - | - | - |
| ODIAGEN | 42.8 | 0.82 | 41.6 | 0.81 | 30.5 | 0.69 | 42.4 | 0.76 | - | - | - | - |
| Volta | 51.6 | 0.86 | 44.6 | 0.82 | - | - | - | - | - | - | - | - |
| BITS-P | 52.1 | 0.85 | 45.0 | 0.83 | 48.7 | 0.83 | 50.6 | 0.81 | 42.2 | 0.76 | 51.9 | 0.80 |
| Ours (leaderboard) | **53.4** | **0.842** | 43.7 | 0.81 | 44.8 | 0.78 | 44.5 | 0.76 | 39.8 | 0.74 | **51.9** | 0.78 |
| Ours† | **54.7** | **0.86** | <u>43.9</u> | **0.83** | <u>46.9</u> | <u>0.81</u> | <u>47.7</u> | <u>0.79</u> | <u>40.3</u> | <u>0.74</u> | **51.9** | **0.93** |

Table 5: English-to-lowres leaderboard Scores for Multimodal translation track across multiple languages (Team 007). †denotes the results after submission deadline using the IndicTrans2 evaluation scripts

| Method | Hi-Ch | | Hi-Test | | Bn-Ch | | Bn-Test | | Ml-Ch | | Ml-Test | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ |
| PALO-7B | 14.8 | 0.585 | 13.3 | 0.567 | 7.9 | 0.469 | 9.6 | 0.464 | 0.1 | 0.001 | 0 | 0 |
| PALO-13B | 15.8 | 0.605 | 14.9 | 0.605 | 6.7 | 0.44 | 7.0 | 0.45 | 0.1 | 0.004 | 0 | 0 |
| Chitranuvad (Eng) | 18.3 | 0.629 | 12.9 | 0.585 | 8.7 | 0.512 | 8.3 | 0.477 | 8.7 | 0.487 | 7.3 | 0.426 |
| Chitranuvad (Eng+Hindi) | 20.0 | **0.698** | 14.8 | 0.653 | 9.4 | 0.537 | 8.9 | 0.494 | 9.2 | 0.511 | 8.6 | 0.466 |
| Chitranuvad (Multilingual) | **25.0** | 0.694 | **19.0** | **0.66** | **11.4** | **0.569** | **9.7** | **0.515** | **12.2** | **0.54** | **10.3** | **0.486** |

Table 6: 0-shot results for Multimodal Machine Translation track as discussed in Section 5.3. *Eng* denotes only English data is used in Stage 1 and 2. *Eng+Hindi* denotes English and Hindi data in Stage 2. As expected, we find the best scores when the training data mix consists of data from the 10 Indic languages.

ident from the scores in Table 6. When we also include Hindi data in the training mix, we find an uplift on the Hindi translation task. Including multi-lingual data in both the stages further showed improvement on all three language translation tasks in the 0-shot setting. We thus use this as the base model in the Stage 3 training. We also evaluate against the open-source baseline of PALO-7/13B (Maaz et al., 2024) in the 0-shot setting. To our surprise, our Chitranuvad model consistently outperforms the 0-shot PALO baseline, even when our model is fine-tuned with English only data in both the stages. We hypothesis that this is because the base LLM Vicuna (Zheng et al., 2024) used in PALO is not inherently multi-lingual in nature.

## 5.4 Other fine-tuning approaches

In this section, we elaborate on key findings with different fine-tuning approaches, with all the results reported in Table 7.

**LoRA vs Full finetuning.** We investigated the efficacy of full fine-tuning versus Low-Rank Adaptation (LoRA) using Visual Genome data. Our experiments (see Table 7) reveal that full fine-tuning consistently outperforms LoRA, i.e. LoRA learns less (Biderman et al., 2024).

**Bi-lingual vs Multi-lingual** For Stage 3 training, we experiment with training specialized models for each language (Hindi and Bengali) compared to multi-lingual setting with a mix of data from all the three Indic languages. We didn't find any im-

provement over multi-lingual model but instead observe catastrophic forgetting (Zhai et al., 2023c; Tong et al., 2024), where the translation abilities of the model in the other languages deteriorate completely. We hypothesize that a mix of multiple languages probably act as regularizaton and enhance the general translation capabilities.

**Do we need second stage training?** Similar to (Tong et al., 2024), we investigate if we even need Stage 2 instruction tuning. We find that our model, if finetuned directly on Visual Genome translation data (i.e. Stage 1 and 3 training only) performed comparable to the previous baselines. Including Stage 2 training provided an uplift in the scores with an added advantage of building a general purpose Multimodal agent.

**Back translation** Back translations, i.e. using the reverse translations have been a popular technique both for data augmentation as well as postprocessing or re-ranking techniques in traditional Machine Translation and Natural Language Generation systems (Sennrich et al., 2015; Li et al., 2015; Agarwal et al., 2018; Edunov, 2018; Graça et al., 2019). This involves re-translating content from the target language to its source language. Thus, apart from the original task of En -> Hi/Bn/Ml, we also included in our training corpus, the task of reverse translation from Hi/Bn/Ml -> En in the Stage 3 training mix. However, in our experiments, we found that this strategy showed decreased performance in terms of automatic metrics.

*Object tags*: person, bowl, oven, dog
*Bbox tag*: dog
*English*: a woman holding a dog

*Hindi GT*: एक स्त्री जो कुत्ता रखती है
*Hindi Translated*: एक महिला एक कुत्ते को पकड़े हुए
*Bengali GT*: একটি মহিলা একটি কুকুর ধরে
*Bengali Translated*: একটি মহিলা একটি কুকুর ধরে
*Malayalam GT*: ഒരു സ്ത്രീ നായയെ പിടിക്കുന്നു
*Malayalam Translated*: ഒരു സ്ത്രീ നായയെ പിടിക്കുന്നു

---

*Object tags*: person, person, person, tennis racket
*Bbox tag*: tennis racket
*English*: A TENNIS RACKET

*Hindi GT*: एक टेनिस रेकेट
*Hindi Translated*: एक टेनिस रैकेट
*Bengali GT*: একটি টেনিস র্যাকেট
*Bengali Translated*: একটি টেনিস র্যাকেট
*Malayalam GT*: ഒരു ടെന്നീസ് റാക്കറ്റ്
*Malayalam Translated*: ഒരു ടെന്നീസ് റാക്കറ്റ്

---

*Object tags*: person, snowboard, backpack
*Bbox tag*: person
*English*: snow on the ground

*Hindi GT*: बर्फ़ जमीन पर
*Hindi Translated*: जमीन पर बर्फ
*Bengali GT*: মাটিতে বরফ
*Bengali Translated*: মাটিতে তুষার
*Malayalam GT*: നിലത്ത് മഞ്ഞ്
*Malayalam Translated*: നിലത്ത് മഞ്ഞ്

---

*Object tags*: dog
*Bbox tag*: dog
*English*: white flower on curtain

*Hindi GT*: पर्दे पर सफेद फूल
*Hindi Translated*: पर्दे पर सफेद फूल
*Bengali GT*: পর্দার উপর সাদা ফুল
*Bengali Translated*: পর্দায় সাদা ফুল
*Malayalam GT*: തിരശ്ശീലയിൽ വെളുത്ത പുഷ്പം
*Malayalam Translated*: വെളുത്ത പുഷ്പം

---

*Object tags*: person, skis
*Bbox tag*: person
*English*: Woman going fast downhill.

*Hindi GT*: तेज गति से जा रही महिला।
*Hindi Translated*: नीचे की ओर तेजी से जा रही महिला।
*Bengali GT*: মহিলা দ্রুত উতরাইয়ের দিকে যাচ্ছে।
*Bengali Translated*: মহিলা দ্রুত নিচে যাচ্ছে।
*Malayalam GT*: വേഗത്തിൽ താഴേക്ക് പോകുന്ന സ്ത്രീ.
*Malayalam Translated*: മലയിറങ്ങുന്ന വേഗത്തിൽ സ്ത്രീ.

Figure 3: English-to-lowres Multimodal Machine Translation track supports translation of source sentence into multiple Indic languages (Hindi, Bengali, Malayalam). We enrich the dataset to include labels of all the identified objects. We show the outputs of our best model which is trained with a mix of multi-lingual data in all the 3 stages.

| Method | Hi-Ch | | Hi-Test | | Bn-Ch | | Bn-Test | | Ml-Ch | | Ml-Test | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ | BLEU ↑ | RIBES ↑ |
| LoRA | 42.1 | 0.721 | 34.5 | 0.770 | 28.3 | 0.69 | 30.4 | 0.669 | 23.2 | 0.61 | 27.0 | 0.601 |
| Bi-lingual (Hi Stage 3) | 53.0 | 0.848 | 43.3 | 0.833 | 0.3 | 0.006 | 0.1 | 0.001 | 0.2 | 0.003 | 0 | 0 |
| Bi-lingual (Bn Stage 3) | 0.3 | 0.005 | 0.1 | 0.001 | 45.4 | 0.797 | 46.6 | 0.781 | 0.3 | 0.003 | 0.1 | 0.001 |
| Only Stage 1, 3 | 53.6 | 0.853 | 43.4 | 0.826 | 45.2 | 0.801 | 46.4 | 0.788 | 38.2 | 0.735 | 50.3 | 0.781 |
| Back Translation mix | 53.8 | 0.856 | 43.6 | 0.828 | 46.0 | 0.806 | 46.8 | 0.792 | 37.5 | 0.729 | 46.3 | 0.738 |

Table 7: Different finetuning strategies for Multimodal Machine Translation track as described in Section 5.4 in the order of discussion. Our Stage 3 full finetuning training performs the best compared to the other training recipes.

# 6 Conclusion

We present Chitranuvad, a multimodal LLM that is adapted for image grounded Machine translation. Our model encodes images using a pre-trained image encoder (Alexey, 2020) and translates the English sentences autoregressively into different Indic languages (Hindi, Bengali, Malayalam) using a pre-trained LLM. Empirically, our model outperforms previous baselines for different tasks. However, we also observed that vision modality had little impact on the translation, echoing the observations from (Grönroos et al., 2018; Lala et al., 2018; Wu et al., 2021; Li et al., 2022).

**Broader Impact:** We believe our work paves way for building next generation assistants which can do multimodal machine translation. We believe these systems can empower different sectors like education, healthcare, banking and financial services, etc. to name a few.

**Limitations and Future Work:** While this work is focused to three Indic languages (Hindi, Bengali, Malayalam), we consider our approach as a first step towards building general purpose multilingual system which can handle various Indic languages. While in our current setup, we freeze the vision encoder during training, recent works have shown that unfreezing the vision encoder with Perceiver Resampler (Jaegle et al., 2021), helps learn better representations (Laurençon et al., 2024; Tong et al., 2024), which we plan to explore in the future.

# Acknowledgements

# References

Idris Abdulmumin, Satya Ranjan Dash, Musa Abdullahi Dawud, Shantipriya Parida, Shamsuddeen Hassan Muhammad, Ibrahim Sa'id Ahmad, Subhadarshi Panda, Ondřej Bojar, Bashir Shehu Galadanci, and Bello Shehu Bello. 2022. Hausa visual genome: A dataset for multi-modal english to hausa machine translation. *arXiv preprint arXiv:2205.01133*.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Shubham Agarwal, Marc Dymetman, and Eric Gaussier. 2018. Char2char generation with reranking for the e2e nlg challenge. *arXiv preprint arXiv:1811.05826*.

Dosovitskiy Alexey. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv: 2010.11929*.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Abhinand Balachandran. 2023. Tamil-llama: A new tamil language model based on llama 2.

Abhijit Bendale, Michael Sapienza, Steven Ripplinger, Simon Gibbs, Jaewon Lee, and Pranav Mistry. 2024. Sutra: Scalable multilingual language model architecture. *arXiv preprint arXiv:2405.06694*.

Dan Biderman, Jose Gonzalez Ortiz, Jacob Portes, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, et al. 2024. Lora learns less and forgets less. *arXiv preprint arXiv:2405.09673*.

Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, et al. 2020. Experience grounds language. *CoRR abs/2004.10151*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child,

Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Christian Buck, Kenneth Heafield, and Bas Van Ooyen. 2014. N-gram counts and language models from the common crawl. In *LREC*, volume 2, page 4.

Ozan Caglayan, Loïc Barrault, and Fethi Bougares. 2016. Multimodal attention for neural machine translation. *arXiv preprint arXiv:1609.03976*.

Ozan Caglayan, Julia Ive, Veneta Haralampieva, Pranava Madhyastha, Loïc Barrault, and Lucia Specia. 2020. Simultaneous machine translation with visual context. In *CoRR abs/2009.07310*.

Iacer Calixto and Qun Liu. 2017. Incorporating global visual features into attention-based neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 992–1003, Copenhagen, Denmark. Association for Computational Linguistics.

Haishun Chen, Ying Wang, Xin Yang, and Jie Li. 2021. Captioning transformer with scene graph guiding. In *2021 IEEE international conference on image processing (ICIP)*, pages 2538–2542. IEEE.

Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023. Sharegpt4v: Improving large multimodal models with better captions. *arXiv preprint arXiv:2311.12793*.

Kyunghyun Cho. 2014. Learning phrase representations using rnn encoder–decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.

Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh M Khapra, and Pratyush Kumar. 2021. Indicbart: A pre-trained model for indic natural language generation. *arXiv preprint arXiv:2109.02903*.

Amulya Dash, Hrithik Raj Gupta, and Yashvardhan Sharma. 2023. Bits-p at wat 2023: Improving indic language multimodal translation by image augmentation using diffusion models. In *Proceedings of the 10th Workshop on Asian Translation*, pages 41–45.

Jean-Benoit Delbrouck and Stéphane Dupont. 2017. An empirical study on the effectiveness of images in multimodal neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 910–919, Copenhagen, Denmark. Association for Computational Linguistics.

Sergey Edunov. 2018. Understanding back-translation at scale. *arXiv preprint arXiv:1808.09381*.

Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30k: Multilingual english-german image descriptions. *arXiv preprint arXiv:1605.00459*.

Desmond Elliott and Ákos Kádár. 2017. Imagination improves multimodal translation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 130–141, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Jerry A Fodor. 1975. *The language of thought*, volume 5. Harvard university press.

Baban Gain, Dibyanayan Bandyopadhyay, and Asif Ekbal. 2021. Iitp at wat 2021: System description for english-hindi multimodal translation task. *arXiv preprint arXiv:2107.01656*.

Jay Gala, Pranjal A Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, et al. 2023a. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307*.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023b. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*.

Jay Gala, Thanmay Jayakumar, Jaavid Aktar Husain, Aswanth Kumar M, Mohammed Safi Ur Rahman Khan, Diptesh Kanojia, Ratish Puduppully, Mitesh M. Khapra, Raj Dabre, Rudra Murthy, and Anoop Kunchukuttan. 2024. Airavata: Introducing hindi instruction-tuned llm.

Lizhao Gao, Bo Wang, and Wenmin Wang. 2018. Image captioning with scene-graph based semantic concepts. In *Proceedings of the 2018 10th international conference on machine learning and computing*, pages 225–229.

Ross Girshick. 2015. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448.

Miguel Graça, Yunsu Kim, Julian Schamper, Shahram Khadivi, and Hermann Ney. 2019. Generalizing back-translation in neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 45–52, Florence, Italy. Association for Computational Linguistics.

Stig-Arne Grönroos, Benoit Huet, Mikko Kurimo, Jorma Laaksonen, Bernard Merialdo, Phu Pham, Mats Sjöberg, Umut Sulubacak, Jörg Tiedemann, Raphael Troncy, et al. 2018. The memad submission to the wmt18 multimodal translation task. *arXiv preprint arXiv:1808.10802*.

Junjun Guo, Junjie Ye, Yan Xiang, and Zhengtao Yu. 2023. Layer-level progressive transformer with modality difference awareness for multi-modal neural machine translation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.

Kshitij Gupta, Devansh Gautam, and Radhika Mamidi. 2021. Vita: Visual-linguistic translation by aligning object tags. *arXiv preprint arXiv:2106.00250*.

Stevan Harnad. 1990. The symbol grounding problem. *Physica D: Nonlinear Phenomena*.

Stevan Harnad. 2003. Minds, machines and turing. In *The Turing Test*, pages 253–273. Springer.

Stevan Harnad. 2005. To cognize is to categorize: Cognition is categorization. In *Handbook of categorization in cognitive science*, pages 21–54. Elsevier.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Yuheng Hu, Lydia Manikonda, and Subbarao Kambhampati. 2014. What we instagram: A first analysis of instagram photo content and user types. In *8th International Conference on Weblogs and Social Media, ICWSM 2014*.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 944–952.

Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. 2021. Perceiver: General perception with iterative attention. In *International conference on machine learning*, pages 4651–4664. PMLR.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Guneet Singh Kohli, Shantipriya Parida, Sambit Sekhar, Samirit Saha, Nipun B Nair, Parul Agarwal, Sonal Khosla, Kusumlata Patiyal, and Debasish Dhal. 2023. Building a llama2-finetuned llm for odia language utilizing domain knowledge instruction set.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73.

Telugu Labs. 2023. Navarsa: Indic llms based on gemmma.

Chiraag Lala, Pranava Swaroop Madhyastha, Carolina Scarton, and Lucia Specia. 2018. Sheffield submissions for wmt18 multimodal translation shared task. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 624–631.

Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. 2024. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*.

M Lewis. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

Bei Li, Chuanhao Lv, Zefan Zhou, Tao Zhou, Tong Xiao, Anxiang Ma, and JingBo Zhu. 2022. On vision features in multimodal machine translation. *arXiv preprint arXiv:2203.09173*.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning. *Advances in neural information processing systems*, 36.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning.

Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. 2024a. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*.

Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. 2024b. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26439–26455.

Muhammad Maaz, Hanoona Rasheed, Abdelrahman Shaker, Salman Khan, Hisham Cholakal, Rao M Anwer, Tim Baldwin, Michael Felsberg, and Fahad S Khan. 2024. Palo: A polyglot large multimodal model for 5b people. *arXiv preprint arXiv:2402.14818*.

Raymond J Mooney. 2008. Learning to connect language and perception. In *AAAI*.

Toshiaki Nakazawa, Chenchen Ding, Raj Dabre, Anoop Kunchukuttan, Nobushige Doi, Yusuke Oda, Ondřej Bojar, Shantipriya Parida, Isao Goto, and Hidaya Mino. 2019. Proceedings of the 6th workshop on asian translation. In *Proceedings of the 6th Workshop on Asian Translation*.

Toshiaki Nakazawa, Kazutaka Kinugawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida, Makoto Morishita, Ondrej Bojar, Akiko Eriguchi, Yusuke Oda, Akiko Eriguchi, Chenhui Chu, and Sadao Kurohashi, editors. 2023. *Proceedings of the 10th Workshop on Asian Translation*. Asia-Pacific Association for Machine Translation, Macau SAR, China.

Toshiaki Nakazawa, Kazutaka Kinugawa, Hideya Mino, Isao Goto, Raj Dabre, Shohei Higashiyama, Shantipriya Parida, Makoto Morishita, Ondrej Bojar, Akiko Eriguchi, Yusuke Oda, Akiko Eriguchi, Chenhui Chu, et al., editors. 2022. *Proceedings of the 9th Workshop on Asian Translation*. International Conference on Computational Linguistics, Gyeongju, Republic of Korea.

Toshiaki Nakazawa, Hideki Nakayama, Chenchen Ding, Raj Dabre, Anoop Kunchukuttan, Win Pa Pa, Ondřej Bojar, Shantipriya Parida, Isao Goto, Hidaya Mino, Hiroshi Manabe, Katsuhito Sudoh, Sadao Kurohashi, and Pushpak Bhattacharyya, editors. 2020. *Proceedings of the 7th Workshop on Asian Translation*. Association for Computational Linguistics, Suzhou, China.

Toshiaki Nakazawa, Hideki Nakayama, Isao Goto, Hideya Mino, Chenchen Ding, Raj Dabre, Anoop Kunchukuttan, Shohei Higashiyama, Hiroshi Manabe, Win Pa Pa, Shantipriya Parida, Ondřej Bojar, Chenhui Chu, Akiko Eriguchi, Kaori Abe, Yusuke Oda, Katsuhito Sudoh, Sadao Kurohashi, and Pushpak Bhattacharyya, editors. 2021. *Proceedings of the 8th Workshop on Asian Translation (WAT2021)*. Association for Computational Linguistics, Online.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Shantipriya Parida and Ondřej Bojar. 2021. Malayalam visual genome 1.0. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.

Shantipriya Parida, Ondřej Bojar, and Satya Ranjan Dash. 2019. Hindi Visual Genome: A Dataset for Multimodal English-to-Hindi Machine Translation. *Computación y Sistemas*, 23(4):1499–1505. Presented at CICLing 2019, La Rochelle, France.

Shantipriya Parida, Subhadarshi Panda, Stig-Arne Grönroos, Mark Granroth-Wilding, and Mika Koistinen. 2022. Silo nlp's participation at wat2022. *arXiv preprint arXiv:2208.01296*.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.

Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Lee Rainie, Joanna Brenner, and Kristen Purcell. 2012. Photos and videos as social currency online. *Pew Internet & American Life Project*.

Arghyadeep Sen, Shantipriya Parida, Ketan Kotwal, Subhadarshi Panda, Ondřej Bojar, and Satya Ranjan Dash. 2022. Bengali visual genome: A multimodal dataset for machine translation and image captioning. In *Intelligent Data Engineering and Analytics*, pages 63–70. Springer.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Improving neural machine translation models with monolingual data. *arXiv preprint arXiv:1511.06709*.

Sk Shahid, Guneet Singh Kohli, Sambit Sekhar, Debasish Dhal, Adit Sharma, Shubhendra Kushwaha, Shantipriya Parida, Stig-Arne Grönroos, and Satya Ranjan Dash. 2023. Odiagenai's participation at wat2023. In *Proceedings of the 10th Workshop on Asian Translation*, pages 46–52.

Attend Show. 2015. Tell: Neural image caption generation with visual attention kelvin xu. *Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio arXiv (2015-02-10) https://arxiv. org/abs/1502.03044 v3*.

Lucia Specia, Stella Frank, Khalil Sima'an, and Desmond Elliott. 2016. A shared task on multimodal machine translation and crosslingual image description. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 543–553.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Chameleon Team. 2024a. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

Krutrim Team. 2024b. Krutrim LLM: Multilingual foundational model for over a billion people. *Under Review*.

Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. 2024. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Emiel van Miltenburg, Merel van de Kerkhof, Ruud Koolen, Martijn Goudbeek, and Emiel Krahmer. 2019. On task effects in NLG corpus elicitation: A replication study using mixed effects modeling. In *INLG*.

Rejin Varghese and M Sambath. 2024. Yolov8: A novel object detection algorithm with enhanced performance and robustness. In *2024 International Conference on Advances in Data Engineering and Intelligent Computing Systems (ADICS)*, pages 1–6. IEEE.

Ashish Vaswani. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762*.

Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. 2023. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.

Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. 2019. Detectron2. https://github.com/facebookresearch/detectron2.

Zhiyong Wu, Lingpeng Kong, Wei Bi, Xiang Li, and Ben Kao. 2021. Good for misconceived reasons: An empirical revisiting on the need for visual context in multimodal machine translation. In *ACL*, pages 6153–6166, Online. Association for Computational Linguistics.

Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, et al. 2024. xgen-mm (blip-3): A family of open large multimodal models. *arXiv preprint arXiv:2408.08872*.

Shaowei Yao and Xiaojun Wan. 2020. Multimodal transformer for multimodal machine translation. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 4346–4350.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023a. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023b. Sigmoid loss for language image pre-training.

Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. 2023c. Investigating the catastrophic forgetting in multimodal large language models. *arXiv preprint arXiv:2309.10313*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.

# Brotherhood at WMT 2024: Leveraging LLM-Generated Contextual Conversations for Cross-Lingual Image Captioning

**Siddharth Betala**
betalas5@gmail.com

**Ishan Chokshi**
ishan.c1671@gmail.com

## Abstract

In this paper, we describe our system under the team name *Brotherhood* for the English-to-Lowres Multi-Modal Translation Task. We participate in the multi-modal translation tasks for English-Hindi, English-Hausa, English-Bengali, and English-Malayalam language pairs. We present a method leveraging multimodal Large Language Models (LLMs), specifically GPT-4o and Claude 3.5 Sonnet, to enhance cross-lingual image captioning without traditional training or fine-tuning.

Our approach utilizes instruction-tuned prompting to generate rich, contextual conversations about cropped images, using their English captions as additional context. These synthetic conversations are then translated into the target languages. Finally, we employ a weighted prompting strategy, balancing the original English caption with the translated conversation to generate captions in the target language.

This method achieved competitive results, scoring 37.90 BLEU on the English-Hindi Challenge Set and ranking 1st and 2nd for English-Hausa on the Challenge and Evaluation Leaderboards, respectively. We conduct additional experiments on a subset of 250 images, exploring the trade-offs between BLEU scores and semantic similarity across various weighting schemes.

## 1 Introduction

Machine translation (MT) is a classic subfield in NLP that investigates the usage of computer software to translate text or speech from one language to another without human involvement (Yang et al., 2020). Machine translation (MT) has seen remarkable advancements in recent years, primarily due to the success of neural approaches (Bahdanau, 2014; Vaswani, 2017). However, these improvements have been predominantly observed in high-resource language pairs, leaving low-resource languages significantly behind (Sennrich and Zhang, 2019; Costa-jussà et al., 2022). The challenges in low-resource MT are multifaceted, including limited parallel corpora, lack of linguistic diversity in training data, and the absence of specialized tools and resources.

One promising direction to address these challenges is the incorporation of visual information into the translation process, known as multimodal machine translation (MMT) (Elliott et al., 2016; Specia et al., 2016; Calixto et al., 2017). The underlying hypothesis is that visual context can provide crucial disambiguating cues, especially for languages with limited textual resources. This approach aligns with the human cognitive process of language understanding, which often relies on multiple sensory inputs (Beinborn et al., 2018).

The Workshop on Machine Translation (WMT) 2024 has presented a shared task on English-to-Low-Resource Multi-Modal Translation, focusing on Hindi (Parida et al., 2019), Bengali (Sen et al., 2022), Malayalam[1], and Hausa (Abdulmumin et al., 2022). This task utilizes variants of the Visual Genome (Krishna et al., 2017) dataset, adapted for these target languages. While these datasets provide a valuable resource for research, they also present unique challenges. First, the quality of translations in low-resource languages can be inconsistent (see Table 1), potentially introducing noise into the training process. Second, the limited size of these datasets makes it difficult to train robust neural models without overfitting. Lastly, the cultural and linguistic nuances of these languages may not be fully

---

[1]https://ufal.mff.cuni.cz/
malayalam-visual-genome

852

**Source English Caption:**
soap is in the dish
**Target Hindi Caption:**
साबुन पकवान में है

Table 1: Illustration of a translation ambiguity in the dataset. The English caption for the bounded region **"soap is in the dish"** is mistranslated into Hindi as "साबुन पकवान में है" (**"soap is in the food"**). The Hindi word "पकवान" typically means **"food dish"**, whereas the visual context clearly indicates that **"dish"** refers to a **soap holder**. This example also highlights the importance of visual context in resolving lexical ambiguities in multimodal translation tasks.

captured in direct translations (Hershcovich et al., 2022) of English captions.

Existing approaches (Calixto and Liu, 2017) to MMT often involve training complex neural architectures on large parallel corpora with paired images. However, such methods may struggle with the limited and potentially noisy data available for low-resource languages. Moreover, these approaches often fail to leverage the rich semantic understanding capabilities of recent large language models (LLMs) along with their proficiency in visual (Liu et al., 2024; Radford et al., 2021) and multilingual (Üstün et al., 2024; Workshop et al., 2022; Dubey et al., 2024; Touvron et al., 2023; Jiang et al., 2023) understanding gained through training on large corpora of data across domains.

In this paper, we propose a pipeline that addresses these challenges by leveraging multi-modal LLMs, specifically GPT-4o[2] and Claude 3.5 Sonnet[3], to enhance cross-lingual image captioning. Our approach uses

instruction-tuned prompting to generate rich, contextual conversations about images using the English captions as context along with the image. These synthetic conversations, comprising detailed descriptions, simple question-answer (QA) pairs, and complex reasoning question-answer pairs, are then translated into the target languages and used to inform the final caption generation in the target language. This method allows us to:

- Mitigate the impact of limited and noisy training data by generating synthetic, high-quality contextual information.

- Leverage the advanced reasoning capabilities of LLMs to provide culturally and linguistically nuanced translations.

- Explore the balance between source fidelity and enhanced description through a weighted prompting strategy.

The main contributions of our work are:

- A pipeline for low-resource MMT that requires no traditional training or fine-tuning.

- A weighted prompting mechanism that calibrates between source caption fidelity and LLM-generated contextual information, facilitating a nuanced balance of translation accuracy, caption diversity, and exhaustive visual description coverage.

- A framework for dataset enrichment through the generation of detailed descriptions and complex reasoning QA pairs to augment existing multimodal datasets.

- Empirical analysis of LLMs' capabilities in direct translation and target language summarization, providing insights into their potential for low-resource languages.

## 2  Dataset

We utilized only the datasets specified by the organizers for the related tasks. However, our use of the GPT-4o and Claude 3.5 Sonnet models places our submissions in the unconstrained track. The provided datasets contain captions in English and the target language, describing

---

[2]gpt-4o-2024-08-06: https://platform.openai.com/docs/models/gpt-4o.

[3]Claude 3.5 Sonnet: https://www.anthropic.com/api.

rectangular regions in images of various scenes. The task involves generating captions in the target language using either the text, the image, or both. Across all languages, the training set consists of 29,000 examples. The dataset is complemented by three test sets: development (D-Test), evaluation (EV-Test), and challenge (CH-Test). Our submissions were evaluated on the EV-Test and CH-Test sets, which contain approximately 1,600 and 1,400 examples, respectively. The development set comprises around 1,000 examples.

Table 2 shows the parallel corpus statistics across the various languages. Table 3 shows that data sources of datasets for each task.

## 3 Methodology

The overall pipeline of our approach is shown in Figure 1.

### 3.1 Preprocessing

For the text data, all utterances are converted to lowercase, and punctuation is removed. The dataset includes images of complete scenes along with the coordinates of the bottom-left corner and the dimensions of the rectangular region corresponding to each caption. This information is used to crop the relevant rectangular regions from the images. Since these images are later used as part of prompts for LLMs, base64 encodings of all images in the EV and CH sets are generated.

### 3.2 Multi-Model Context Generation in English with a Fusion approach

In this step , we leverage the capabilities of two large language models (LLMs) - GPT-4o and Claude 3.5 Sonnet - to generate rich, contextual conversations about the input image and its associated English caption. This process involves two key stages: individual LLM processing and conversation fusion.

We separately prompt GPT-4o and Claude 3.5 using Prompt-1 from Table 4. The format of the conversation and prompt design is inspired by an example prompt from Liu et al. (2024). Both models are given the same input: the cropped image and its English caption. This parallel processing allows us to capitalize on the unique strengths of each model.

After obtaining separate conversations from GPT-4o and Claude 3.5, we employ a fusion



Figure 1: Overview of the system pipeline. Prompt-1, detailed in Table 4, and Prompt-2, detailed in Table 5, respectively. The weight x% must be specified during querying.

model, implemented using the GPT-4o API,

854

| Set | Sentences | English | Hindi | Malayalam | Bengali | Hausa |
|---|---|---|---|---|---|---|
| **Train** | 28930 | 143164 | 145448 | 107126 | 113978 | 140981 |
| **D-Test** | 998 | 4922 | 4978 | 3619 | 3936 | 4857 |
| **E-Test** | 1595 | 7853 | 7852 | 5689 | 6408 | 7736 |
| **C-Test** | 1400 | 8186 | 8639 | 6044 | 6657 | 8752 |

Table 2: Parallel Corpus Token statistics for each dataset split across different languages.

| Task | Source |
|---|---|
| English→Hindi | HindiVisualGenome1.1 (Parida et al., 2019) |
| English→Malayalam | MalayalamVisualGenome1.0[4] |
| English→Bengali | BengaliVisualGenome1.0 (Sen et al., 2022) |
| English→Hausa | HausaVisualGenome1.0 (Abdulmumin et al., 2022) |

Table 3: Tasks and their corresponding dataset sources.

to integrate these conversations. This fusion process is designed to create a single, coherent dialogue that encompasses three response types: **(1) Short QA pairs, (2) Detailed descriptions, and (3) Complex reasoning QA**. The prompt has been designed to generate such types of responses to ensure that the responses can capture the context effectively. The fusion prompt is carefully crafted to ensure that the final conversation retains the most relevant and insightful elements from both initial conversations.

This fusion approach is aimed at using complementary strengths of the models for error mitigation and rich context consolidation. Recent work has shown the advantage of such fusion and ensembling approaches in mitigating hallucination across tasks such as machine translation, definition modeling, and paraphrase generation (Mehta et al., 2024). Additionally, such an approach has a flexibility advantage as it allows for future integration of additional LLMs or specialized models.

### 3.3 Translation of Context to Target Languages

This stage involves translating the rich contextual information generated in English to the target languages. This step is essential for preserving the nuanced understanding developed in the previous stages while adapting it to the linguistic and cultural context of each target language. For the translation of the generated conversations, we employ different approaches based on the target language:

- **Hindi, Bengali, and Malayalam:** For these Indic languages, we utilize the IndicTrans2 (Gala et al., 2023) model. This state-of-the-art translation model is specifically designed for Indian languages and has demonstrated strong performance in multilingual translation tasks.

- **Hausa**: Due to the limited availability of specialized translation models for Hausa, we leverage the GPT-4o's translation capabilities.

### 3.4 Weighted Prompt-Based Caption Generation

The final stage of our pipeline employs a weighted prompting strategy to generate the target language caption, balancing fidelity to the original English caption with the rich contextual information derived from our LLM-generated conversations. We utilize GPT-4o API for this crucial step, employing a carefully crafted prompt (Prompt-2, detailed in Table 5) that takes two primary inputs along with the weight value:

1. The original English caption (weight: 100-x%)

2. The translated conversation in the target language (weight: x%)

The weighting mechanism allows us to control the influence of each input on the final caption. This approach offers flexibility in balancing between direct translation fidelity and contextual enrichment. For our submissions we provide equal weightage to the given English caption and the generated conversation in the target language.

## 4 Results

The BLEU score (Papineni et al., 2002) serves as the primary metric for evaluating model performance on the leaderboard, complemented

Table 4: Prompt-1: This illustrates the prompt construction process for generating the conversation response using an AI visual assistant. The code snippet generates a prompt for GPT-4 to create a conversation based on an image and caption. The conversation includes three response types: a short Q&A, a detailed description, and complex reasoning. The generated response is capped to a maximum of 1024 tokens for this step.

Table 5: Prompt-2: This prompt guides the generation of a caption in the target language by balancing the original English caption with additional context, weighted according to specified ratios. The generated caption should reflect the natural fluency of the target language and match the original caption's length. The generated response is capped to a maximum of 300 tokens for this step.

by the RIBES metric (Isozaki et al., 2010) for a more comprehensive assessment. Table 6 presents our results across all language pairs and datasets. Notably, our approach achieves significant success in the English-Hausa task, securing the 1st position on the Challenge Leaderboard and the 2nd position on the Evaluation Leaderboard. For the English-Hindi task, we demonstrate competitive performance, obtaining strong BLEU and RIBES scores on both the evaluation and challenge sets. Our method, which requires no training, yields competitive results across all 8 submissions, underscoring its effectiveness and versatility. However, we observe relatively limited performance in the Bengali, Hausa, and Malayalam tasks. This performance discrepancy likely reflects the current limitations of state-of-the-art LLMs, whose training data may not adequately represent low-resource languages and their unique semantic characteristics.

## 4.1 Weight Tuning Analysis and Performance Metrics

To thoroughly evaluate our weighted prompt-based approach, we conducted extensive experiments across different weight combinations for each target language. This section presents our findings and analyzes the impact of various weight distributions on caption quality and semantic preservation. We employed three pri-

| System and Task | BLEU | RIBES | Position |
|---|---|---|---|
| **English→Hindi** | | | |
| MMEVMM22en-hi | 29.70 | 0.725 | 5th |
| MMCHMM22en-hi | 37.90 | 0.796 | 3rd |
| **English→Malayalam** | | | |
| MMEVMM22en-ml | 15.10 | 0.411 | 4th |
| MMCHMM22en-ml | 13.60 | 0.428 | 3rd |
| **English→Bengali** | | | |
| MMEVMM22en-bn | 22.10 | 0.575 | 5th |
| MMCHMM22en-bn | 21.70 | 0.644 | 4th |
| **English→Hausa** | | | |
| MMEVMM22en-ha | 17.70 | 0.580 | 2nd |
| MMCHMM22en-ha | 21.10 | 0.637 | 1st |

Table 6: Summary of results for various English-to-target language multimodal tasks. The table shows BLEU, RIBES scores, and positions for different tasks.

| Weight | BLEU | Sem. Sim. | Norm. Sem. |
|---|---|---|---|
| 0 | 28.56 | 0.9238 | 0.9738 |
| 10 | 23.10 | 0.8978 | 0.9715 |
| 20 | 21.46 | 0.8866 | 0.9726 |
| 30 | 18.73 | 0.8504 | 0.9723 |
| 40 | 16.52 | 0.8381 | 0.9800 |
| 50 | 16.51 | 0.8503 | 0.9747 |
| 60 | 13.71 | 0.7972 | 0.9850 |
| 70 | 13.01 | 0.8103 | 0.9822 |
| 80 | 11.57 | 0.8066 | 0.9784 |
| 90 | 11.26 | 0.8158 | 0.9761 |
| 100 | 10.22 | 0.7893 | 0.9722 |

Table 7: Results for Hindi (hi) with varying context weights. The table shows Average BLEU, Semantic Similarity (hi-hi), and Normalized Semantic Similarity (en-hi).

mary metrics to assess the performance of our model under different weight configurations:

| Weight | BLEU | Sem. Sim. | Norm. Sem. |
|--------|------|-----------|------------|
| 0 | 12.24 | 0.7575 | 0.9380 |
| 10 | 16.47 | 0.7379 | 0.9426 |
| 20 | 11.67 | 0.6992 | 0.9215 |
| 30 | 14.29 | 0.6189 | 0.9437 |
| 40 | 10.06 | 0.5780 | 0.9532 |
| 50 | 8.83 | 0.5998 | 0.9611 |
| 60 | 8.43 | 0.5446 | 0.9570 |
| 70 | 7.10 | 0.5091 | 0.9559 |
| 80 | 8.50 | 0.5460 | 0.9617 |
| 90 | 8.81 | 0.5185 | 0.9583 |
| 100 | 6.43 | 0.5091 | 0.9474 |

Table 8: Results for Malayalam (ml) with varying context weights. The table shows Average BLEU, Semantic Similarity (ml-ml), and Normalized Semantic Similarity (en-ml).

| Weight | BLEU | Sem. Sim. | Norm. Sem. |
|--------|------|-----------|------------|
| 0 | 24.23 | 0.9348 | 0.9617 |
| 10 | 20.07 | 0.9266 | 0.9699 |
| 20 | 17.21 | 0.9072 | 0.9677 |
| 30 | 14.53 | 0.8957 | 0.9811 |
| 40 | 14.98 | 0.8607 | 0.9681 |
| 50 | 13.92 | 0.8670 | 0.9653 |
| 60 | 10.11 | 0.8319 | 0.9829 |
| 70 | 10.29 | 0.8380 | 0.9925 |
| 80 | 10.23 | 0.8545 | 0.9827 |
| 90 | 11.24 | 0.8498 | 0.9789 |
| 100 | 9.35 | 0.8486 | 0.9653 |

Table 9: Results for Bengali (bn) with varying context weights. The table shows Average BLEU, Semantic Similarity (bn-bn), and Normalized Semantic Similarity (en-bn).

| Weight | BLEU | Sem. Sim. | Norm. Sem. |
|--------|------|-----------|------------|
| 0 | 51.58 | 0.6187 | 0.9126 |
| 10 | 46.38 | 0.5599 | 0.9629 |
| 20 | 45.36 | 0.5358 | 0.9505 |
| 30 | 42.23 | 0.5031 | 0.9598 |
| 40 | 37.58 | 0.4829 | 0.9701 |
| 50 | 36.14 | 0.4712 | 0.9556 |
| 60 | 30.48 | 0.4366 | 0.9667 |
| 70 | 28.68 | 0.4344 | 0.9708 |
| 80 | 29.20 | 0.4407 | 0.9757 |
| 90 | 29.41 | 0.4361 | 0.9642 |
| 100 | 32.01 | 0.4442 | 0.9738 |

Table 10: Results for Hausa (ha) with varying context weights. The table shows Average BLEU, Semantic Similarity (ha-ha), and Normalized Semantic Similarity (en-ha).

- **BLEU Score:** Bilingual Evaluation Understudy (BLEU) (Papineni et al., 2002) measures the similarity between the generated caption and the reference caption. It provides a quantitative measure of translation quality.

- **Semantic Similarity (Sem. Sim.):** We use cosine similarity between sentence embeddings to measure the semantic closeness of the generated caption to the reference caption in the target language. This metric is calculated using the SentenceTransformer's SentenceBERT (Reimers and Gurevych, 2019) model 'distiluse-base-multilingual-cased-v1', which provides multilingual sentence embeddings.

- **Normalized Semantic Similarity (Norm. Sem.):** This metric compares the semantic similarity of the generated caption to the English source with that of the reference translation to the English source as the ratio of the two. It helps assess how well the generated caption preserves the meaning of the original English caption relative to the reference translation.

To ensure a comprehensive evaluation of our approach across various scenarios, we conducted an in-depth analysis on a diverse subset of the corpus. This subset comprises 250 image-caption pairs, randomly selected from the training, development, evaluation, and challenge sets. Tables 7, 8, 9, and 10 present the results for Hindi, Malayalam, Bengali, and Hausa, respectively. The 'Weight' column represents the percentage weight given to the translated conversation, with the complementary weight assigned to the original English caption. The reported metrics are averaged over all 250 data points for each weight configuration.

Generally, BLEU scores decrease as more weight is given to the translated conversation. This suggests that higher weights on the original caption produce translations more closely aligned with the reference.

The semantic similarity between the generated and reference captions in the target language tends to decrease with increasing weight on the translated conversation. This indicates that while the generated captions may become more descriptive, they may deviate from the reference in terms of semantic content.

Interestingly, the normalized semantic similarity remains relatively stable across weight distributions. This suggests that our approach consistently preserves the semantic relation-

ship between the English source and the generated caption, regardless of the weight distribution.

While the highest BLEU scores are generally achieved with lower weights on the translated conversation, there's a trade-off between BLEU score and the richness of the generated caption. A balanced approach (e.g., 50-50 weighting) often provides a good compromise between translation accuracy and contextual enrichment.

Notably, setting a 100% weight for the original English caption allows us to evaluate GPT-4's zero-shot cross-lingual transfer capabilities in direct translation tasks from English to Hindi, Bengali, Malayalam, and Hausa. Conversely, assigning a 100% weight to the additional context—which consists of the translated LLM-generated conversation in the target language—enables us to assess the model's abstractive summarization abilities in these non-English languages. This analysis provides insights into the models' multilingual competence and their capacity for language understanding and generation across diverse linguistic contexts.

## 5 Conclusion

Key strengths of our approach include its training-free nature, which avoids propagating errors from potentially flawed datasets, and its flexibility in balancing source fidelity with enhanced descriptiveness through a weighted prompting strategy. The method's multilingual capability and rich context generation offer promising avenues for dataset enrichment and improvement in low-resource languages. In Section 4.1, we demonstrated how our weighted prompting strategy serves as a probe for assessing LLMs' capabilities in zero-shot cross-lingual transfer for direct translation tasks, as well as their abstractive summarization abilities in low-resource target languages. However, we acknowledge limitations such as reliance on LLM APIs, potential for hallucination, and computational intensity. The challenge of evaluating enhanced descriptions with traditional metrics like BLEU also highlights the need for more comprehensive evaluation methods. Future work should focus on:

- Conducting human evaluations to better assess caption quality and appropriateness.

- Analyzing specific cases of significant improvements or detractions from original captions.

- Exploring applications in dataset error correction and enhancement.

- Investigating performance across diverse image types and caption complexities.

By addressing these areas, we aim to further refine and expand the capabilities of our approach, potentially leading to more robust and versatile multimodal translation systems. This work represents a step towards bridging the gap between high-resource and low-resource languages in multimodal machine translation, opening new possibilities for cross-lingual image understanding and dataset enrichment.

## References

Idris Abdulmumin, Satya Ranjan Dash, Musa Abdullahi Dawud, Shantipriya Parida, Shamsuddeen Muhammad, Ibrahim Sa'id Ahmad, Subhadarshi Panda, Ondřej Bojar, Bashir Shehu Galadanci, and Bello Shehu Bello. 2022. Hausa visual genome: A dataset for multi-modal English to Hausa machine translation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6471–6479, Marseille, France. European Language Resources Association.

Dzmitry Bahdanau. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Lisa Beinborn, Teresa Botschen, and Iryna Gurevych. 2018. Multimodal grounding for language processing. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2325–2339, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Iacer Calixto and Qun Liu. 2017. Incorporating global visual features into attention-based neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 992–1003, Copenhagen, Denmark. Association for Computational Linguistics.

Iacer Calixto, Qun Liu, and Nick Campbell. 2017. Doubly-attentive decoder for multi-modal neural machine translation. In *Proceedings of the*

*55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1913–1924, Vancouver, Canada. Association for Computational Linguistics.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30k: Multilingual english-german image descriptions. *arXiv preprint arXiv:1605.00459*.

Jay Gala, Pranjal A Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, et al. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307*.

Daniel Hershcovich, Stella Frank, Heather Lent, Miryam de Lhoneux, Mostafa Abdou, Stephanie Brandl, Emanuele Bugliarello, Laura Cabello Piqueras, Ilias Chalkidis, Ruixiang Cui, Constanza Fierro, Katerina Margatina, Phillip Rust, and Anders Søgaard. 2022. Challenges and strategies in cross-cultural NLP. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6997–7013, Dublin, Ireland. Association for Computational Linguistics.

Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic evaluation of translation quality for distant language pairs. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 944–952.

AQ Jiang, A Sablayrolles, A Mensch, C Bamford, DS Chaplot, D de las Casas, F Bressand, G Lengyel, G Lample, L Saulnier, et al. 2023. Mistral 7b (2023). *arXiv preprint arXiv:2310.06825*.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.

Rahul Mehta, Andrew Hoblitzell, Jack O'keefe, Hyeju Jang, and Vasudeva Varma. 2024. Halunlp at semeval-2024 task 6: Metacheckgpt-a multi-task hallucination detection using llm uncertainty and meta-models. In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, pages 342–348.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Shantipriya Parida, Ondřej Bojar, and Satya Ranjan Dash. 2019. Hindi visual genome: A dataset for multi-modal english to hindi machine translation. *Computación y Sistemas*, 23(4):1499–1505.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Nils Reimers and Iryna Gurevych. 2019. Sentencebert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Arghyadeep Sen, Shantipriya Parida, Ketan Kotwal, Subhadarshi Panda, Ondřej Bojar, and Satya Ranjan Dash. 2022. Bengali visual genome: A multimodal dataset for machine translation and image captioning. In *Intelligent Data Engineering and Analytics: Proceedings of the 9th International Conference on Frontiers in Intelligent Computing: Theory and Applications (FICTA 2021)*, pages 63–70. Springer.

Rico Sennrich and Biao Zhang. 2019. Revisiting low-resource neural machine translation: A case study. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 211–221, Florence, Italy. Association for Computational Linguistics.

Lucia Specia, Stella Frank, Khalil Sima'An, and Desmond Elliott. 2016. A shared task on multimodal machine translation and crosslingual image description. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 543–553.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Ahmet Üstün, Viraat Aryabumi, Zheng Yong, Wei-Yin Ko, Daniel D'souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. 2024. Aya model: An instruction finetuned open-access multilingual language model. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15894–15939, Bangkok, Thailand. Association for Computational Linguistics.

A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*.

BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.

Shuoheng Yang, Yuxin Wang, and Xiaowen Chu. 2020. A survey of deep learning techniques for neural machine translation. *arXiv preprint arXiv:2002.07526*.

# TIM-UNIGE Translation into Low-Resource Languages of Spain for WMT24

**Jonathan Mutal** and **Lucía Ormaechea**
TIM, University of Geneva
40 Boulevard du Pont-d'Arve – Geneva, Switzerland
Jonathan.Mutal@unige.ch, Lucia.OrmaecheaGrijalba@unige.ch

## Abstract

We present the results of our constrained submission to the WMT 2024 shared task, which focuses on translating from Spanish into two low-resource languages of Spain: Aranese (spa-arn) and Aragonese (spa-arg). Our system integrates real and synthetic data generated by large language models (e.g., BLOOMZ) and rule-based Apertium translation systems. Built upon the pre-trained NLLB system, our translation model utilizes a multistage approach, progressively refining the initial model through the sequential use of different datasets, starting with large-scale synthetic or crawled data and advancing to smaller, high-quality parallel corpora. This approach resulted in BLEU scores of 30.1 for Spanish to Aranese and 61.9 for Spanish to Aragonese.

## 1 Introduction

This work presents the results of our constrained submission for the *Translation into Low-Resource Languages of Spain* shared task at WMT24.[1] The task involves translating from Spanish into two low-resource languages spoken in the northeast of the Iberian Peninsula: Aragonese (spa-arg) and Aranese (spa-arn).

Despite the existence of monolingual corpora for these languages, parallel data from Spanish to Aragonese is extremely scarce, amounting to only about 60,000 parallel sentences in OPUS (Tiedemann, 2016). In the case of Aranese, fewer than a thousand parallel sentences are available (FLORES+, Guzmán et al., 2019). In addition to that, these Romance languages are notable for their graphemic instability. Although proposals for orthographic standardization (Estudio de Filología Aragonesa, 2010) and official recognition (Boletín Oficial del Estado, 2006) have been intro-

duced, the absence of a commonly accepted writing system has hindered the development of machine translation (MT) systems into Aragonese and Aranese (Forcada, 2020).

A few previous works have explored MT for these language combinations. For instance, Apertium MT systems (Forcada et al., 2011) provided translations for the above-mentioned pairs using a rule-based approach, achieving better results than neural-based MT systems (Oliver, 2020). Similarly, Cortés et al. (2012) complemented Apertium with an additional orthographic module, and proposed a bidirectional spa-arg MT system. More recently, a multilingual MT model (No Language Left Behind, NLLB Team et al., 2022) included under-resourced Iberian languages like Asturian in its training set. However, it did not cover Aranese or Aragonese.

Given the characteristics of this low-resource scenario, we addressed the translation from Spanish into Aragonese and Aranese using a multilingual multistage approach. The multilingual aspect involved leveraging data from linguistically related languages (such as Occitan for Aranese translation), and employing multilingual pre-trained models (specifically, NLLB[2]) to facilitate generalization across different languages. The multistage approach was designed to consecutively enhance translation performance in the target languages using increasingly specific fine-tuning data sets.

Additionally, we applied data augmentation techniques to increase the volume of relevant data in our training set. This involved: *i*) resorting to LLMs within the constraint of one thousand million parameters (in particular, BLOOMZ[3]) to synthetically create more data in the target languages, and *ii*) producing aligned data through Apertium systems on

---

[2]Particularly, the following model: https://huggingface.co/facebook/nllb-200-distilled-600M.

[3]Specifically: https://huggingface.co/bigscience/bloomz-560m.

the basis of real and synthetic monolingual data from both sides of the languages pairs.

The structure of this paper is as follows: Section 2 describes the methods employed to gather parallel and monolingual data for our experiments. Section 3 introduces the multistage fine-tuning approach. In Section 4, we discuss the experiments conducted on both language combinations and the results obtained. Lastly, Section 5 summarizes our findings and suggests directions for future research.

## 2 Data

To train our MT `spa-arn` and `spa-arg` models, we first compiled parallel data from OPUS and FLO-RES+ (FLORES+$_{DEV}$) bilingual corpora. We then compiled monolingual data from two sources: *i*) we sub-sampled 19 million sentences from Wikimedia and NLLB datasets available in OPUS for Spanish, and *ii*) we collected all monolingual corpora for Aragonese, Aranese and Occitan from OPUS and PILAR (Galiano-Jiménez et al., 2024) when available. Table 1 details the number of segments in the bilingual corpora, and Table 2 reports the segments counts for each monolingual corpus. The notation "k" denotes thousands, and "M" signifies millions. A ✗ indicates the absence of available data.

| *Corpus* | spa-arn | spa-oci | spa-arg |
|---|---|---|---|
| OPUS | ✗ | 1.11M* | 60k |
| FLORES+$_{DEV}$ | 997 | 997** | 997 |

Table 1: Number of parallel segments for the available bilingual dataset. *CCMATRIX was not utilized. **These sentences were not used in any experiment.

| *Corpus* | spa | arn | oci | arg |
|---|---|---|---|---|
| OPUS | 19M | ✗ | 739k | 213k |
| PILAR | ✗ | 322k* | ✗ | 84k |

Table 2: Number of monolingual segments for each available dataset. *Monolingual paragraphs were not utilized.

### 2.1 Synthetic Monolingual Data

We generated synthetic monolingual data in Aranese using BLOOMZ (Muennighoff et al., 2023). To do so, we fine-tuned BLOOMZ with the monolingual data (i.e., PILAR) using a causal language modeling objective, which involves predicting the next token in a sequence. We used

a learning rate of $5 \times 10^{-5}$ with an early stopping mechanism based on accuracy with a patient value of 5. As for the validation data, we randomly picked $1,000$ segments extracted from the same data distribution.

To generate new sentences in the target language, we took the beginnings of sentences in FLORES+$_{DEV}$. Then, the model completed segments from varying numbers of input words (ranging from 1 to 60 words) and generated up to a maximum of 65 tokens. We produced $59,820\,(997 \times 60)$ sentences in Aranese using multinomial sampling. All other generation hyperparameters were set to their default values.[4]

### 2.2 Synthetic Parallel Data

Using the monolingual and synthetic data described above, we produced parallel data through Apertium systems (see Table 3). The following strategies were employed to synthetically create parallel sets:

- **Forward translation** (Burlot and Yvon, 2018). We generated synthetic Aranese, Occitan and Aragonese from monolingual Spanish (see Table 2).

- **Backtranslation** (Sennrich et al., 2016). We backtranslated the segments from monolingual Occitan, Aranese, and Aragonese. We also backtranslated synthetic segments in Aranese produced by BLOOMZ (see Section 2.1).

| *Strategy* | *Corpus* | spa-arn | spa-oci | spa-arg |
|---|---|---|---|---|
| FT | OPUS | 20M | 20M | 20M |
| BT | OPUS | ✗ | 1.8M | 273k |
| | PILAR | 322k | ✗ | ✗ |
| | BLOOMZ | 59k | ✗ | ✗ |
| *Total* | | 20.3M | 21.8M | 20.2M |

Table 3: Training data synthetically generated using forward translation (FT) and backtranslation (BT).

## 3 Approach

Our approach, termed "multistage fine-tuning" involves sequentially refining a model using multiple datasets arranged in a specific order – a method proven to improve performance in machine translation for low-resource language pairs (Dabre et al., 2019).

---

[4]See documentation: `https://huggingface.co/docs/transformers/en/main_classes/text_generation`.

| System | Stage | Data | BLEU↑ | ChrF↑ | TER↓ |
|--------|-------|------|-------|-------|------|
| Apertium | - | - | 28.8 | 49.4 | 72.3 |
| MarianNMT | 1 | OPUS+PILAR (38M) | 25.0 | 47.1 | 76.4 |
| Helsinki-NLP | 1 | OPUS+PILAR (38M) | 22.3 | 45.6 | 81.9 |
| NLLB | 1 | OPUS+PILAR (38M) | 29.0 | 49.4 | 72.3 |
| NLLB | 2.i | PILAR | 28.2 | 48.8 | 73.0 |
| NLLB | 2.ii | PILAR+BLOOMZ | 28.9 | 49.2 | 72.5 |
| NLLB | 3.i | FLORES+$_{DEV}$ | *30.0 | *49.7 | *71.8 |
| NLLB | 3.ii | FLORES+$_{DEV}$ | **\*30.1** | **\*49.8** | **\*71.5** |

Table 4: BLEU, ChrF and TER calculated on the test data for `spa-arn`. Scores with * are significantly better than the baseline Apertium with $p < 0.01$, calculated using paired approximate randomization with $10,000$ trials.

In this work, the models were initially trained using large-scale synthetic or crawled data aiming to match or surpass the performance of the open-source Apertium MT systems. Following this, the models underwent further fine-tuning with smaller, high-quality parallel corpora to improve their performance.

Performance comparisons for the initial models were conducted among three systems: *i*) a model built from scratch using MarianNMT (Junczys-Dowmunt et al., 2018); *ii*) a fine-tuned Helsinki-NLP model with $\approx$72M parameters; and *iii*) a fine-tuned large language model, NLLB, trained on 200 different languages with a larger number of parameters (600M). This enabled us to identify the best performing model for the first stage.

## 4 Experiments and Results

All our systems are Encoder-Decoder models based on the Transformer architecture (Vaswani et al., 2017). The models were trained until convergence, with training progress monitored using BLEU score each $5,000$ steps and an early stopping patience value of 10 using FLORES+$_{DEV}$ as validation data. The details of the training procedure and the results obtained for validation are detailed in Appendix A and B.

In the following sections, we describe the evaluation setup as well as the experiments and results obtained for each language pair.

### 4.1 Evaluation Setup

We evaluated our models using the FLORES+ test data ($1,012$ sentences). We calculated accuracy-based metrics BLEU (Papineni et al., 2002) and ChrF (Popović, 2015), and also computed an error-based metric, i.e., Translation Error Rate (TER, Snover et al., 2006). All metrics were calculated

using the Sacrebleu implementation (Post, 2018).[5] We used paired approximate randomization with $10,000$ trials to calculate the level of significance of the results.

We compared the performance of our models with Apertium MT systems, which are strong baselines for these language pairs.

### 4.2 Spanish-Aranese

For this specific language pair, we had almost no parallel sentences, but we did have a larger corpus of parallel sentences from a linguistically close language, Occitan (see Tables 1 and 2). To leverage the non-negligible quantity of data in this language, we built an MT model using all available data in Occitan and Aranese. In previous experiments, we observed that fine-tuning NLLB with multilingual data (i.e., Spanish-Aranese and Spanish-Occitan) outperformed its bilingual version (i.e., Spanish-Aranese). We also observed that using special tokens to differentiate the two languages is beneficial, and thus used them whenever possible. Appendix A.1 and A.2 show the results of these experiments.

Consequently, in the first stage, the models leveraged all available multilingual data from the OPUS and PILAR (including also synthetic data produced by forward and backtranslation), comprising roughly 42M sentences in Occitan and Aranese. We excluded sentences longer than 100 tokens, resulting in a total of 38M segments. We deliberately omitted synthetic data from BLOOMZ and the validation set to mitigate the risk of overfitting and ensure generalization in the first stage.

In the second stage, the NLLB model, identified

---

[5]The signatures are:
nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp
nrefs:1|case:mixed|eff:yes|nc:6|nw:0|space:no
nrefs:1|case:lc|tok:tercom|norm:no|punct:yes|asian:no.

| System | Stage | Data | BLEU↑ | ChrF↑ | TER↓ |
|---|---|---|---|---|---|
| Apertium | - | - | 61.1 | 79.3 | 27.2 |
| MarianNMT | 1 | OPUS (15M) | 58.2 | 77.8 | 29.9 |
| Helsinki-NLP | 1 | OPUS (15M) | 57.5 | 77.2 | 30.5 |
| NLLB | 1 | OPUS (15M) | *60.5 | *79.0 | *27.7 |
| NLLB-Post-Edition | 2 | FLORES+$_{DEV}$ | 61.0 | *78.9 | 27.2 |
| NLLB-Translation | 2 | FLORES+$_{DEV}$ | *61.9 | 79.5 | *26.8 |

Table 5: BLEU, ChrF and TER calculated on the test data for spa-arg. Scores with * are significantly better or worse than the baseline Apertium with $p < 0.01$, calculated using paired approximate randomization with $10,000$ trials.

as the most performing from the first stage, was fine-tuned using two different data combinations: PILAR data (2.i), and a combination of PILAR and synthetic Aranese data from BLOOMZ (2.ii). To mitigate the risk of overfitting in 2.ii, we fine-tuned the model using a fixed number of steps to reach a slightly higher validation BLEU score than the model trained in 2.i. During the third and final stage, the two models from the previous stage underwent $7,500$ additional training steps on the FLORES+$_{DEV}$ (3.i and 3.ii).

**Results** Results from the first stage showed that NLLB slightly outperformed the Apertium spa-arn system by $0.2$ BLEU points, although this improvement was not statistically significant. MarianNMT and Helsinki-NLP performed worse than Apertium, which appears to agree to the findings in Oliver (2020). Interestingly, MarianNMT outperformed Helsinki-NLP, which might indicate that knowledge acquired during pre-training does not help to the task at hand. The underlying reasons for this discrepancy should be explored in future research.

The most performing model, NLLB (stage 3.ii), which was trained through a three-stage process, surpassed all previous models, improving the Apertium systems by $1.3$ BLEU points and $0.4$ ChrF points, and reduced the TER by $0.8$ points.

The results indicate that the multistage approach enhance model performance. They also underscore the importance of a high-capacity model pre-trained on a diverse set of languages to improve translation from Spanish to Aranese. Additionally, the findings suggest that integrating synthetic data generated by BLOOMZ is beneficial in the third stage of fine-tuning (NLLB 3.i vs. NLLB 3.ii).[6]

### 4.3 Spanish-Aragonese

The model training for spa-arg was conducted in two stages. In the first stage, we used all OPUS-based synthetic data from Spanish to Aragonese to fine-tune NLLB.[7] This initial corpus amounted to roughly 20M parallel sentences, but we later filtered out the source or target sentences exceeding 100 tokens, which resulted in 15M pairs. With this set, we achieved comparable performance to the Apertium MT system in the validation data.

In the second stage, the model was fine-tuned with a lower learning rate, and utilized the FLORES+$_{DEV}$ in two different approaches:

- **Translation**, using as source the original sentences in Spanish.

- **Post-Edition (PE)**, using the Aragonese generated by the Apertium rule-based system as the source to train a post-edition model (apertium_arg-arg).

**Results** The experiments indicate that the performance is superior for translation tasks compared to post-edition tasks. Specifically, our optimal system, NLLB-Translation, surpassed the Apertium baseline by $0.8$ BLEU points and reduced the translation error rate by $0.4$ points.

Regarding the PE model, we assumed that a system trained using apertium_arg-arg could only help correct the mistakes made by such rule-based approach and thus improve its performance. Surprisingly, the resulting model (NLLB-Post-Edition) did not outperform the rule-based system, and instead degraded its results (see Table 5). One possible explanation for this is that the NLLB model from stage 1 was trained on spa-arg translation

---

[6]We also fine-tuned the resulting NLLB model from the first stage with FLORES+$_{DEV}$ data using $7,500$ steps. It underperformed the systems from the third stage.

[7]In previous experiments, we observed that PILAR was not helpful for the spa-arg task, so we decided to exclude it from the training set in our final models.

data rather than post-edition data. Further experiments need to be conducted in order to better understand the behavior of the PE model.

On another note, the results obtained for the fine-tuned Helsinki-NLP model revealed that the knowledge gained during pre-training does not appear to improve the results on the task. As can be observed, the model trained from scratch (MarianNMT) slightly outperforms the small-scale fine-tuned one (Helsinki-NLP), verified by the paired bootstrap statistical test.

## 5 Conclusions

Our experiments demonstrate the potential of combining synthetic data with multilingual pre-trained models to improve translation from Spanish into Iberian low-resource languages like Aranese and Aragonese. By leveraging data from linguistically related languages and employing a multistage approach, the `spa-arn` model achieved a BLEU score of 30.1, while the `spa-arg` model (NLLB-Translation) achieved 61.9 BLEU points. Our findings also indicate that the NLLB model, which benefited from a large number of pre-trained languages and high model capacity, delivered the best performances.

While these results are promising, we have identified several avenues for future research. One key area is to explore the impact of the ratio of real vs. synthetic data for training, as it can help evaluate how changes in data composition influence automatic metrics. Additionally, we plan to investigate the integration of external resources, such as dictionaries (Institut d'Estudis Aranesi, 2019) and orthographic standards (Academia Aragonesa de la Lengua, 2023), to determine whether these can further enhance the performance of our models.

# References

Academia Aragonesa de la Lengua. 2023. *Ortografía de l'aragonés*. Academia Aragonesa de la Lengua.

Boletín Oficial del Estado. 2006. Ley Orgánica 6/2006, de 19 de julio, de reforma del Estatuto de Autonomía de Cataluña. Art. 6.5.

Franck Burlot and François Yvon. 2018. Using monolingual data in neural machine translation: a systematic study. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 144–155, Brussels, Belgium. Association for Computational Linguistics.

Juan Pablo Martínez Cortés, Jim O'Regan, and Francis Tyers. 2012. Free/open source shallow-transfer based machine translation for Spanish and Aragonese. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2153–2157, Istanbul, Turkey. European Language Resources Association (ELRA).

Raj Dabre, Atsushi Fujita, and Chenhui Chu. 2019. Exploiting multilingualism through multistage finetuning for low-resource neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1410–1416, Hong Kong, China. Association for Computational Linguistics.

Estudio de Filología Aragonesa. 2010. *Propuesta ortografica provisional de l'Academia de l'Aragonés*. Edicions Dichitals de l'Academia de l'Aragonés, Zaragoza, Spain.

Mikel Forcada. 2020. Building machine translation systems for minor languages: challenges and effects. In *Revista de Llengua i Dret, Journal of Language and Law*, volume 73, pages 1–20.

Mikel Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis Tyers. 2011. Apertium: A free/open-source platform for rule-based machine translation. *Machine Translation*, 25:127–144.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Pan-iberian language archival resource.

Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc'Aurelio Ranzato. 2019. The FLORES evaluation datasets for low-resource machine translation: Nepali–English and Sinhala–English. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6098–6111, Hong Kong, China. Association for Computational Linguistics.

Institut d'Estudis Aranesi. 2019. *Diccionari der aranés*. Acadèmia Aranesa dera Lengua Occitana.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No Language Left Behind: Scaling Human-Centered Machine Translation.

Antoni Oliver. 2020. Traducción automática para las lenguas románicas de la península ibérica. *Studia Romanica et Anglica Zagrabiensia*, 65:367–375.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Jörg Tiedemann. 2016. OPUS – Parallel Corpora for Everyone. In *Proceedings of the 19th Annual Conference of the European Association for Machine Translation: Projects/Products*. Baltic Journal of Modern Computing.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

## A Training Setup NLLB and Helsinki-NLP

We employed the Adam optimizer with a batch size of 16. We used 50 warm-up steps, and the number of beams was set to 5. The maximum sequence length was set to 100, and the remaining hyper-parameters were left unchanged[8], except for the learning rates which are reported in the following sections. All experiments were conducted using the Transformers library (Wolf et al., 2020) and the University of Geneva HPC clusters, Baobab and Yggdrasil. We used a fixed seed (111) for reproducibility purposes.

### A.1 Helsinki-NLP Results

Given the absence of Aragonese or Aranese as targets in any of the existing OPUS-based Helsinki-NLP MT models, we decided to fine-tune them using different target languages. More specifically, our goal was to determine which of the available Romance languages (namely, Galician, Catalan, French, Italian and Romanian) would be most relevant for the `spa-arg` and `spa-arn` tasks.

After conducting an initial round of experiments, we observed that a geographically close language, Italian (i.e., `Helsinki-NLP/opus-mt-es-it`), most aided the translation into Aragonese on the validation set. Similarly, Catalan (i.e., `Helsinki-NLP/opus-mt-es-ca`) proved to be the most helpful target language for Aranese translation. For this language combination, we also conducted experiments to evaluate the potential gain from the use of two dedicated special tokens for Aranese and Occitan. Specifically, we used `<arn>` for Aranese and `<ca>` for Occitan.

| LR | BLEU | ChrF |
|---|---|---|
| $1 \times 10^{-5}$ | 59.1 | 78.6 |
| $2 \times 10^{-5}$ | 61.9 | 80.1 |
| $3 \times 10^{-5}$ | 62.1 | 80.3 |
| $4 \times 10^{-5}$ | 62.1 | 80.3 |
| $5 \times 10^{-5}$ | **62.2** | **80.3** |

Table 6: Results of Helsinki-NLP `spa-arg` models on validation data with different learning rates.

Once we selected the most relevant model for each language pair, we used different learning rates to fine-tune them for our task at hand. Table 6 reports the BLEU and ChrF results for `spa-arg`

translation. Table 7 shows the results for the two versions of our `spa-arn` models: one that uses a single special token (`<ca>`) and another one that distinguishes between the two languages with distinct special tokens (`<ca>|<arn>`). All experiments were conducted using the Trainer class.[9]

| | LR | BLEU | ChrF |
|---|---|---|---|
| Helsinki-NLP$_{<ca>}$ | $1 \times 10^{-5}$ | 26.0 | 52.7 |
| | $2 \times 10^{-5}$ | 26.5 | 53.2 |
| | $3 \times 10^{-5}$ | 24.8 | 52.2 |
| | $4 \times 10^{-5}$ | 25.8 | 52.8 |
| Helsinki-NLP$_{<ca>|<arn>}$ | $1 \times 10^{-5}$ | **29.7** | **54.9** |
| | $2 \times 10^{-5}$ | 28.6 | 54.3 |
| | $3 \times 10^{-5}$ | 28.8 | 54.2 |
| | $4 \times 10^{-5}$ | 29.0 | 54.9 |

Table 7: Results of Helsinki-NLP `spa-arn` models on validation data with different learning rates and different special token configurations.

### A.2 NLLB Results

To generate Aranese, we used the Occitan special token (`oci_Latn`) in the target, which is presumably the closest language to Aranese covered by NLLB. Similarly to the Helsinki-NLP models, we used the Italian special token (`ita_Latn`) for Aragonese.

| | LR | BLEU | ChrF |
|---|---|---|---|
| NLLB-Bi$_{<oci>}$ | $9 \times 10^{-6}$ | 37.7 | 59.9 |
| | $1 \times 10^{-5}$ | 37.7 | 59.9 |
| | $3 \times 10^{-5}$ | 37.6 | 59.8 |
| NLLB-Multi$_{<oci>}$ | $9 \times 10^{-6}$ | 29.5 | 55.0 |
| | $1 \times 10^{-5}$ | 28.3 | 54.3 |
| | $3 \times 10^{-5}$ | 26.5 | 53.2 |
| NLLB-Multi$_{<oci>|<cat>}$ | $9 \times 10^{-6}$ | 37.8 | 60.0 |
| | $1 \times 10^{-5}$ | **38.1** | **60.1** |
| | $3 \times 10^{-5}$ | 37.9 | 60.0 |

Table 8: Results of NLLB `spa-arn` bilingual (NLLB-Bi$_{<oci>}$) and multilingual models (NLLB-Multi$_{<oci>}$ and NLLB-Multi$_{<oci>|<cat>}$) on validation data with different learning rates and special token configurations.

For Aranese translation, we carried out experiments to evaluate the gain of using a dedicated special token for Aranese and Occitan. In particular, we compared the performance of a multilingual model trained with Aranese and Occitan using the same token (`oci_Latn`), NLLB-Multi$_{<oci>}$, and another model using two special tokens: one for Aranese (`oci_Latn`) and a different one for

---

[8]Refer to: `https://huggingface.co/docs/autotrain/en/seq2seq_params`.

[9]Refer to: `https://huggingface.co/docs/transformers/main_classes/trainer`.

Occitan (`cat_Latn`), NLLB-Multi$_{<oci>|<cat>}$. We also assessed the performance of a bilingual model trained only with Spanish-Aranese data for comparison purposes (NLLB-Bi$_{<oci>}$). Table 8 shows the results on the validation data for the three approaches, indicating that the use of special tokens to differentiate the language is beneficial, and so is including Occitan in the training set.

| Data | PILAR | | PILAR+BLOOMZ | |
|------|-------|-----|--------------|-----|
| LR | BLEU | ChrF | BLEU | ChrF |
| $1 \times 10^{-8}$ | 35.2 | 57.5 | 38.1 | 60.1 |
| $5 \times 10^{-8}$ | 36.0 | 58.4 | 38.0 | 60.0 |
| $1 \times 10^{-6}$ | 37.7 | 59.9 | 39.2 | 60.5 |
| $9 \times 10^{-6}$ | 37.4 | 59.6 | **39.9** | **60.9** |

Table 9: Results of NLLB on stage two with PILAR and BLOOMZ on validation data with different learning rates.

Table 9 shows the results of NLLB on stage two and Table 10 shows the results of NLLB on `spa-arg`.

| LR | BLEU | ChrF |
|----|------|------|
| $5 \times 10^{-7}$ | 64.2 | 81.4 |
| $1 \times 10^{-6}$ | 63.6 | 81.1 |
| $3 \times 10^{-6}$ | 65.2 | 81.9 |
| $9 \times 10^{-6}$ | 65.2 | 81.9 |
| $1 \times 10^{-5}$ | **65.4** | **82.0** |
| $3 \times 10^{-5}$ | 65.3 | 81.9 |

Table 10: Results of NLLB-Baseline `spa-arg` on validation data with different learning rates.

## B MarianNMT Setup and Results

| LR | BLEU | ChrF |
|----|------|------|
| $3 \times 10^{-5}$ | 26.6 | 53.2 |
| $5 \times 10^{-5}$ | 29.6 | 54.9 |
| $3.5 \times 10^{-4}$ | **30.5** | **55.5** |
| $3 \times 10^{-3}$ | n.a.n | n.a.n |

Table 11: Results of MarianNMT `spa-arn` models on validation data with different learning rates.

We used the default hyperparameters from the Marian toolkit (Junczys-Dowmunt et al., 2018) to train the models.[10] We conducted all experiments employing three random seeds and averaging the results measured by the automatic metrics. This

---

[10]Refer to: `https://marian-nmt.github.io/docs/cmd/marian/`.

| LR | BLEU | ChrF |
|----|------|------|
| $3 \times 10^{-5}$ | **55.7** | **78.6** |
| $5 \times 10^{-5}$ | 53.3 | 77.7 |
| $3.5 \times 10^{-4}$ | 50.9 | 76.8 |
| $3 \times 10^{-3}$ | n.a.n | n.a.n |

Table 12: Results of MarianNMT `spa-arg` models on validation data with different learning rates.

approach is intended to reduce the variability of results inherent to individual models randomly initialized.

Tables 11 and 12 present the results for `spa-arn` and `spa-arg` across different learning rates. The notation "n.a.n" indicates that the model diverged at that particular learning rate.

# TAN-IBE participation in the Shared task:
# Translation into Low-Resource Languages of Spain

**Antoni Oliver**

TAN-IBE team - Universitat Oberta de Catalunya (UOC)

`aoliverg@uoc.edu`

## Abstract

This paper describes the systems presented by the TAN-IBE team into the WMT24 Shared task Translation into Low-Resource Languages of Spain. The aim of this joint task was to train systems for Spanish-Asturian, Spanish-Aragonese and Spanish-Aranesian. Our team presented systems for all three language pairs and for two types of submission: for Spanish-Aragonese and Spanish-Aranese we participated with constrained submissions, and for Spanish-Asturian with an open submission.

## 1 Introduction

The TAN-IBE team, consisting of participants from the project *TAN-IBE: Neural Machine Translation for the Romance Languages of the Iberian Peninsula* (Oliver et al., 2023), has developed systems for all language pairs in the task and participated with two types of submissions: Spanish-Asturian (open submission), Spanish-Aragonese (constrained submission), and Spanish-Aranese (constrained submission).

The principal concern of the team was not only the scarcity of resources for these language pairs but also the inadequate quality of the available resources. In order to address these issues, we decided to:

- Clean the existing parallel corpora using a tool developed during the project that rechecks the language of the segments and calculates the cosine similarity between source and target segments.
- Create parallel corpora from Wikipedia for all three language pairs.
- Experiment with the use of backtranslation.
- Experiment with the use of synthetic corpora.
- Experiment with the use of multilingual systems.

For some of the language pairs and direction of the shared task, a freely available rule-based MT system exists: Apertium[1] (Forcada et al., 2011;

Khanna et al., 2021). In previous research (Oliver, 2020), it was demonstrated that Apertium achieves highly competitive quality, and that it is challenging for a neural system to achieve superior quality results. Consequently, this system has been employed to create backtranslated and synthetic corpora. Furthermore, Apertium will be employed as a reference system to facilitate the evaluation of the trained systems prior to submission to the shared task. The specific versions of Apertium used for each language pair are: Spanish–Aragonese 0.6.0,[2] Spanish–Aranese 1.0.8,[3] Spanish–Asturian 1.1.1.[4]

## 2 Tools

To train the NMT systems we have used the marian-nmt[5] (Junczys-Dowmunt et al., 2018). All the systems have been trained with a Transformer-big configuration. To calculate the subwords units, we have used SentencePiece[6] (Kudo and Richardson, 2018). Additional information on other training parameters can be found in the subsections for each system.

In order to create the parallel corpora and to clean and preprocess the corpora, several components of the MTUOC project[7] (Oliver and Alvarez, 2023) have been employed. It should be noted that several of these components have been developed during the course of the TAN-IBE project. The components that have been used are as follows:

- To create parallel corpora from the Wikipedia: MTUOC-WikipediaDump[8] and MTUOC-

---

[1] `https://apertium.org/`

[2] `https://github.com/apertium/apertium-spa-arg/releases/tag/v0.6.0`

[3] `https://github.com/apertium/apertium-oc-es/releases/tag/v1.0.8`

[4] `https://github.com/apertium/apertium-spa-ast/releases/tag/v1.1.1`

[5] `https://marian-nmt.github.io/`

[6] `https://github.com/google/sentencepiece`

[7] `https://mtuoc.github.io/`

[8] `https://github.com/mtuoc/MTUOC-WikipediaDumps`

aligner.[9]

- To clean both the existing and the newly created parallel corpora: MTUOC-clean-parallel-corpus,[10] that performs several cleaning operations; and MTUOC-PCorpus-rescorer,[11] that rechecks the language of the segments and calculates the cosine similarity using a multilingual model, as SBERT, for example.
- To select from a large corpus the most similar segments from a large corpus: MTUOC-corpus-combination.[12]
- To preprocess the parallel corpora to train the systems: MTUOC-corpus-preprocessing.[13]

For some cleaning operations the language of the segments should be detected. As Asturian, Aranese and Aragonses are underrepresented in available language detection models, we decided to develop our own language detection model using fasttext[14] (Joulin et al., 2016). We trained a model able to detect the following languages: Aragonese, Aranese, Asturian, Catalan, English, French, Galician, Occitan, Portuguese and Spanish. For all the languages, except for Aranese, we used 400K segments from the Wikipedia. As no Wikipedia is available for Aranese, we used 297,557 segments from the PI-LAR corpus.[15]

The trained model performs similarly to Idiomata Cognitor[16] (Galiano-Jiménez et al., 2024), with the difference than Italian is included in this tool. In Table 1, the precisions calculated using the FLORES+ dev corpus for Idiomata Cognitor and our trained language detection model are presented.

We decided to train our own language detection model because fasttext models integrate seamlessly in our corpus cleaning scripts, and the same training strategy can be used in future experiments with language not present in Idomata Cognitor.

## 3 Existing resources

Table 2 provides an overview of the existing corpora that have been employed for system training,

---

| Language | Idiomata Cognitor | TAN-IBE fasttext model |
|---|---|---|
| Spanish | 0.95 | 0.95 |
| Catalan | 1.00 | 0.99 |
| Aragonese | 0.96 | 1.00 |
| Aranese | 0.96 | 1.00 |
| Occitan | 0.94 | 0.93 |
| Asturian | 0.99 | 0.98 |
| Galician | 0.98 | 0.99 |
| French | 1.00 | 1.00 |
| Portuguese | 1.00 | 0.99 |

Table 1: Precision on language detection for Idiomata Cognitor and our trained fasttext model.

accompanied by the number of segments in the original corpora and the number of segments resulting from the cleaning process peformed with MTUOC-PCorpus-rescorer. As previously stated, this tool performs a second language detection of the segments using fasttext and calculates a cosine similarity between the source and target segments using SBERT. As the default language detection model used in fasttext have been trained with underrepresented texts for Asturian, Aragonese and Occitan (Aranese), we decided to retrain a language model for the cleaning of corpora for these languages, as explained in section 2. Please, note that there are no available parallel corpora for Aranese, and in the table we state the figures for the Spanish-Occitan parallel corpus used.

| Langs | Corpus | Raw | Clean |
|---|---|---|---|
| spa-ast | NLLB | 6,470,015 | 504,532 |
| spa-arg | Wikimatrix | 33,724 | 16,456 |
| spa-oci | NLLB | 925,448 | 108,440 |

Table 2: Size of the existing corpora in segments used for training the systems

## 4 Newly created resources

As the available corpora for the working language pairs are clearly insufficient to train NMT systems we have created a new parallel corpus from Wikipedia dumps. To this end, we have developed a series of scripts, which are freely available at the MUTOC-WikipediaDumps repository, that are capable of:

- Extract all the text from the Wikipedia dump, along with a file containing the titles of the articles. This process is performed for the smaller Wikipedias, in this case the Asturian, Aragonese and Aranese.
- Translate the list of titles of the extracted articles into Spanish using the langlinks database dump.

- Extract the text of the articles of the larger Wikipedia, in our case the Spanish one, restricting the extracted articles to the titles in the translated title list.

It should be noted that this process is carried out separately for Spanish-Asturian, Spanish-Aragonese and Spanish-Occitan. Once the text has been obtained, it is segmented and the resulting segments are deduplicated and a file is generated for each language pair, containing all the source segments and a file for all the target segments. For the creation of the Wikipedia corpora we have used the dumps of first of May of 2024. All the created corpora can be downloaded from Github.[17]

To align the files we use a bitext mining strategy using SBERT, implemented in the MTUOC-Aligner. The alignment process gives sets of aligned files that contain source segment, target segment and a margin score. The resulting aligned files are cleaned using MTUOC-PCorpus-rescorer, using a confidence of 0.75 for language detection and SBERT score. In table 3 illustrates the raw and cleaned number of parallel segments obtained through this process.

| Langs | Raw | Clean | Backtrans. |
|---|---|---|---|
| spa-ast | 2,194,031 | 1,193,264 | 4,023,140 |
| spa-arg | 153,863 | 32,374 | 386,666 |
| spa-oci | 295,476 | 55,112 | 1,151,205 |

Table 3: Size of the Wikipedia corpora created

Additionally, the extracted text from Wikipedia in Asturian, Aragonese and Occitan has been employed to generate back-translated Spanish-Asturian, Spanish-Aragonese and Spanish-Occitan corpora. In order to achieve this, it is necessary to have access to machine translation systems that are capable of functioning in the opposite direction. Fortunately, Apertium provides translation systems for Aragonese-Spanish and Occitan-Spanish, but not for Asturian-Spanish. To address this gap, we have trained a Transformer Neural system using the cleaned NLLB and the Wikipedia corpus created for this purpose. This system was then used for backtranslation. In subsection 5.1 more details on this system, along with evaluation figures are presented. Table 3 shows the number of segments in the backtranslated Wikipedia corpus.

As previously stated, there are no available corpora for the Spanish-Aranese language pair, and there is no Wikipedia version for Aranese either.

Therefore, no parallel resources for this language pair can be employed. For this language pair, we have utilised the monolingual data available in the PILAR (*Pan-Iberian Language Archival Resource*) and backtranslated it into Spanish using Apertium. This process yielded 297,557 backtranslated parallel segments.

In order to augment the number of parallel corpora, we devised a method for the generation of synthetic data utilising the Spanish-Catalan Paracrawl corpus, which had been previously subjected to cleaning procedures. The 13 million Spanish segments were translated into Asturian, Aragonese and Aranese using Apertium. For each translated segment, a confidence score was calculated as the ratio of the unknown words (marked with an asterisk by Apertium) to the total number of words.

In the case of the multilingual system under experimentation, the parallel corpora presented in Table 4 were also employed. With regard to the Spanish-Catalan and Spanish-Galician languages, a parallel corpus was created from Wikipedia, as previously described.

| Langs | Corpus | Raw | Clean |
|---|---|---|---|
| spa-cat | Paracrawl | 17,238,953 | 13,931,594 |
| spa-cat | Wikipedia | 5,790,903 | 2,586,448 |
| spa-fra | MultiParacrawl | 39,026,138 | - |
| spa-fra | WikiMatrix | 905,761 | - |
| spa-glg | Paracrawl | 1,879,649 | - |
| spa-glg | Wikipedia | 1,697,307 | 729,840 |
| spa-por | MultiParacrawl | 26,181,054 | - |
| spa-por | WikiMatrix | 923,725 | - |

Table 4: Additional parallel corpora used to train the multilingual system

## 5 Trained systems and evaluation

In this section we will present all the trained systems for this shared task, along with evaluation figures using the FLORES+ devtest sets developed by the organisers of the shared task. To calculate the evaluation metrics (BLEU, TER and chrF will be presented), we have used Sacrebleu[18] (Post, 2018). In Appendix 8 we present the metric signatures for these metrics.

### 5.1 Neural Asturian-Spanish for backtranslation

In order to create back-translated corpora for Aragonese and Aranese, the Apertium system was

---

[17] https://github.com/mtuoc/WikipediaCorpora

[18] https://github.com/mjpost/sacrebleu

employed, given that it is available for these language pairs. However, the Asturian-Spanish pair is not available in Apertium. Consequently, a neural system was trained using Marian and the following parallel corpora: NLLB cleaned (see table 2) and the newly created Wikipedia corpus (see table 3). A transformer configuration using SentencePiece with a vocabulary size of 32K has been used.

Table 5 presents the evaluation metrics for the trained systems, along with the metrics for the Apertium system for Aragonese-Spanish and Aranese-Spanish, using the FLORES+ devtest sets.[19] For purposes of comparison, the evaluation figures for all reverse translation directions are also provided. It should be noted that Apertium systems are available for all reverse directions.

| System | BLEU | chrF | TER |
|---|---|---|---|
| Marian ast-spa | 24.0 | 53.5 | 57.4 |
| Apertium arg-spa | 61.4 | 77.9 | 26.7 |
| Apertium oci-aran–spa | 26.7 | 47.2 | 70.8 |
| Apertium spa-ast | 17.0 | 50.8 | 80.4 |
| Apertium spa-arg | 61.1 | 79.3 | 27.2 |
| Apertium spa–oci-aran | 28.8 | 49.4 | 72.3 |

Table 5: Evaluation figures for the systems used to create backtranslated corpora.

A number of conclusions can be drawn from Table 5. The Apertium systems for Aragonese-Spanish and Spanish-Aragonese achieve highly comparable results, as do the pairs Aranese-Spanish and Spanish-Aranese. Consequently, given that the Marian system for Asturian-Spanish attains superior outcomes compared to the reverse system Spanish-Asturian, it can be inferred that the quality of the training system is analogous to, and even surpasses, that of the Apertium systems utilised for backtranslation.

The Asturian-Spanish NMT system was employed to backtranslate the segments from the NLLB corpus that were identified as Asturian, but with translations in languages other than Spanish or with Spanish translations with SBERT scores below 0.75. The resulting backtranslated Spanish-Asturian corpus comprises 2,084,594 segments.

### 5.2 Basic neural Spanish-Asturian system

The basic neural system for Spanish-Asturian has been trained using the same configuration than the Asturian-Spanish for backtranslation, that is: NLLB cleaned (see table 2) and the newly cre-

ated Wikipedia corpus (see table 3). A transformer configuration using SentencePiece with a vocabulary size of 32K has been used. In table 6 we can observe the evaluation figures for this system. It can be observed that this basic system attains inferior results in comparison to Apertium (see Table 5). The BLEU score is 15.3, whereas Apertium achieves 17.0.

### 5.3 Basic neural Spanish-Aragonese system

The basic neural system for Spanish-Aragonese has been trained using the same configuration than the Asturian-Spanish, using the following corpora: Wikimatrix cleaned (see table 2), the newly created Wikipedia corpus (see table 3, and the newly created backtranslated Wikipedia corpus (see table 3).

In table 6 we can observe the evaluation figures for this system. The evaluation results obtained by this system are significantly lower (18.8 BLEU) than the obtained by Apertium (see Table 5) (61.1 BLEU).

### 5.4 Multilingual system

As no parallel corpora are available for Aranese, and the available for Aragonese are very small, we experimented with multilingual systems to see whether the multilingual configuration may produce good results for these two languages, and may also improve the results obtained for Asturian.

We have trained a multilingual system from Spanish to the following languages: Aragonese, Aranese, Asturian, Catalan, Galician, French, Occitan and Portuguese. To train this system we have used the following corpora:

- **Spanish-Aragonese**: WikiMatrix cleaned (see table 2), newly created Wikipedia corpus (see table 3), Wikipedia backtranslated using Apertium (see table 3).
- **Spanish-Aranese**: Pilar backtranslated using Apertium.
- **Spanish-Asturian**: NLLB cleaned (see table 2), Wikipedia clean (see table 3) and the backtranslated corpus described in 5.1
- **Spanish-Catalan**: 10 M segments automatically selected form Paracrawl cleaned (see table 4)
- **Spanish-Galician**: all segments available in Paracrawl cleaned (see table 4)
- **Spanish-French**: 10 M segments automatically selected form MultiParacrawl cleaned (see table 4)

---

[19]Apertium v 3.9.4; linguistic data versions: spa-ast v.1.1.1-1; spa.arg v.0.6.0-1; es-oc v1.0.8-1

- **Spanish-Occitan**: NLLB cleaned (see table 2), newly created Wikipedia corpus (see table 3), Wikipedia backtranslated using Apertium (see table 3).
- **Spanish-Portuguese**: 10 M segments automatically selected form MultiParacrawl cleaned (see table 4)

To calculate the SentencePiece model for the multilingual system we have randomly selected 1M segments from each language pair, except for Aragonese and Aranese, for which we have used all the available segments. The final multilingual system uses an ensemble of the three best check points regarding BLEU DETOK.

Table 6 presents the evaluation figures for the multilingual system. With regard to Spanish-Asturian, the multilingual configuration exhibits an increase of 0.9 BLEU points in comparison to the basic system. However, this figure remains 0.8 BLEU points below the level achieved by Apertium. In the case of Spanish-Aragonese, the multilingual system exhibits a noteworthy enhancement of 13.3 BLEU points. Nevertheless, it remains considerably distant (29 BLEU points) from the performance of Apertium. In the case of Spanish-Aranese, the multilingual configuration yielded a BLEU score that was 8.7 points lower than that obtained with Apertium.

| System | langs. | BLEU | chrF | TER |
|---|---|---|---|---|
| Basic | spa-ast | 15.3 | 48.0 | 77.5 |
| Basic | spa-arg | 18.8 | 51.5 | 67.2 |
| Multilingual | spa-ast | 16.2 | 50.0 | 75.8 |
| Multilingual | spa-arg | 32.1 | 65.3 | 48.0 |
| Multilingual | spa–oci-aran | 20.1 | 44.5 | 77.8 |
| Synth. val Flores | spa-ast | 16.3 | 50.6 | 77.1 |
| Synth. val Flores | spa-arg | 57.2 | 78.1 | 29.4 |
| Synth. val Flores | spa–oci-aran | 26.9 | 48.8 | 72.7 |
| Backt. val Flores | spa-ast | 18.0 | 51.6 | 74.5 |
| Apertium | spa-ast | 17.0 | 50.8 | 80.4 |
| Apertium | spa-arg | 61.1 | 79.3 | 27.2 |
| Apertium | spa–oci-aran | 28.8 | 49.4 | 72.3 |

Table 6: Evaluation figures for the different systems trained and Apertium

### 5.5 Synthetic val Flores Spanish-Asturian

This systems uses the same configuration as the Basic Spanish-Aranese system, but adding the synthetic corpus from Paracrawl described in subsection 4. In this system we use the Flores dev corpus for validation. The final system uses an ensemble of the model corresponding the the 3 best checkpoints using the BLEU DETOK metric.

From table 6 we can observe that with this configuration we improve the basic system by 1 BLEU point, but we are still below Apertium.

### 5.6 Synthetic val Flores Spanish-Aragonese

This systems uses the same configuration as the Basic Spanish-Aragonese system, but adding the synthetic corpus from Paracrawl described in subsection 4. In this system we use the Flores dev corpus for validation. The final system uses an ensemble of the model corresponding the the 3 best checkpoints using the BLEU DETOK metric.

This configuration achieves an impressive improvement of 38.4 BLEU points in comparison with the basic configuration (see table 6), but Apertium keeps an advantage of 3.9 BLEU points. At this point, we can try to explain two key observations: why does Apertium achieve such strong evaluation metrics for Spanish-Aragonese? And why does the system trained with synthetic corpora created by Apertium show such a remarkable improvement? The answer may lie in how the Flores corpus was developed for this language pair. Since it was generated through machine translation from the Spanish Flores using Apertium, the system has a clear advantage when calculating automatic evaluation metrics.

This is the final submission for the Open systems for Spanish-Asturian, with id 568.

### 5.7 Synthetic val Flores Spanish-Aranese

This system have been trained using the backtranslated Pilar corpus and the synthetic corpus from Paracrawl described in subsection 4. For the validation set we have used the Flores dev corpus. The final system uses an ensemble of the model corresponding the the 3 best checkpoints using the BLEU DETOK metric.

This system achieves an improvement of 6.8 BLEU points in comparison with the multilingual configuration (see table 6), but it is still 1.9 BLEU points below Apertium.

This is the final submission for the Open systems for Spanish-Aranese, with id 610.

### 5.8 Backtranslation val Flores Spanish-Asturian

We have followed the same configuration than the Basic Spanish-Asturian, but using also the backtranslated Wikipedia corpus. For the validation set we have used the Flores dev corpus. The final system uses an ensemble of the models corresponding

the the 3 best checkpoints using the BLEU DETOK metric.

This system achieves an improvement of 2.7 BLEU points in comparison with the basic system, and 1.7 BLEU points in comparison with the synthetic one (see table 6), and also outperforms Apertium by 1 BLEU point.

This is the final submission for the Open systems for Spanish-Asturian, with id 568.

# 6 Final submissions

Table 7 contains a comprehensive overview of the systems that have been submitted to the shared task. With regard to the Spanish-Asturian task, an open system has been submitted, whereas for Spanish-Aragonese and Spanish-Aranese, constrained systems have been submitted.

# 7 Energy consumption report

The training scripts generate a log with the timestamp and the GPU consumption every second throughout the entirety of the training process. This enables the calculation of the total training time and the approximate energy consumption in kWh. The total training time and the consumption of each of the two GPU units utilized, along with the total consumption, are presented in Table 8. The final submitted systems are highlighted in bold.

As we can see, all the systems are trained in short times (from 1 h. 10 m. to 3 h. 23 m.) in a relatively modest computer wit two NVIDIA RTX A5000 GPU units with 24 GB each, a AMD Ryzen Threadripper PRO 3945WX CPU with 12-Cores and 64 GB of RAM. As the training times are short, the energy consumption is very low in all trainings, ranging from 0.481 to 1.292 kWh.

# 8 Conclusions and future work

In this paper we have presented the systems that the TAN-IBE team have submitted to the WMT24 Shared Task Translation into Low-Resource Languages of Spain. We have presented an open system for the Spanish-Asturian language pair, and constrained systems for Spanish-Aragonese and Spanish-Aranese.

The primary challenge in completing the task was the unavailability of high-quality parallel corpora for the specified language pairs. Fortunately, all the language pairs in question have an Apertium system, and Apertium is also available for all the reverse language pairs except Asturian-Spanish.

This enabled us to conduct experiments with synthetic and backtranslated corpora. To perform backtranslation experiments for Spanish-Asturian, we trained a basic neural Asturian-Spanish system.

Additionally, monolingual and parallel corpora were generated from Wikipedia dumps for Spanish to Asturian, Aragonese and Occitan (given the unavailability of an Aranese Wikipedia).

Regarding the training strategies, we experimented with bilingual and multilingual systems. While multilingual systems demonstrated enhanced performance relative to basic systems, the use of synthetic and backtranslated corpora yielded superior outcomes.

In the period preceding the conclusion of the TAN-IBE project in July 2025, it is our intention to undertake the following actions:

- During the course of the project, a corpus of monolingual and bilingual texts in Asturian, Aragonese, Aranese and Spanish has been compiled. The next stage is to process and align these texts in order to increase the number of available parallel segments.
- We plan to train new systems using the parallel corpora and the training techniques presented in this paper. The quality of the resulting systems will then be evaluated in order to ascertain whether the inclusion of the new parallel texts has had a positive impact.
- Furthermore, we intend to learn from the other participants in the shared task and attempt to reproduce the training techniques that have yielded the most favourable outcomes, utilising the newly created parallel corpora.
- We plan to develop neural systems for Spanish to the other languages of the TAN-IBE project, namely, Portuguese, Galician, Asturian, Argonese, Catalan and Aranese. These systems will be freely released.
- We also plan to train a multilingual system able to translate to and from all the languages of the TAN-IBE project.

After the completion of the TAN-IBE project, we plan to increase the size of the Apertium monolingual and transfer dictionaries for Spanish to Asturian, Aragonese and Aranese using automatic techniques that make use of monolingual and parallel corpora. The quality of the Apertium systems is noteworthy, and the enhancement of the dictionaries has the potential to further optimise the efficiency of the Apertium systems for the generation

| Submission | Type | Section | langs. | BLEU | chrF | TER |
|---|---|---|---|---|---|---|
| #568 | Open | 5.8 | spa-ast | 18.0 | 51.6 | 74.5 |
| #584 | Constrained | 5.6 | spa-arg | 57.2 | 78.1 | 29.4 |
| #610 | Constrained | 5.7 | spa–oci-aran | 26.9 | 48.8 | 72.7 |

Table 7: Information about the systems submitted to the shared task.

| System | langs. | Section | Time. | GPU0 (kWh) | GPU1 (kWh) | Total (kWh) |
|---|---|---|---|---|---|---|
| Backtranslation | ast-spa | 5.1 | 2 h. 51 m. | 0.594 | 0.598 | 1.191 |
| Basic | spa-ast | 5.2 | 3 h. 23 m. | 0.704 | 0.709 | 1.414 |
| Basic | spa-arg | 5.3 | 1 h. 38 m. | 0.332 | 0.339 | 0.671 |
| Multilingual | spa-MULT | 5.4 | 3 h. 7 m. | 0.627 | 0.643 | 1.270 |
| Synthetic | spa-ast | 5.5 | 3 h. 13 m. | 0.645 | 0.650 | 1.292 |
| **Synthetic** | **spa-arg** | **5.6** | **1 h. 31 m.** | **0.302** | **0.305** | **0.606** |
| **Synthetic** | **spa–oci-aran** | **5.7** | **2 h. 53 m.** | **0.577** | **0.583** | **1.160** |
| **Backtranslation** | **spa-ast** | **5.8** | **1 h. 10 m.** | **0.240** | **0.241** | **0.481** |

Table 8: Total time and energy consumption for all the trainings.

of synthetic and backtranslated corpora.

## Acknowledgements

## References

Mikel L Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine translation*, 25:127–144.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Idiomata cognitor.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Tanmai Khanna, Jonathan N Washington, Francis M Tyers, Sevilay Bayatlı, Daniel G Swanson, Tommi A Pirinen, Irene Tang, and Hector Alos i Font. 2021. Recent advances in apertium, a free/open-source rule-based machine translation platform for low-resource languages. *Machine Translation*, 35(4):475–502.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Antoni Oliver. 2020. Traducción automática para las lenguas románicas de la península ibérica. *Studia Romanica et Anglica Zagrabiensia: Revue publiée par les Sections romane, italienne et anglaise de la Faculté des Lettres de l'Université de Zagreb*, 65:367–375.

Antoni Oliver and Sergi Alvarez. 2023. Training and integration of neural machine translation with MTUOC. In *Proceedings of the 1st Workshop on Open Community-Driven Machine Translation*, pages 5–13, Tampere, Finland. European Association for Machine Translation.

Antoni Oliver, Mercè Vàzquez, Marta Coll-Florit, Sergi Álvarez, Víctor Suárez, Claudi Aventín-Boya, Cristina Valdés, Mar Font, and Alejandro Pardos. 2023. TAN-IBE: Neural machine translation for the romance languages of the iberian peninsula. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 495–496.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

## Appendix - Metric signatures

- BLEU nrefs:1 | bs:1000 | seed:12345 | case:mixed | eff:no | tok:13a | smooth:exp | version:2.3.1
- chrF2 nrefs:1 | bs:1000 | seed:12345 | case:mixed | eff:yes | nc:6 | nw:0 | space:no | version:2.3.1
- TER nrefs:1 | bs:1000 | seed:12345 | case:lc | tok:tercom | norm:no | punct:yes | asian:no | version:2.3.1

# Enhaced Apertium System: Translation into Low-Resource Languages of Spain
## Spanish–Asturian

**Sofía García González**
**imaxin**|software
Rúa dos Salgueiriños de Abaixo, 11, L-6, 15703
Santiago de Compostela, A Coruña
sofia.garcia@imaxin.com

## Abstract

We present the Spanish–Asturian Apertium translation system, which has been enhanced and refined by our team of linguists for the shared task: Low Resource Languages of Spain of this WMT24 under the closed submission. While our system did not rank among the top 10 in terms of results, we believe that Apertium's translations are of a commendable standard and demonstrate competitiveness with respect to the other systems.

## 1 Introduction

In this shared task: Translation into Low-Resource Languages of Spain, we present an enhancement of the machine translator *Eslema* system for Spanish–Asturian pair (Viejo et al., 2008). We present our system under the closed submission.

Asturian is not an officially recognised language in the Spanish state. The 1981 Statute of Autonomy of Asturias makes only passing reference to Asturian, citing the need to protect and disseminate it. However, it does not accord the language the same privileges as are enjoyed by other official languages of Spain, such as Galician, Catalan or Basque. Subsequently, on 23 March 1998, the Asturian Parliament enacted the Law on the Use and Promotion of Asturian. The aforementioned legislation stipulates that all citizens are entitled to utilise Asturian in verbal and written communication, and that such communication shall be deemed valid. Furthermore, it acknowledges the necessity for the dissemination of Asturian in educational and media contexts (Galán y González, 2015). Consequently, Asturian is categorised as a minority and Low-Resource Language (LRL), exhibiting a paucity of resources and a diminished presence in Natural Language Processing (NLP) relative to other co-official languages of Spain.

The objective of this shared task is to develop innovative systems and data resources for this low-resource language, the Aragonese and the Aranese. In light of the aforementioned considerations, we present an enhacement of Apertium (Forcada et al., 2011), the foundational system of the current Spanish–Asturian MT translator, *Eslema*. Based on this open-source system, a series of grammatical, syntactic and lexical improvements have been implemented in order to participate in this shared task. While the results obtained have not been sufficient to maintain a position within the top 10, they have been noteworthy.

## 2 Eslema

*Eslema* was initiated as a project of the University of Oviedo in 2004 with the objective of assembling corpora in Asturian language. The Asturian Philology Group (*Seminariu de Filoloxía Asturiana*) within the Spanish Philology Department was responsible for its research, with the aim of compiling texts of diverse typology, format and historical periods (Viejo et al., 2008).

Subsequently, at the conclusion of 2008, the Principality of Asturies (*Conseyería de Cultura del Principáu d'Asturies*) assumed the economic responsibility for the establishment of the regulatory framework for the development of a rule-based machine translator for the Spanish–Asturian language pair. This project was carried out in collaboration between the University of Oviedo and the Apertium community. The report published in early 2010 about this translator acknowledged that its functionality was satisfactory, especially in the Spanish–Asturian direction. However, it was also noted that the software still presented some residual problems that would be solved in subsequent updates (Universidad de Oviedo, 2010).

Nevertheless, it is important to recognise that, despite the best efforts of the developers to rectify all potential errors, no machine translator can be considered perfect. In particular, Rule Based Machine Translators require ongoing maintenance and revision to ensure optimal performance. Subsequently,

878

*Eslema* has been provided to the Administration and citizens free of charge until nowadays. Right now there is a new version of the machine translator in the Linguistic Policy (*Política Llingüística*) webpage.[1]

## 3 Apertium

Apertium is a free/open-source platform for Rule-Based Machine Translation (RBMT) (Forcada et al., 2011). It is designed to provide high-quality translation tools for several LRL pairs. Apertium was initially developed as part of a project developed by the Alacant University and different public and private Spanish institutions and companies such as **imaxin**|software[2] and El Huyar.[3] But since then, it has evolved into a collaborative endeavour involving developers, linguists and researches (Khanna et al., 2021). The platform is constructed on a rule-based translation system, which relies on predefined linguistic rules and a modular architecture. In the most recent versions of Apertium, these modules are divided into two monolingual packages and one bilingual package for each language pair. The following subsections will provide an explanation of the monolingual packages 3.1 and the bilingual package 3.2.[4] These modules constitute the various components of the Apertium engine pipeline. Modifications have been made to them with the objective of enhancing the original *Eslema* translator.

### 3.1 Monolingual Packages

There is a monolingual package for each language in the language-pair. For example, in this case, an Asturian package and a Spanish package. Each of these packages is formed by a dictionary 3.1.1, a post-generator 3.1.2 and a Constraint Grammar 3.1.3.

### 3.1.1 Monolingual Dictionary

Monolingual dictionaries[5] serve the function of regulatory modules for the system's lexicon. The dictionary is comprised of two principal sections.

The initial one is the paradigm section, wherein paradigms are defined as patterns or models that delineate the potential declensions of each term, contingent on its category or morphology. The subsequent section is where the lexicon is incorporated. Each novel word is introduced in the format of an entry, which encompasses the term in its fundamental form and the paradigm ascribed to it.

### 3.1.2 Post-generator

The post-generator[6] is a module that is employed to rectify minor spelling issues in each language. To illustrate, in languages where contractions or the use of apostrophes is prevalent, these orthographic phenomena are regulated by the post-generator. It is typically a module that remains consistent for each language and does not necessitate significant updates.

### 3.1.3 Constraint Grammar

As posited by Bick and Didriksen (2015), the Constraint Grammar (CG) may be conceived of as a declarative whole of contextual possibilities and impossibilities for a language or genre. However, in programming terms, it is implemented procedurally as a set of consecutively iterated rules that add, remove or select tagged-encoded information.

In Apertium, each language package is equipped with a CG tool, which serves to clarify the source text. One illustrative example of a CG rule is as follows[7]:

The rule "SELECT VERB IF (1 (det))" indicates that the verb category must be selected whenever the following word is a determiner.

This tool is of vital importance in an RBMT engine, as disambiguation errors can lead to significant translation errors. Therefore, the more developed the CG is, the more accurate the translation will be. For this particular pair, the Spanish CG[8] has been used, which was previously created by the Apertium community and, due to errors detected, has had to be modified on some occasions.

### 3.2 Bilingual Package

There is one bilingual package for each language pair. Each bilingual package is formed by a bilin-

---

[1] https://politicallinguistica.asturias.es/eslema

[2] https://imaxin.com/gl/

[3] https://www.elhuyar.eus/eu

[4] The majority of the information pertaining to these modules has been derived from the Apertium Wiki: https://wiki.apertium.org/wiki/Main_Page

[5] See: https://wiki.apertium.org/wiki/Monodix_basicshttps://wiki.apertium.org/wiki/Monodix_basics

[6] See: https://wiki.apertium.org/wiki/Post-generator

[7] See: https://wiki.apertium.org/wiki/Constraint_Grammar

[8] See: https://github.com/apertium/apertium-spa/blob/master/apertium-spa.spa.rlx

gual dictionary 3.2.1, the transfer rules 3.2.2 and the lexical selection rules 3.2.3.

### 3.2.1 Dictionary

In bilingual dictionaries,[9] the terms of both languages are aligned with one another. As a general rule, in this type of engine, each term in the source language can only have one correspondence in the target language. In other words, a term in the source language will be translated by the same term in all contexts, with the exception of specific cases which will be addressed in the subsequent modules.[10]

### 3.2.2 Transfer Rules

Transfer rules[11] are employed to oversee the most intricate structural divergences between two languages, whether pertaining to syntax, morphology, or grammar. To illustrate, transfer rules facilitate the rearrangement of a sentence in the target language, the alteration or insertion of tags by category, and other modifications that enhance the coherence of the target language. In essence, this module is responsible for managing the majority of complex changes that are contingent upon grammatical or lexical context.

### 3.2.3 Lexical Selection Rules

In instances where a term in the source language has two or more potential translations in the target language, the lexical transfer rules module[12] is responsible for selecting one or the other option, depending on the surrounding context. This context may be either grammatical or lexical in nature.

## 4 Dependencies

The dependencies of our translation system are presented in the following list. It is imperative that all modules are in place for the correct functioning of the pair: Apertium-3.8.3, lttoolbox-3.7.1, apertium-lex-tools-0.4.2 and cg3-3.9.

---

[9]See:https://wiki.apertium.org/wiki/Bilingual_dictionary

[10]The latest versions of Apertium include the Lexical Selection Rules module 3.2.3, which enables to assign a specific meaning and translation to the target language depending on the context. This module will be explained in subsection.

[11]See:https://wiki.apertium.org/wiki/A_long_introduction_to_transfer_rules

[12]See:https://wiki.apertium.org/wiki/Constraint-based_lexical_selection_module

## 5 Methodology

In order to participate in this shared task, the team at **imaxin|software** has utilized the open-source translator published in 2010 by *Eslema*[13] and made available on the Apertium project website,[14] to enhance it in the morphological 5.1 and lexical 5.2 linguistic areas. The implementation of these alterations and enhancements was overseen at all times by a team of linguists with expertise in Asturian, over a period of 18 months, during which not only the Spanish–Asturian direction was considered, but also the Asturian–Spanish.

### 5.1 Morphological Enhacement

With regard to morphological errors, three principal categories may be identified. Firstly, this pair presented a multitude of disambiguation issues. To illustrate, the preposition *para* (for) in Spanish was frequently analyzed as the third person singular of the verb *parar* (to stop). This resulted in errors such as: *Ir para casa* (go home) in Spanish was translated to Asturian as *ir para casa* innstead of *Ir pa casa*. These types of errors were corrected in a generic way by making use of the Constraint Grammar that had already been created by the Apertium community. However, it was also necessary to create new rules for specific cases.

Furthermore, the paradigms created for different grammatical categories contained various errors, either because the term had been assigned the wrong paradigm or because the assigned paradigm contained errors in its definition. To illustrate, Asturian verbs ending in *-ñir* or *-xir* (e.g. *teñir* (to dye ) and *dirixir* (to address)) exhibited erroneous conjugation of the third person singular present indicative and subjunctive forms. This resulted in the generation of incorrect forms, such as *tiñió* (dyed) or *dirixió* (addressed), rather than the intended *tiñó* and *dirixó*. This was due to an erroneous assignment of the paradigm. It was thus necessary to create a specific paradigm for this type of verbs. Furthermore, a considerable number of Asturian verbs with enclitic pronouns were not correctly translated, resulting in the translation of their infinitive form instead of the expected conjugation. To address this issue, it was imperative to rectify the verb paradigms, which, due to inconsistencies with the paradigms of the Spanish dictionary, led to this

---

[13]https://eslema.it.uniovi.es/comun/traductor.php

[14]https://github.com/apertium/apertium-spa-ast

type of error.

Finally, the transfer rules for this pair were found to contain numerous errors, which resulted in a significant decline in the quality of the translation. These errors manifested in various ways, including inconsistencies in gender or number between related nouns and adjectives, the absence of verb conjugation or declension, incorrect disambiguation, and other issues that affected the whole translation.

The aforementioned examples illustrate the work that has been carried out on the monolingual and bilingual packages. As a result, 172 transfer selection rules were corrected for both translation directions and 11 paradigms for each monolingual dictionary.

## 5.2 Lexical Enhacement

The dictionaries have been expanded to include new terms drawn from a number of fields, including administration, toponymy from both Asturias and Spain, and anthroponymy. In total, 3100 new terms have been incorporated into the dictionaries, with the inclusion of each new term informed by the preferences of the Dictionary of the Asturian Academy (*Diccionariu de la Academia de la Llingua Asturiana*) (DALLA[15]). Indeed, one of the most significant alterations implemented at the generic level within the Apertium dictionaries has been the selection of the cultured form in lieu of the vocalised form of the term. To illustrate, the choice of *-pt-* instead of *-ut-* in terms such as *conceptu/conceutu*, the choice of *-ps-* instead of *-us-* in terms like *cápsula/cáusula*, the choice of *-cd-* instead of *-ud-* in terms such as *anécdota/anéuduta* and the preference for the intervocalic *-x-* rather than the *-s-* found in terms such as *exame/esame* are examples of the aforementioned changes. Similarly, the *-zar* ending is preferred for verbs such as *forzar/forciar*, in contrast to the *-ciar* ending. Otherwise, as mentioned above, the DALLA shape was always preferred in all cases where there were two possibilities.

## 6 Results

Table 1 presents the BLEU (Papineni et al., 2002) and chrF++ (Popović, 2015) scores received from the OCELoT system. The table includes the ten best systems presented to this shared task and our own system, identified by the ID 580.

---

[15]See:https://www.diccionariu.alladixital.org/

| ID | BLEU | chrF++ |
|---|---|---|
| **576** | **23.2** | **55.2** |
| 606 | 19.8 | 52.2 |
| 574 | 19.7 | 52.2 |
| 528 | 18.4 | 52.1 |
| 609 | 19.8 | 52.1 |
| 551 | 18.2 | 51.6 |
| 557 | 17.9 | 51.6 |
| 629 | 18.0 | 51.6 |
| 568 | 18.0 | 51.6 |
| 564 | 18.0 | 51.6 |
| 580 | 17.6 | 51.2 |

Table 1: The best 10 scores obtained in the OCELoT system in the WMT24 Shared Task: Low Resource Languages in Spain (Spanish–Asturian) and the enhaced Apertium system, ID 580.

## 7 Analysis

In order to elucidate the outcomes yielded by our system and the top ten in this shared task, it is essential to examine the functioning of the BLEU and chrF++ metrics, on the one hand, and the FLORES test, on the other.

From one perspective, BLEU and chrF++ are lexical-based metrics that rely on a reference corpus to assess the quality of a translation. In essence, BLEU assesses the quality of the system by comparing the MT output with the reference test token by token at sentence or corpus level (Papineni et al., 2002). In contrast, the chrF++ metric functions in a comparable manner, albeit by comparing character by character rather than token by token (Popović, 2015). Both metrics have been the subject of criticism on the grounds of their reliance on a reference corpus, which presents a significant challenge for low-resource languages, such as Asturian. In the absence of the requisite test datasets, the evaluation with these metrics is often impractical. Furthermore, these metrics fail to account for the inherent variability and versatility of languages. In many cases, multiple translations may be equally valid for a given source sentence. However, these metrics treat any deviation from the reference corpus as an error, leading to artificially low metrics when the deviation is linguistically correct (Lee et al., 2023).

In light of the aforementioned considerations, it is pertinent to highlight that the FLORES+ corpus, which serves as the basis for the evaluation of the systems in this shared task, comprises 3001

English sentences extracted from Wikimedia and translated manually by linguists into 200 minority languages, including Asturian. Subsequently, these translations are subjected to automatic revision and post-editing as required (Costa-jussà et al., 2022). It should be noted that the parallel corpora generated by this project are not direct translations. However, the translated text maintains the meaning of the original sentence while deviating from the structure of the source language. This may be due to the fact that the corpus was generated from English. Furthermore, as stipulated in the terms of reference for this shared task, the Asturian corpus has been duly revised by the Asturian Academy for use as a reference corpus in this shared task.

For illustrative purposes, three examples can be found in Table 2. In this table it can be found the original sentence in Spanish from the FLORES+ devtest in the first column, the original sentence in Asturian from the FLORES+ devtest in the second column and the version of the same sentence in English in the third column. This version in English is also taken from FLORES+ devtest. It is evident from these sentences that the translation from Spanish to Asturian is not a literal one, but rather a free rendering. In some cases, the meaning may even change. For instance, in the first sentence, the English meaning is retained in the Spanish sentence, but is lost in the Asturian translation. The direct translation of the sentence from Asturian to English would be: "They all ran back **when the accident happened**". "From where the accident had happened" and "when the accident happened" is not meaning the same. Furthermore, information can also be lost, as evidenced by the second sentence. Once more, the Asturian translation does not convey the same information as the original English sentence and the Spanish translation. In this instance, information is lost. Rather than referencing the navigable canals, the translation merely states that they are located inland, and instead of indicating that they are an optimal destination for "holidays", the translation simply uses the word *viaxes* (travels) which is not an equivalent expression. Furthermore, as evidenced in the third sentence, this phenomenon also occurs in the context of English–Spanish sentence translation. In such instances, the order of the sentence may undergo a change from English to Spanish, even when there is no necessity to align with the grammatical conventions of the target language. Even minor alterations such as this one have a deleterious impact on evaluation using lexical-based metrics.

These discrepancies within the parallel test corpora give rise to suboptimal results in metrics such as BLEU or chrF++, particularly in instances where the translated text may not be technically incorrect. This is not only the case for our system, but for all of them. The highest metric for BLEU is 23.2, while for chrF++ it is 55.2. These results are considerably low.

In regard to the results obtained by our system, it can be stated that Apertium, as a RBMT, produces translations that are literal in nature. In other words, unless it is a syntactic or grammatical feature intrinsic to the target language, the structure of the source language will always be replicated. In the case of Spanish and Asturian, which are two closely related languages, the MT output produced by Apertium will invariably adhere to the structure of Spanish, rather than exhibiting a more Asturian-specific structure. For illustrative purposes, consider the sentences presented in Table 3. This table presents the same sentences as in Table 2, with the second column displaying the translations generated by our system instead of the FLORES+ devtest Asturian sentences. The examples illustrate that the translation produced by Apertium preserves the structure of the source sentence in Spanish, but the translations are all accurate. It should be noted, however, that the system itself is not without limitations. It should first be noted that a word in Spanish has only one possible translation into Asturian, irrespective of context. While this can be managed in some specific cases, it may result in the translation failing in other sentences. It is possible that the inflexibility of this system, which does not always permit adaptation to context, may have had an adverse effect on our results, extending beyond the aforementioned metrics and the test employed. Nevertheless, we consider the output of our system to be a satisfactory translation that could be competitive with other systems despite the results. However, it would be necessary to carry out more tests in order to go deeper and identify the aspects in which the quality of our translation system could be improved, since the RBMT systems, as already mentioned, require constant revision and improvement.

Finally, and this is a strong point of our proposal, this type of system does not have a high computational consumption like Statistical Machine

| Original Sentences in Spanish | Original Sentences in Asturian | Original Sentences in English |
|---|---|---|
| *Todos volvieron corriendo desde el lugar del accidente.* | *Volvieron p'atrás corriendo cuando ocurrió l'accidente.* | They all ran back from where the accident had happened. |
| *Los canales navegables internos pueden ser una buena temática para las vacaciones.* | *Les canales d'interior son un bon tema de viaxe.* | Inland waterways can be a good theme to base a holiday around. |
| *En Inglaterra, las vías de tren ya se habían instalado hacia el siglo XVI.* | *Les primeres víes foron construyíes n'Inglaterra nel sieglu XVI.* | Wagonways were built in England as early as the 16th Century. |

Table 2: Illustrative sentences of the FLORES test dataset taken from the FLORES+ devtest in Spanish, Asturian and English languages.

| Original Sentences in Spanish | Apertium MT output | Original Sentences in English |
|---|---|---|
| *Todos volvieron corriendo desde el lugar del accidente.* | *Toos volvieron corriendo dende'l llugar del accidente.* | They all ran back from where the accident had happened. |
| *Los canales navegables internos pueden ser una buena temática para las vacaciones.* | *Les canales navegables internos puen ser una bona temática pa les vacaciones.* | Inland waterways can be a good theme to base a holiday around. |
| *En Inglaterra, las vías de tren ya se habían instalado hacia el siglo XVI.* | *N'Inglaterra, les víes de tren yá s'instalaren escontra'l sieglu XVI.* | Wagonways were built in England as early as the 16th Century. |

Table 3: Apertium MT output from the Spanish–Asturian translation of three sentences taken from FLORES+ devtest in their Spanish and English version.

Translation (SMT) and Neural Machine Translation (NMT) in its training, as signalled by Shterionov and Vanmassenhove (2023). Comparing the quality produced by this system with its consumption, both for training/development and for use, Apertium is still a competitive system for low-resource languages such as Asturian. Furthermore, it is also more cost-effective to produce than other types of systems. Therefore, it is essential to consider the trade-off between quality, consumption and price in order to assess the performance of the different systems.

## 8   Conclusions

In conclusion, although our enhanced Apertium system has not yet achieved a position among the top ten systems in this shared task, the results obtained in terms of machine translation quality have been exemplary. As previously stated in the Section 7, the test employed and the metrics utilized do not permit an accurate assessment of the quality of a MT system. Additionally, it is noteworthy that a RBMT exhibits a markedly lower consumption rate in comparison to NMTs. Consequently, we regard our system as being competitive with those submitted to this shared task, although it still necessitates further enhancements and revisions.

## References

Eckhard Bick and Tino Didriksen. 2015. Cg-3—beyond classical constraint grammar. In *Proceedings of the 20th Nordic Conference of Computational Linguistics (NODALIDA 2015)*, pages 31–39.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejía González, Prangthip Hansanti, John Hoffman, Semarley Jarret, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Mikel L Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine translation*, 25:127–144.

Inaciu Galán y González. 2015. Asturianu sos: una güeyada sobre la situación de la llingua asturiana y les sos perspectives de futuru. *Luenga & fablas: publicazión añal de rechiras, treballos e decumentazión arredol de l'aragonés ea suya literatura*, (19):67–72.

Tanmai Khanna, Jonathan N Washington, Francis M Tyers, Sevilay Bayatlı, Daniel G Swanson, Tommi A

Pirinen, Irene Tang, and Hector Alos i Font. 2021. Recent advances in apertium, a free/open-source rule-based machine translation platform for low-resource languages. *Machine Translation*, 35(4):475–502.

Seungjun Lee, Jungseob Lee, Hyeonseok Moon, Chanjun Park, Jaehyung Seo, Sugyeong Eo, Seonmin Koo, and Heuiseok Lim. 2023. A survey on evaluation metrics for machine translation. *Mathematics*, 11(4):1006.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Dimitar Shterionov and Eva Vanmassenhove. 2023. *The Ecological Footprint of Neural Machine Translation Systems*, volume 4, pages 185–213. Springer Nature Switzerland AG, Switzerland. 25 pages, 3 figures, 10 tables Copyright © 2023, The Author(s), under exclusive license to Springer Nature Switzerland AG.

Universidad de Oviedo. 2010. Traductor automático castellano-asturiano-castellano algunos datos.

Xulio Viejo, Roser Saurı, and Angel Neira. 2008. Eslema. towards a corpus for asturian. In *Collaboration: Interoperability between people in the creation of language resources for less-resourced languages. A SALTMIL workshop*.

# Universitat d'Alacant's Submission to the WMT 2024 Shared Task on Translation into Low-Resource Languages of Spain

**Aarón Galiano-Jiménez,**[†] **Víctor M. Sánchez-Cartagena,**[†]
**Juan Antonio Pérez-Ortiz,**[*†] **Felipe Sánchez-Martínez**[†]

[†]Dep. de Llenguatges i Sistemes Informàtics, Universitat d'Alacant, Spain

[*]Valencian Graduate School and Research Network of Artificial Intelligence, ValgrAI

{aaron.galiano,vm.sanchez,japerez,fsanchez}@ua.es

## Abstract

This paper describes the submissions of the Transducens group of the Universitat d'Alacant to the WMT 2024 Shared Task on Translation into Low-Resource Languages of Spain; in particular, the task focuses on the translation from Spanish into Aragonese, Aranese and Asturian. Our submissions use parallel and monolingual data to fine-tune the NLLB-1.3B model and to investigate the effectiveness of synthetic corpora and transfer-learning between related languages such as Catalan, Galician and Valencian. We also present a many-to-many multilingual neural machine translation model focused on the Romance languages of Spain.

## 1 Introduction

Spain is home to several languages, each with different levels of representation in neural machine translation (NMT) technologies and availability of training data. For example, Spanish (spa) has abundant data resources and is included in many multilingual translation models and large language models (LLM). Other languages, such as Catalan (cat) and Galician (glg), are relatively well-supported and have enough data to train NMT models from scratch. However, languages such as Asturian (ast), Aragonese (arg) and Aranese (arn) face significant challenges due to the limited availability of data needed to train these systems.

Despite these challenges, the linguistic similarity between some of these languages simplifies their integration into multilingual translation models, allowing them to benefit from transfer-learning from more widely represented languages; an example of this is the inclusion of Asturian in NLLB-200 (NLLB Team et al., 2022). In addition, shallow-transfer rule-based machine translation (MT) systems such as Apertium (Forcada et al., 2011), exist for some of these languages, includ-



Figure 1: Submitted models for the shared task. Blue in *src* or *tgt* indicates text generated by MT. The Denoising and Mix models were trained for a single translation direction, whereas the Many2Many model was trained with multiple language pairs and in both translation directions for each pair. xxx represents any of the target languages: Aragonese, Aranese and Asturian.

ing Spanish–Asturian[1], Spanish–Aragonese[2] and Occitan–Spanish[3]; the development of rule-based systems does not require large amounts of training data, but linguistic knowledge to construct dictionaries and translation rules.

Our approach to developing NMT systems for the WMT 2024 Shared Task on Translation into Low-Resource Languages of Spain (Sánchez-Martínez et al., 2024) for Aragonese, Aranese and Asturian is based on pre-trained models that include similar languages, such as NLLB-200, to incorporate languages that were not originally seen during training. Given the scarcity of corpora for Asturian, Aragonese and Aranese, we used Apertium to generate synthetic corpora. This involved translating monolingual corpora into Spanish and vice versa to generate additional resources for fine-

---

[1]https://github.com/apertium/apertium-spa-ast
[2]https://github.com/apertium/apertium-spa-arg
[3]https://github.com/apertium/apertium-oci-spa

tuning NLLB-200.

By combining pre-trained multilingual models, synthetic data generation and rule-based translation systems, we aim to improve the quality and accessibility of NMT for these low-resource languages.

All of our submissions, shown in Figure 1, are classified as *open* because they are based on the NLLB-200 model with 1.3B parameters, exceeding the 1B parameters limit for *constrained* submissions. However, we only used the corpora allowed for constrained submissions, as described in Section 3.

## 2 State-of-the-Art Methods Used

Our submission makes use of well-established techniques such as denoising pre-training (Lewis et al., 2020), transfer-learning (Zoph et al., 2016), back-translation (Sennrich et al., 2016) and sequence-level knowledge distillation (Kim and Rush, 2016). The languages involved in this shared task are related, making knowledge transfer by a multilingual model an effective solution. This transfer can be achieved by bilingual fine-tuning of a pre-trained model (Zoph et al., 2016), using the knowledge it already possesses, or by multilingual training from scratch (Bommasani et al., 2021), using multiple languages simultaneously during the training process. In our approach, we fine-tune the NLLB-200 model and investigate the effects of training with only one specific translation direction with different types of data, as well as adding more languages during training.

A common technique is to pretrain a system on monolingual corpora with the denoising task[4] to learn language generation, and then train the system on parallel bilingual corpora for translation (Lewis et al., 2020). It is well known that combining both tasks simultaneously improves the translation results (Kamboj et al., 2022). Since NLLB-200 was trained in this way (NLLB Team et al., 2022), it is reasonable to use the same technique for fine-tuning to leverage the available monolingual corpus.

Another method of exploiting monolingual corpora is to create synthetic parallel corpora. For this purpose, we used the rule-based MT systems built using the Apertium platform (Forcada et al., 2011).

| Corpus | Sentences | Src words | Tgt words |
|---|---|---|---|
| spa-arg | 33,723 | 3,706,154 | 3,589,002 |
| spa-arn | 85,491 | 14,720,677 | 14,266,772 |
| spa-ast | 45,506 | 6,844,424 | 6,663,424 |
| spa | 500,000 | 62,004,331 | — |
| arg | 24,675 | 2,718,855 | — |
| arn | 229,886 | 29,110,670 | — |
| ast | 38,868 | 5,504,371 | — |
| spa-cat | 559,805 | 91,543,160 | 88,057,754 |
| spa-glg | 184,861 | 30,716,538 | 28,753,332 |
| spa-val | 287,403 | 52,836,299 | 53,137,411 |

Table 1: Number of sentences and words in each of the corpora used.

Specifically, we used Catalan–Spanish, Aragonese–Spanish, Aranese–Spanish, Spanish–Aragonese, Spanish–Aranese and Spanish–Asturian systems. Translating from source (spa) to target and using this synthetic corpus as a target for training is a type of sequence-level knowledge distillation (Lai et al., 2021; Yu et al., 2021). In contrast, translating from the target language and using this synthetic corpus as the source for training is called back-translation (Sennrich et al., 2016). The latter has the advantage that potential translation errors in the synthetic corpus do not affect the generation of the target language, as the synthetic corpus is used as input during training rather than as the desired output.

## 3 Data

We used only the corpora allowed for the constrained submissions: Opus[5] and PILAR (Galiano-Jiménez et al., 2024)[6]. For the development set, we used the FLORES+[7] (NLLB Team et al., 2022) dev versions (997 sentences) for Spanish (spa), Aranese (arn), Aragonese (arg) and Asturian (ast) (Pérez-Ortiz et al., 2024). The specific details of the corpora used are described below and shown in Table 1.

### 3.1 Parallel Corpora

**Aragonese and Asturian:** We used parallel corpora with Spanish available in OPUS,[8] consisting of 33,723 Spanish–Aragonese sentences and 45,506 Spanish–Asturian sentences.

---

[4]The denoising task is a self-supervised learning strategy that helps models learn effective representations from monolingual data by training them to restore original sentences from corrupted inputs.

[5]https://opus.nlpl.eu/
[6]https://github.com/transducens/PILAR
[7]https://github.com/openlanguagedata/flores
[8]https://opus.nlpl.eu/

**Aranese:** We used the parallel Catalan-Aranese corpus available in PILAR (85,491 sentences) and translated the Catalan part into Spanish using Apertium. This corpus was not included in the synthetic data description in section 3.3, as it was not automatically translated from or into Aranese.

## 3.2 Monolingual Corpora

We used the monolingual corpora available in PILAR, which contains 24,675 sentences in Aragonese, 229,886 in Aranese and 38,868 in Asturian.

## 3.3 Synthetic Corpora Generation

We used Apertium to generate different synthetic corpora. By translating the above monolingual corpora into Spanish, we created two corpora (there is no Asturian–Spanish system for Apertium) that were used for back-translation. Since these corpus sizes are relatively small for training NMT models, we also translated 500,000 Spanish sentences extracted from Paracrawl[9] into Aranese[10], Aragonese[11] and Asturian[12]. This resulted in synthetic corpora with the target language as the synthetic part.

## 3.4 Additional data

For our multilingual system, we used all the above corpora and added Wikimedia Spanish–Catalan (559,805 sentences) and Spanish–Galician (184,861 sentences) corpora from OPUS, as well as the Spanish–Valencian (val) corpus available in PILAR, making 287,403 sentences.

## 4 Methodology

We trained several models for each language pair to analyse the effects of different types of corpora and transfer-learning between different tasks and languages. Below we describe the methodology we used and the different systems we trained.

## 4.1 Model Architecture and Baseline

All our models are based on a fine-tuning of the NLLB-200 (NLLB Team et al., 2022) model with

---

9[https://opus.nlpl.eu/ParaCrawl/es&ca/v9/ParaCrawl](https://opus.nlpl.eu/ParaCrawl/es&ca/v9/ParaCrawl)

10[https://github.com/apertium/apertium-oc-es/releases/tag/v1.0.8](https://github.com/apertium/apertium-oc-es/releases/tag/v1.0.8)

11[https://github.com/apertium/apertium-spa-arg/releases/tag/v0.6.0](https://github.com/apertium/apertium-spa-arg/releases/tag/v0.6.0)

12[https://github.com/apertium/apertium-spa-ast/releases/tag/v1.1.1](https://github.com/apertium/apertium-spa-ast/releases/tag/v1.1.1)

---

1.3 billion parameters. We chose this model because it is a transformer (Vaswani et al., 2017) pre-trained with 200 languages and specialised in translation tasks. These 200 languages include Spanish and Asturian, which correspond to one of the translation directions in this shared task, as well as other related languages, such as Catalan, Galician and Occitan.

As a baseline, we used Apertium to compare the effectiveness of a rule-based system with a neural-based system.

## 4.2 Training Approaches

In this section, we describe the different approaches we used to train the translation models for each language pair.

**Bilingual parallel:** For each translation direction, we trained a specific model using only the parallel corpus available for that language pair. These models correspond to the `Parallel` row in tables 2 and 3.

**Bilingual Parallel + Monolingual:** We trained a model for each translation direction by combining the translation task using the parallel corpus with a denoising task using monolingual data of the target language. This approach helps to improve translation quality by exploiting additional monolingual resources. These models correspond to the `Denoising` row in tables 2 and 3.

**Bilingual Synthetic Generated with Apertium:** For each translation direction, we trained a model using only synthetic parallel corpora generated by translating the Spanish monolingual corpora into the target language using Apertium. These models correspond to the `Synthetic` row in tables 2 and 3.

**Bilingual Parallel with Synthetic and Back-translation:** For each translation direction, we trained a model using a combination of parallel corpora, synthetic corpora generated with Apertium, and back-translation. For back-translation, we used Apertium to translate the monolingual target language data into Spanish. These models correspond to the `Mix` row in tables 2 and 3.

**Multilingual Parallel with Synthetic and Back-translation:** This model extends the previous approach by training a single model on all three translation directions. This multilingual training allows

the model to benefit from common linguistic features across languages. This model corresponds to the `Multilingual` row in tables 2 and 3.

**Multilingual Many-to-Many:** This model extends the multilingual approach by incorporating parallel corpora from related languages and including both translation directions for each language pair. This model can translate between Spanish, Aragonese, Aranese, Asturian, Catalan, Galician and Valencian. This model corresponds to the `Many2Many` row in tables 2 and 3.

### 4.3 Evaluation Metrics

To measure the quality of the translation models, we used BLEU (Papineni et al., 2002)[13] and chrF2 (Popović, 2015)[14] scores on the translation of the development set. For the NLLB-200 based models, we translated using beam search (Graves, 2012) with a beam size of 5.

## 5 Experiments and Results

In this section, we present the experimental setup, including hyperparameters and training configurations. We compare the performance of each system against the baseline provided by Apertium and discuss the results for each translation direction.

### 5.1 Experimental Setup

All our translation models were trained using the Transformers (Wolf et al., 2020) library on an Nvidia A100 GPU with 40 GB of memory. We extended the library to support simultaneous training with multiple tasks, including different datasets with different languages and combining translation and denoising tasks[15]. To balance the training data from datasets of different sizes, we used temperature upsampling with $1/T = 0.3$, following the approach used in NLLB-200 training (NLLB Team et al., 2022).

Due to GPU memory constraints, we used a batch size of 16 and accumulated gradients as many times as the number of different datasets used in the training to ensure that all tasks were seen before updating the model weights. The learning rate was set to 5e-5, and we used the AdamW optimizer (Loshchilov and Hutter, 2017) with $\beta_1$=0.9,

| Model | spa-arg | spa-arn | spa-ast |
|---|---|---|---|
| Apertium | 66.0 | 38.0 | 17.1 |
| Parallel | 41.4 | 34.4 | 17.9 |
| Denoising* | 41.6 | 35.7 | 17.8 |
| Synthetic | 65.3 | 37.6 | 17.0 |
| Mix* | 65.1 | 37.8 | 17.0 |
| Multilingual | 64.8 | 37.5 | 17.0 |
| Many2Many* | 65.2 | 37.9 | 17.0 |

Table 2: BLEU scores on the FLORES+ dev. Models marked with an asterisk are those we submitted for the Shared Task.

$\beta_2$=0.999 and $\varepsilon$=$10^{-8}$. Models were trained for a maximum of 100 epochs with early stopping, and evaluations were performed every 1000 training steps. The stopping criterion was based on the BLEU score on the development set, with a patience of 6 evaluations.

We used the NllbTokenizer[16] class for corpus segmentation. This tokenizer uses the Sentence-Piece (Kudo and Richardson, 2018) model used by NLLB-200 and applies language tokens to both the source and target texts.

NLLB-200 uses language tokens in both the source and target sentences. When training with a new language, it is possible to use the language token of a similar language, but this eliminates the possibility of translating to or from the language of the original token. In our case, we added new language tokens for Aragonese, Aranese and Valencian. To avoid learning the embeddings for these tokens from scratch, we initialised them with the embeddings of the most similar languages included in NLLB-200. Specifically, we initialized the embeddings for Aragonese and Valencian with the Catalan embedding, and the embedding for Aranese with the Occitan embedding[15].

### 5.2 Results

Tables 2 and 3 show the results of the translation models in terms of BLEU and chrF scores, respectively. The first row corresponds to the Apertium baseline and the remaining rows show the results of each trained model.

For the shared task, we submitted specific Denoising and Mix models for each translation direction, and the Many2Many model for all three

---

[13] SacreBLEU BLEU signature: `nrefs:1 | case:mixed | eff:no | tok:13a | smooth:exp | version:2.0.0`

[14] SacreBLEU chrF2 signature: `nrefs:1 | case:mixed | eff:yes | nc:6 | nw:0 | space:no | version:2.0.0`

[15] The code is available at `https://github.com/transducens/transformers-multilingual-training`

[16] `https://github.com/huggingface/transformers/blob/v4.42.0/src/transformers/models/nllb/tokenization_nllb.py`

| Model | spa-arg | spa-arn | spa-ast |
|---|---|---|---|
| Apertium | 82.2 | 60.0 | 50.7 |
| Parallel | 70.8 | 57.9 | 50.8 |
| Denoising* | 69.5 | 58.6 | 50.7 |
| Synthetic | 81.9 | 59.9 | 50.7 |
| Mix* | 81.8 | 59.9 | 50.8 |
| Multilingual | 81.8 | 59.8 | 50.8 |
| Many2Many* | 81.9 | 60.0 | 50.8 |

Table 3: chrF2 scores on the FLORES+ dev. Models marked with an asterisk are those we submitted for the Shared Task on Translation into Low-Resource Languages of Spain.

| Model | spa-arg | spa-arn | spa-ast |
|---|---|---|---|
| Denoising | 37.8 | 27.0 | 17.4 |
| Mix | 60.2 | 28.5 | 16.9 |
| Many2Many | 59.8 | 28.5 | 16.8 |

Table 4: BLEU scores on the WMT 2024 Shared Task on Translation into Low-Resource Languages of Spain test.

directions [17]. This allows us to compare a model trained only on the available corpus for each language, another that incorporates a synthetic corpus, and one that includes multiple translation directions and additional languages. The results on the test set of the task, which correspond to the FLORES+ devtest versions of these languages, are shown in Tables 4 and 5.

When analysing the results, it is important to consider how the development sets were created (Pérez-Ortiz et al., 2024). The Asturian sentences were first professionally translated from English by Meta (NLLB Team et al., 2022) and then revised by academics. In contrast, the Aragonese and Aranese sentences were first machine translated from Spanish using Apertium, then manually edited by language specialists and finally reviewed by academics. This means that the development sets for Aragonese and Aranese may be biased towards the results produced by Apertium.

We conducted paired significance tests to determine whether the submitted models outputs were significantly different despite the similarity of some results. Specifically, we calculated paired approximate randomisation (Riezler and Maxwell, 2005) as implemented by SacreBLEU on the devtest using BLEU and chrF2. The results indicated that the

---

[17]The Many2Many model is available at https://huggingface.co/Transducens/IbRo-nllb

| Model | spa-arg | spa-arn | spa-ast |
|---|---|---|---|
| Denoising | 67.5 | 48.3 | 50.7 |
| Mix | 78.9 | 49.3 | 50.9 |
| Many2Many | 78.8 | 49.3 | 50.9 |

Table 5: chrF2 scores on the WMT 2024 Shared Task on Translation into Low-Resource Languages of Spain test.

differences between the Mix and Many2Many models for Asturian and Aranese were not statistically significant, whereas the differences between all the other pairs of models were statistically significant.

**Effect of adding a monolingual corpus:** The results show a minimal difference when adding a monolingual corpus compared to training only with parallel corpora (rows 2 and 3 of Tables 2 and 3). However, this could be due to the amount of training data available. The improvement in the quality of the translations for the Aranese direction is particularly remarkable, since it is the largest monolingual corpus.

**Effect of using the synthetic corpus produced by Apertium:** The increase in performace for the synthetic models compared to the Denoising models for Aragonese and Aranese can be explained both by the difference in data volume and by the bias of the development sets. Conversely, there is a decrease in the results for Asturian, suggesting a bias in the other sets.

The combination of parallel and synthetic corpora in both target and source (Mix models) shows minimal variation in the results. Again, this may be due to the difference in the proportion of the corpus generated by the spa-xxx translation and that generated by the xxx-spa translation.

**Effect of multilingual training:** Combining multiple translation directions in the same training session complicates the learning task for the model, but also increases the amount of data available. Adding more languages slightly improves the results compared to training with only the languages of the common task.

## 6 Conclusions

Overall, the results highlight the critical role of training data volume in the development of effective NMT models. The challenge with large neural models lies in the insufficient amount of training

data available for low-resource languages, which limits the full potential of these architectures.

However, rule-based systems remain a viable option for these languages, although they require linguistic expertise to build. The use of these systems to generate synthetic corpora is proving beneficial in integrating low-resource languages into neural translation models and exploiting the advantages they offer.

## Acknowledgments

## References

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, and et al. 2021. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258.

Mikel L. Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Gema Ramírez-Sánchez Felipe Sánchez-Martínez, and Francis M. Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine Translation*, 25(2):127–144. Special Issue: Free/Open-Source Machine Translation.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Pan-iberian language archival resource.

Alex Graves. 2012. Sequence transduction with recurrent neural networks. *CoRR*, abs/1211.3711.

Samta Kamboj, Sunil Kumar Sahu, and Neha Sengupta. 2022. DENTRA: Denoising and translation pre-training for multilingual machine translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1057–1067, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Yoon Kim and Alexander M. Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327, Austin, Texas. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Wen Lai, Jindřich Libovický, and Alexander Fraser. 2021. The LMU Munich system for the WMT 2021 large-scale multilingual machine translation shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 412–417, Online. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *arXiv*, abs/2207.04672.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia,

Pennsylvania, USA. Association for Computational Linguistics.

Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Víctor M. Sánchez-Cartagena, Miquel Esplà-Gomis, Aarón Galiano-Jiménez, Antoni Oliver, Claudi Aventín-Boya, Cristina Valdés, Alejandro Pardos, and Juan Pablo Martínez. 2024. FLORES+ datasets for Aragonese, Aranese, Asturian and Valencian. In *Proceedings of the Ninth Conference on Machine Translation*, Miami. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Stefan Riezler and John T. Maxwell. 2005. On some pitfalls in automatic evaluation and significance testing for MT. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 57–64, Ann Arbor, Michigan. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Felipe Sánchez-Martínez, Juan Antonio Pérez-Ortiz, Aarón Galiano-Jiménez, and Antoni Oliver. 2024. Findings of the WMT 2024 Shared Task on Translating into Low-Resource Languages of Spain: Blending rule-based and neural systems. In *Proceedings of the Ninth Conference on Machine Translation (WMT24)*, Miami, Florida, USA. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Zhengzhe Yu, Daimeng Wei, Zongyao Li, Hengchao Shang, Xiaoyu Chen, Zhanglin Wu, Jiaxin Guo, Minghan Wang, Lizhi Lei, Min Zhang, Hao Yang, and Ying Qin. 2021. HW-TSC's participation in the WMT 2021 large-scale multilingual translation task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 456–463, Online. Association for Computational Linguistics.

Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1568–1575, Austin, Texas. Association for Computational Linguistics.

# Samsung R&D Institute Philippines @ WMT 2024 Low-resource Languages of Spain Shared Task

**Dan John Velasco**[★,a]   **Manuel Antonio Rufino**[★,a]   **Jan Christian Blaise Cruz**[a,b]

[a]Samsung R&D Institute Philippines   [b]MBZUAI

{dj.velasco,ma.rufino}@samsung.com, jan.cruz@mbzuai.ac.ae

[★]**Equal Contribution**

## Abstract

This paper details the submission of Samsung R&D Institute Philippines (SRPH) Language Intelligence Team (LIT) to the WMT 2024 Low-resource Languages of Spain shared task. We trained translation models for Spanish to Aragonese, Spanish to Aranese/Occitan, and Spanish to Asturian using a standard sequence-to-sequence Transformer architecture, augmenting it with a noisy-channel reranking strategy to select better outputs during decoding. For Spanish to Asturian translation, our method reaches comparable BLEU scores to a strong commercial baseline translation system using only constrained data, backtranslations, noisy channel reranking, and a shared vocabulary spanning all four languages.

## 1   Introduction

This paper details our constrained system for translating from Spanish to Aragonese (spa→arg), Aranese/Occitan (spa→arn), and Asturian (spa→ast) for the WMT24 Shared Task: Translation into Low-Resource Languages of Spain. We trained standard sequence-to-sequence Transformer architecture (Vaswani et al., 2017) from scratch combined with heavy data preprocessing (Cruz, 2023), data augmentation via backtranslation (Sennrich et al., 2016a), and noisy channel reranking (Yee et al., 2019) to achieve performance that is comparable to Apertium (Khanna et al., 2021) v3.9.6 for spa→ast. We present ablation results of the effect of data augmentation via backtranslation and noisy channel reranking with respect to BLEU scores. Furthermore, we analyzed the generated translations and we found that the model learned to regurgitate, i.e. repeat with minor modifications, the source Spanish sentences for the spa→ast case. We also identify rarely occurring characters that the model failed to learn. Lastly, we

also investigated the effect of the length of back-translated data on improving model performance.

## 2   Methodology

### 2.1   Environment

For preprocessing, training, and generation, we used fairseq 0.12.2 and PyTorch 1.12.1. The training was done on either 2x NVIDIA Quadro GPUs or 8x NVIDIA P100 GPUs. We used Apertium[1] v3.9.6 for generating baseline results and generating backtranslated (BT) data whenever available for the language pair.

### 2.2   Data Preprocessing

We trained on the OPUS dataset (Tiedemann, 2016) for all language pairs. The data preprocessing pipeline utilizes the ratio-based filters and embedding-based filters of Samsung R&D Institute Philippines' WMT23 entry (Cruz, 2023). The dataset statistics before and after preprocessing can be found in Table 1.

For the parallel data, the data preprocessing pipeline are as follows: remove exact duplicate parallel data → ratio-based filters → embeddings-based filters. The ratio-based filters remove sentences based on sentence length, token length, character to token ratio, pair token ratio, and pair length ratio. Exact details on these criteria are explained in (Cruz, 2023). Similar to last year's paper, we tokenized and detokenized sentences using Sacre-Moses[2] before and after running our filters, respectively. The embeddings-based filter filters data based on the cosine similarity of a sentence pair using LaBSE (Feng et al., 2022). Using the methodology of (Cruz, 2023), pairs with a cosine similarity $0.7 \leq s \leq 0.96$ are kept.

For monolingual data, we combined the monolingual data of the target language and the target

---

[b]Work done while at Samsung R&D Institute Philippines

[1]https://wiki.apertium.org/wiki/Install_Apertium_core_using_packaging

[2]https://github.com/alvations/sacremoses

| source→target | Pairs | Words (source) | Words (target) | % Kept |
|---|---|---|---|---|
| spa→arg | 58,284 | 746,567 | 733,985 | 100 |
| spa→arg Filtered | 21,362 | 181,523 | 190,724 | 36.6 |
| spa→arg Filtered + BT[†] | 81,195 | 849,031 | 857,111 | - |
| spa→ast | 13,393,052 | 310,197,263 | 298,687,582 | 100 |
| spa→ast Filtered | 620,168 | 6,495,284 | 6,442,051 | 4.6 |
| spa→ast Filtered + BT | 920,168 | 11,758,510 | 11,629,822 | - |
| spa→arn | 2,032,440 | 14,046,448 | 13,891,193 | 100 |
| spa→arn Filtered | 779,615 | 4,807,268 | 5,020,187 | 38.4 |
| spa→arn Filtered + BT[†] | 1,079,615 | 8,835,974 | 9,040,705 | - |

Table 1: Statistics of OPUS parallel data before and after filtering and the addition of backtranslated data (BT). The % Kept is the percentage of pairs left after filtering ("-" means not applicable). † means BT data was generated via Apertium.

side of parallel data from OPUS and then removed exact duplicates. We used this monolingual data to train language models for each target language.

After preprocessing of parallel and monolingual data, we apply train and validation split of 95% and 5%, respectively.

Lastly, for the training corpus of the tokenizer, we combined the filtered parallel data of all three language pairs. We used this combined data to learn a shared BPE (Sennrich et al., 2016b) vocabulary that spans Spanish, Aragonese, Aranese, and Asturian consisting of 31,960 tokens using SentencePiece (Kudo and Richardson, 2018). This shared vocabulary was used by all models for generating submissions to WMT24. We used this approach as the four languages belong to the same language family.

## 2.3 Augmenting Data with Backtranslation

We augmented the filtered training data using backtranslation (Sennrich et al., 2016a). For each language pairs for both source→target (except spa→ast) and target→source directions, Apertium 3.9.6 was used to generate BT data. Due to the lack of direct translation support for some language pairs in Apertium, the translation for arg→spa went through the following translation path: Aragonese → Catalan → Interlingua → Spanish[3]. For arn→spa, it goes through Aranese → Catalan → Spanish[4].

Translation from Asturian to Spanish is not supported by Apertium. Alternatively, we used the ast→spa model that was originally intended for noisy channel reranking (NCR), a technique which will be explained in Section 2.5, to generate BT

data. For decoding, we used combined top-k and nucleus sampling:

$$\sum_{i=0}^{\delta_k} P(\hat{y}_i^{(T)}|x; \hat{y}^{(T-1)}) \cdot \delta_{temp} \leq \delta_p \qquad (1)$$

where $\delta_k$ is the top values considered for top-k sampling, $\delta_{temp}$ is temperature, $\delta_p$ is the maximum total probability for nucleus sampling. For these hyperparameters, we used the same values as (Cruz, 2023) which are as follows: $\delta_k = 50$, $\delta_{temp} = 0.7$, and $\delta_p = 0.93$.

Once the BT data for each language pairs and translation direction are generated we took a subset in different ways. For BT data for training Direct Translation Models (spa→arg/ast/arn), we used all the generated BT data for spa→arg since it's less than 300K. For BT data of spa→ast and spa→arn, we keep the longest 300K sentences.

For BT data for training Channel Translation Models (arg/ast/arn→spa), we used all the BT data for arg→spa since it's less than 100K. For ast→spa, we randomly sampled 100K sentences. Due to time constraints, we did not generate BT data for arn→spa.

## 2.4 Model Training

For each language pair, we trained three types of models: a **Direct Translation Model**, a **Channel Model**, and a **Language Model** which will be detailed in the following subsections. These three models will be combined via Noisy Channel Reranking (Yee et al., 2019) which will be explained in Section 2.5.

### 2.4.1 Direct Translation Models

For each direct translation models (spa→arg/arn/ast), we trained encoder-decoder

---

[3]Apertium language codes: arg-cat→cat-ina→ina-spa
[4]Apertium language codes: oc_aran-ca→cat-spa

Transformer architecture (Vaswani et al., 2017) from scratch with and without BT data. We used the large variant of transformers which has 213M parameters[5]. We describe two training configurations: **tf-large60k** which was trained for 60,000 steps of which 3,000 are warmup steps, and **tf-large100k** which was trained for 100,000 steps of which 10,000 are warmup steps. Training settings with "-plusbt" suffix indicates that the model was trained on a mixture of provided training data and BT data. Otherwise, it indicates the model is trained only on the provided training data. For example, **tf-large100k-plusbt** means the model was trained on the mixture of provided training data and BT data for 100,000 steps of which 10,000 are warm up steps.

For both settings and all language directions, unless stated otherwise, we used the same hyperparameters in Table 2. For generating WMT24 submissions, we used models trained on **tf-large100k-plusbt** setting as our Direct Translation Model.

### 2.4.2 Channel Translation Models

For the channel translation models (arg/arn/ast→spa), we used the same architecture and hyperparameters as the direct translation models, except it was trained on **tf-large60k-plusbt** setting, batch size/max tokens of 10,000, and learning rate of 7e-4 (arg→spa and arn→spa) and 5e-5 (ast→spa). These were used as channel models for noisy channel reranking which is explained further in Section 2.5 and for performing hyperparameter sweeps of noisy channel reranking parameters detailed in Section 2.6.

### 2.4.3 Language Models

We trained monolingual language models for Aragonese, Aranese, and Asturian from scratch using the decoder-only part of the original Transformer architecture as described in (Vaswani et al., 2017). We used the base variant which has 65M parameters[6]. For all languages, we used Adam optimizer (Kingma and Ba, 2017) with $\beta_1$=0.90, $\beta_2$=0.98. We trained for a maximum of 250,000 steps of which 4,000 are warmup steps. The warmup initial learning rate is 1e-7 and the max learning rate is 5e-4 and then decayed following an Inverse Square root learning rate schedule. The batch size / max tokens is 40,000, and the dropout

---

[5]Fairseq model code: transformer_wmt_en_de_big
[6]Fairseq model code: transformer_lm

| Training Hyperparameters | |
|---|---|
| Vocab Size | 31,960 |
| Tied Weights | Yes |
| Dropout | 0.3 |
| Attention Dropout | 0.1 |
| Weight Decay | 0.0 |
| Label Smoothing | 0.1 |
| Optimizer | Adam |
| Adam Betas | $\beta_1$=0.90, $\beta_2$=0.98 |
| Adam $\epsilon$ | $\epsilon$=1e-6 |
| Learning Rate | 5e-5 |
| LR Schedule | Inverse Sqrt |
| Batch Size | 8,000 tokens |

Table 2: Fixed hyperparameters for direct translation models.

is 0.1. These models were used in noisy channel reranking which is explained further in Section 2.5 and for performing hyperparameter sweeps of noisy channel reranking parameters detailed in Section 2.6.

### 2.5 Noisy-Channel Reranking

Similar to (Cruz, 2023), we experimented with using Noisy Channel Reranking (Yee et al., 2019) to improve translations. This works by using a direct translation model (source→target), channel model (target→source) and a monolingual language model (target only) to rescore every candidate translation token during beam search decoding. The score of the candidate translation token $\hat{y}_i^{(T)}$ at time step $T$ is recomputed using a linear combination of all three models:

$$P(\hat{y}_i^{(T)}|x; \hat{y}^{(T-1)})' = \frac{1}{t} \log(P(y|\hat{x}^{(T-1)}) + \frac{1}{s}[\delta_{ch} \log(P(x|\hat{y}^{(T-1)}) \quad (2) + \delta_{lm} \log(P(\hat{y}^{(T-1)}))]$$

where $t$ is the length of target sentence $y$ and $s$ is the length of source sentence $x$ which serves as debiasing terms. The $\delta_{ch}$ and $\delta_{lm}$ are weights of the channel model and language model, respectively, which controls the influence of the models to the final score. For this paper, both $\delta_{ch}$ and $\delta_{lm}$ were set to 0.5

### 2.6 Hyperparameter Sweeping

Similar to (Cruz, 2023), we utilized a Bayesian hyperparameter search to find an optimal value for

| Setting | BLEU | | | | | |
|---------|------|------|------|------|------|------|
| | FLORES+ dev | | | WMT24 Test | | |
| | spa→arg | spa→ast | spa→arn | spa→arg | spa→ast | spa→arn |
| Apertium (baseline) | **70.3** | 22.6 | **42.4** | - | - | - |
| No BT; No NCR | 18.3 | 23.9 | 8.7 | 13.4 | 16.8 | 7.2 |
| No BT; w/ NCR | 21.3 | 24.0 | 8.7 | 16.5 | 16.9 | 7.2 |
| w/ BT; No NCR | 35.4 | **24.4** | 14.4 | 26.7 | **17.5** | **7.7** |
| w/ BT; w/ NCR | 37.1 | 24.3 | 13.7 | **28.2** | 17.2 | 7.2 |

Table 3: BLEU scores of various system configurations compared to Apertium. BT and NCR denotes backtranslated data and noisy channel reranking, respectively. Highest score per language pair are in bold.

| Model configuration | BLEU |
|---------------------|------|
| tf-base100k w/o NCR | 19.3 |
| tf-base100k w/ NCR | 20.7 |
| tf-base100k-plusbt w/o NCR | 36.4 |
| tf-base100k-plusbt w/ NCR | **37.6** |
| tf-large100k w/o NCR | 18.3 |
| tf-large100k w/ NCR | 21.3 |
| tf-large100k-plusbt w/o NCR | 35.4 |
| tf-large100k-plusbt w/ NCR | **37.1** |

Table 4: Ablation results for spa→arg. NCR denotes noisy channel reranking.

| Setting | BLEU | | |
|---------|------|-----|------|
| | whole | mid | long |
| no-BT (baseline) | 8.2 | 8.4 | 8.2 |
| short-BT | 10.5 | 9.6 | 10.4 |
| mid-BT | 11.6 | 10.7 | 11.6 |
| long-BT | **14.3** | **11.4** | **14.3** |

Table 5: BLEU scores per length group of BT data. long-BT outperforms all other settings in all test setups.

length penalty. The length penalty sweep was performed for 137 iterations sampling from a uniform distribution with minimum 0.0 and maximum 2.0. Hyperparameter sweeping was performed using the **tf-large60k-plusbt** direct translation models with noisy channel reranking enabled on the Spanish to Aragonese language pair. Translations for the hyperparameter sweep were generated from the copy of FLORES+ (Team et al., 2022) found in the PILAR (Galiano-Jiménez et al., 2024) repository[7]. The results of this sweep were used on all language pairs. We performed the sweep on spa→arg only and on a **tf-large60k-plusbt** model due to hardware and time constraints. Our sweeps showed that setting length penalty to 1.726 is optimal.

# 3 Results and Discussion

In this section, we discuss the results of our experiments and discuss our findings. Experiments were performed using the copy of FLORES+ (Team et al., 2022) found in the PILAR (Galiano-Jiménez et al., 2024) repository were computed using Sacre-BLEU[8] (Post, 2018).

For all translations, we used the following decoding hyperparameters: top_k=50, top_p= 0.93, temperature=0.7, beam=5. Additional hyperparameters are specified per experiment.

## 3.1 Comparison Against Baselines

We compare our system against Apertium 3.9.6. Results are listed in Table 3. We observe that Apertium yields the highest BLEU score for spa→arg and spa→arn. For spa→ast, the systems trained with BT data both outperform the Apertium baseline.

Our method performs worst on the spa→arn language pair while it performs best for spa→arg. However for both of these pairs our system is outperformed by Apertium. From this we can conclude that our current pipeline cannot overcome the low resource nature of these language pairs in order to close the gap with Apertium. For spa→ast, we were able to outperform Apertium with a difference of 1.8 BLEU.

## 3.2 Ablations

We perform an ablation study by varying model size, use of BT data, and use of noisy channel reranking. Due to hardware and time constraints, we only perform our ablations in the spa→arg direction. Results are summarized in Table 4.

We observe that the addition of BT data and

---

[7]https://github.com/transducens/PILAR

[8]SacreBLEU signature: nrefs:1|case:mixed|eff:no|tok:flores101|smooth:exp|version:2.4.2

Figure 1: Sequence length distribution of the target side of the train (filtered) and test set per language pair. spa-arn has the least overlap between train and test and has the most short examples in training data which hints why the BLEU score is relatively lower compared to other language pairs.

noisy channel reranking resulted in an increased BLEU score. Using both strategies yields the highest BLEU score for both base and large model sizes. It is notable that the base model with both BT data and noisy channel reranking yields the highest BLEU score in our ablation study. We speculate that this be due to the large model having too many parameters for the given task or a lack of data. Another reason is that spa-arn is relatively easier compared to other language pairs because the source and target are more similar with each other as shown in Table 6. More experiments are needed to confirm these.

### 3.3 Adding Longer Examples Improves BLEU Better than Shorter Examples

For the spa→arn baseline model (No BT), we observed a BLEU score of 8.7 on FLORES+ dev set. One possible explanation for the low score is the mismatch between the length distribution of training and test data. We observed that the training data is comprised mostly of short examples while the FLORES+ dev set is relatively longer (see Figure 1). We hypothesize that adding longer examples to the training set will improve BLEU score, especially on longer examples.

To provide evidence for the hypothesis, we generated BT data of size 100,000 for different length groups namely, short-BT (1-10 words), mid-BT (11-20 words), long-BT (20+ words). We mixed the BT data with the training data then trained a model for each setup. We trained each model for 50,000 of which 5,000 are warmup steps. We used the same training hyperparameters as in Table 2. For fair comparison, we trained a baseline model (no-BT) using the same training hyperparameters. For generating the translations, we did not use noisy channel reranking and we fix the

length penalty to 1.0. The results are summarized in Table 5.

The result shows that long-BT gives an absolute BLEU score improvement of +6.1 over baseline, followed by mid-BT (+3.4), and then short-BT (+2.3). This tells us that while augmenting with BT data generally improves the performance, strategically adding more long examples can give the most improvements in a resource-constrained setting. To strengthen this claim further, we performed a fine-grained test by grouping FLORES+ dev set by length groups (mid/long). For this experiment, we did not include the short length group because it only contains 3 examples after grouping. The results shows that long-BT gives the most improvements on mid and long test groups, followed by mid-BT and short-BT (see Table 5). This suggests that training on longer sequences also improves performance on shorter sequences.

While this experiment shows empirical results that adding longer examples improves the overall BLEU score better than adding shorter examples, it does not say something about the quality and diversity of the text. It is possible that these findings might not hold if the long examples are of low quality. Another possible explanation on why long-BT outperforms its shorter counterparts is because, with the same number of examples of 100,000, long-BT contains more tokens than short-BT and mid-BT. To further solidify the claim that adding longer examples improves the overall BLEU score better than adding shorter examples, more experiments are needed where total token count per length group are equal or close to each other.

| Language Pair | JS (generated) ↑ | ED (generated) ↓ | JS (ground truth) ↑ | ED (ground truth) ↓ |
|---|---|---|---|---|
| spa→arg | 0.34 | 0.59 | **0.34** | **0.57** |
| spa→arn | 0.23 | 0.69 | 0.13 | 0.83 |
| spa→ast | **0.44** | **0.44** | 0.23 | 0.76 |

Table 6: Average Jaccard similarity (JS) and average normalized edit distance (ED) between source and generated translations and ground truth translations. Results confirm our observatoin that our system is regurgitating Spanish source sentences in the spa→ast direction. Results also suggest that the Spanish and Aragonese sentences in the FLORES+ dev set are more similar to each other compared to others.

### 3.4 Regurgitation of Spanish Sentences in Generated Translations

We observed that our model was producing some translations that were only slightly altered versions of the source Spanish sentence. To empirically evaluate the extent of this problem for our system, we compare the BPE tokenized source Spanish sentences of the FLORES+ dataset from PILAR to the corresponding generated translations made by our system and the corresponding ground truth. We compared this system via two metrics: Jaccard similarity (JS) and normalized edit distance (ED). To compute the two metrics between two BPE encoded sentences $S_1$ and $S_2$, we get the set of tokens of each sentence $T_1$ and $T_2$ and compute Jaccard similarity as

$$JS = \frac{|T_1 \cap T_2|}{|T_1 \cup T_2|}$$

and normalized edit distance as

$$ED = \frac{D(S_1, S_2)}{\max(|S_1|, |S_2|)}$$

where $D(S_1, S_2)$ denotes token-level Levenshtein distance between BPE encoded sentences $S_1$ and $S_2$. We divide by the maximum length between $S_1$ and $S_2$ to ensure that we get a value between 0 and 1. Results of this analysis are summarized in Table 6.

We observe that spa→ast shows the highest average Jaccard similarity and the lowest normalized edit distance among language pairs for generated translations; however, the corresponding metrics for the spa→ast ground truth translations tell a different story. Ground truth translations for spa→ast show a lower Jaccard similarity and a higher normalized edit distance, indicating that we may be regurgitating Spanish sentences.

Below is a sample of a Spanish sentence together with a generated Asturian translation which exhibits regurgitation and the corresponding ground truth translation. Notice how the generated translation is closer in similarity to the Spanish sentence than the correct Asturian translation. In the below example, "S -" is the source spanish sentence, "H -" is the generated Asturian translation, "T -" is the ground truth Asturian translation, "J -" is the jaccard similarity compared to the source Spanish sentence, and "E -" is the normalized edit distance compared to the source Spanish sentence. All sentences are BPE encoded.

```
S - _Apenas _pas adas _las _11: 00 _h , _los
    _integrantes _de _la _manifestación _bloque
    aron _la _circulación _del _car ril _de
    _White h all _que _va _hacia _el _norte .

H - _Ap enes _pasa es _les _11: 00 _h , _los
    _integrantes _de _la _manifestación _blo qui
    aron _la _circulación _del _car ril _de
    _White h all _que _va _escontra ' l _norte .
J - 0.553
E - 0.303

T - _X usto _depués _de _les _11: 00, _los
    _manifestantes _blo qui aron _el _trá ficu
    _nel _sentíu _norte _en _White h all .
J - 0.244
E - 0.833
```

For spa→arg, Jaccard similarity and normalized edit distance are similar for both generated translations and ground truth translations. We note that this language pair has the highest Jaccard similarity and lowest normalized edit distance between its source Spanish sentences and ground truth Aragonese translations. This indicates that there is a degree of similarity between the Spanish and Aragonese sentences in the dataset which may explain why the spa→arg model exhibited the highest BLEU score in our baseline comparison. We provide a sample below where the source Spanish sentence is similar to the ground truth Aragonese translation.

```
S - _En _el _partido , _Nadal _acumul ó _un _8
    8% _de _puntos _ne tos _y _ganó _76 _en _el
    _primer _servicio .

H - _En _o _parti to , _Nadal _acumul ó _un _8
    8% _de _puntos _ne tos _y _ganó _76 _en _o
```

```
    _primer _servicio .
J - 0.792
E - 0.174

T - _En _o _parti u , _Nadal _acumul ó _un  _8
    8% _de _puntos _ne tos _y _ganó _76 _en _o
    _primer _servicio .
J - 0.792
E - 0.174
```

We plot the histogram of Jaccard similarity and normalized edit distance for all language pairs in Figures 2, 3, and 4.

### 3.5 Character Set Analysis

We observe that our generated translations do not contain all characters present in the ground truth as shown in Table 7. For all languages, the missing characters are present in the training data with the exception of Õ for Asturian and Aragonese. All missing characters constitute less than 1% of the training data which may explain why they were not learned by our models.

## 4 Conclusions

We detailed our constrained system for translating from Spanish to Aragonese (spa→arg), Aranese/Occitan (spa→arn), and Asturian (spa→ast). These systems were trained from scratch on constrained data, augmented by backtranslated (BT) data. Translations were further improved by utilizing Noisy Channel Reranking. This approach outperformed Apertium on the spa→ast translation direction. Our ablation study for spa→arg showed that utilizing backtranslation and noisy channel reranking improves BLEU score. However, more experiment is needed for other language pairs. Our ablation experiment also suggests that smaller models are capable enough for spa→arg, at least for this train and test set.

We investigated the cause of low BLEU score for spa→arn despite having more data (after filtering) than spa→arg and spa→ast. We linked it to the train-test mismatch of spa-ast data in terms of sequence length. We also found that adding longer backtranslated data improves overall BLEU score even in shorter sequences.

Lastly, we observed that our model for spa→ast was regurgitating Spanish sentences in Asturian translations and that characters with low frequencies in the training data are not being learned by our models.

## Limitations

We are unable to evaluate whether the translations we generate are syntactically or semantically sound due to the fact that none of us speak Spanish, Aragonese, Asturian, or Aranese/Occitan.

## References

Jan Christian Blaise Cruz. 2023. Samsung R&D institute Philippines at WMT 2023. In *Proceedings of the Eighth Conference on Machine Translation*, pages 103–109, Singapore. Association for Computational Linguistics.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Pilar.

Tanmai Khanna, Jonathan N. Washington, Francis M. Tyers, Sevilay Bayatlı, Daniel G. Swanson, Tommi A. Pirinen, Irene Tang, and Hèctor Alòs i Font. 2021. Recent advances in apertium, a free/open-source rule-based machine translation platform for low-resource languages. *Machine Translation*, 35(4):475–502.

Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016b. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

| Language | Missing Characters | Frequency in Training Data | # Characters in Training Data |
|---|---|---|---|
| Aragonese | » « ] & õ [ Õ | 1,583 | 4,713,907 |
| Aranese/Occitan | & « » Ç Ò Õ *U+0301* ' | 96,233 | 36,014,337 |
| Asturian | " Õ Ú *U+1E24* ḥ – — ' | 56,798 | 42,897,857 |

Table 7: Characters present in ground truth translations but missing in generated translations together with their frequency in training data compared to the total number of characters in training data. Unicode symbol code in italics listed when a character is unsupported by LᴬTEX. All missing characters constitute less that 1% of the training data.



(a)                                    (b)

Figure 2: Distribution of Jaccard similarity and normalized edit distance for spa→arg of source sentences vs generated translations and ground truth translations. We can see that the distributions for both Jaccard similarity and normalized edit distance almost entirely overlap. Taken together with the means from Table 6, these show that any regurgitation our model exhibits can also be seen in the ground truth test data.



(a)                                    (b)

Figure 3: Distribution of Jaccard similarity and normalized edit distance for spa→arn of source sentences vs generated translations and ground truth translations. We can see in (a) that while Jaccard similarity of generated translations vs. source Spanish sentences is higher compared to that of ground truth translations vs. source Spanish sentences, they both tend to be less than 0.4. In (b), we see that while normalized of generated translations vs. source Spanish sentences is lower compared to that of ground truth translations vs. source Spanish sentences, they both tend to be greater than 0.6. This indicate low amounts of regurgitation in the case of our spa→arn system.

|  (a)  |  (b)  |

Figure 4: Distribution of Jaccard similarity and normalized edit distance for spa→ast of source sentences vs generated translations and ground truth translations. We see in (a) that the Jaccard similarity of generated Asturian translations compared to source Spanish sentences is higher than that of ground truth translations compared to source sentences. In (b), we see that the normalized edit distance of generated translations compared to source sentences is lower than that of ground truth vs. source sentences. This indicates that our model is regurgitating more Spanish words rather than translating to Asturian.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Jörg Tiedemann. 2016. OPUS – parallel corpora for everyone. In *Proceedings of the 19th Annual Conference of the European Association for Machine Translation: Projects/Products*, Riga, Latvia. Baltic Journal of Modern Computing.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Kyra Yee, Yann Dauphin, and Michael Auli. 2019. Simple and effective noisy channel modeling for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5696–5701, Hong Kong, China. Association for Computational Linguistics.

# Back to the Stats: Rescuing Low Resource Neural Machine Translation with Statistical Methods

**Velayuthan Menan[1], Dilith Jayakody[1], Nisansa de Silva[1],**
**Aloka Fernando[1], Surangika Ranathunga[2]**
[1]Dept. of Computer Science & Engineering, University of Moratuwa,10400, Sri Lanka
[2]Massey University, Palmerston North, 4443, New Zealand
{velayuthan.22,dilith.18,NisansaDdS,alokaf}@cse.mrt.ac.lk
s.ranathunga@massey.ac.nz

## Abstract

This paper describes our submission to the *WMT24 shared task for Low-Resource Languages of Spain* in the Constrained task category. Due to the lack of deep learning-based data filtration methods for these languages, we propose a purely statistical-based, two-stage pipeline for data filtration. In the primary stage, we begin by removing spaces and punctuation from the source sentences (Spanish) and deduplicating them. We then filter out sentence pairs with inconsistent language predictions by the language identification model, followed by the removal of pairs with anomalous sentence length and word count ratios, using the development set statistics as the threshold. In the secondary stage, for corpora of significant size, we employ a Jensen-Shannon divergence-based method to curate training data of the desired size. Our filtered data allowed us to complete a two-step training process in under 3 hours, with GPU power consumption kept below 1 kWh, making our system both economical and eco-friendly. The source code, training data, and best models are available on the project's GitHub page[1].

## 1 Introduction

We[2] participated in the **Constrained submission** category of the WMT24 shared task for Low-Resource Languages of Spain (Sánchez-Martínez et al., 2024), focusing on the **Spanish-Asturian** language pair. For the **Constrained submission** category, we are limited to using only the resources provided on the official shared task site[3], and all models utilized must not exceed 1 billion parameters.

Previous shared tasks on data filtering have used deep learning-based scoring methods like LASER (Heffernan et al., 2022) and LaBSE (Feng et al., 2022), as well as sentence alignment methods such as SentAlign (Steingrimsson, 2023) and Vecalign (Thompson and Koehn, 2019). However, these methods often fail with low-resource languages (LRL) due to a lack language of support.

Following prior work (Cruz and Sutawika, 2022; Vegi et al., 2022; Zhang, 2023), we focus on statistical data filtration and sampling techniques to curate our datasets, ensuring our method is not limited to specific languages. Given our compute resource constraints, we design our pipeline to utilize small dataset sizes, enabling a larger volume of experiments. As noted in Ranathunga et al. (2024), randomly sampling a large corpus and training on that sample yields sub-optimal results. Therefore, we use the Jensen-Shannon Divergence (JSD) (Lu et al., 2020) to filter subsets from large corpora (see Section 3.1.1). We favor JSD over Kullback-Leibler (KL) divergence and higher-order domain discriminators due to its symmetric property and relatively simple implementation.

In addition to data filtration, we experiment with a two-step training schema. First, we train the entire model on a larger filtered dataset. Then, in the second step, we freeze the encoder layers and fine-tune the model on the filtered dataset for fine-tuning. This approach proved effective for the task. We select models with under 1 billion parameters for all experiments to adhere to the rules of the *Constrained task*.

Our key contributions are:

- We propose a two-stage data filtration system that can be applied to any language. This system includes statistical data filtration methods for bilingual and monolingual data, along with a Jensen-Shannon divergence-based filtration method.

---

- We achieve competitive results in a compute resource-constrained environment (Table 4).

- Based on our experiments, we show that fine-tuning a multilingual translation model for a high-resource source language and a low-resource target language is most effective when high-quality monolingual target data is leveraged, and the encoder is frozen to preserve the source language knowledge while training the model.

- We maintain an effective training time well under 3 hours and keep the total GPU power consumption of our best model (training + fine-tuning) under 1 kWh, resulting in a minimal carbon footprint and enhancing the eco-friendliness of our training schema.

## 2 Related Work

Bilingual parallel data curated from web-mined corpora are prone to various types of noise. Kreutzer et al. (2022) investigated the issue of noise in web-mined corpora by analyzing a sample of 100 sentence pairs, providing evidence of the problem. In a similar vein, Khayrallah and Koehn (2018) examined how different types of noise in parallel training data impact the quality of neural machine translation systems. Building on this body of work, Ranathunga et al. (2024) further established that data quality is more critical than data quantity, particularly for low-resource languages. These studies collectively highlight the importance of addressing *bad* data in web-mined corpora.

Handling noisy parallel data in machine translation had been extensively studied, with various methods proposed in the literature. Below, we present these methods, grouping them into two main categories: (1) deep learning-based approaches and (2) statistical-based approaches. Additionally, we include methods that do not fit broadly into these categories under the section *Other Approaches*.

**Deep Learning-Based Approaches.** Recent advances in deep learning have introduced several methods to handle unclean data. Zhang (2023) proposed a denoising approach by pretraining on corrupted data and regenerating the original content. They utilized text corruption techniques as proposed by Lewis et al. (2020), including token masking, sentence permutation, document rotation, token deletion, and text infilling. They pretrained

models on synthetic and monolingual data and fine-tuned them on clean parallel corpora, achieving translation perplexity scores by training two models and analyzing prediction difficulties. Ensembling methods were also employed to enhance performance. Chaudhary et al. (2019) utilized LASER and an ensemble of scoring methods to check the similarity between embeddings and cross-entropy scores for both directions, penalizing significant differences. Abdulmumin et al. (2022) developed a binary classifier to predict translation accuracy, collecting positive data from a gold standard dataset and negative data from the worst LASER alignment scores. Unfortunately, majority of the methods and models mentioned above do not support Asturian, Aragonese, or Aranese.

**Statistical Approaches.** Steingrímsson et al. (2023) applied rules such as filtering sentences with three tokens or less, ensuring 60% or more token overlap between languages, and requiring at least 70% alphabetical characters in both sentences. Vegi et al. (2022) introduced constraints such as filtering sentences where the source or target exceeds 800 characters, where the length ratio is greater than 2.5 or less than 0.4, or where words exceed 10 characters. Cruz and Sutawika (2022) extended these rules to include sentences with too many contiguous punctuations (three or more), a large percentage of numbers or punctuations, and additional filtering criteria. Minh-Cong et al. (2023) proposed building a dictionary using MGiza++ and clean parallel corpora, translating source sentences, calculating edit distances, and iteratively training NMT models, assuming the availability of clean parallel corpora.

**Other Approaches.** Other approaches for handling unclean data include additional filtering rules proposed by Vegi et al. (2022), such as removing sentences that are empty or identical between source and target. Cruz and Sutawika (2022) suggested removing sentences with missing punctuations in one language, sentences containing HTML or URLs, ensuring numbers appear in both source and target, and deduplicating data after preprocessing. Steingrímsson et al. (2023) recommended using a language filter to ensure both languages are in the top two predictions and removing near-duplicate pairs.

## 3 Methodology

**Training Datasets:** For training, we use bilingual datasets from OPUS[4] CCMatrix and WikiMedia (Spanish-Asturian), subjected to the filtration outlined in section 3.1.1. We also use monolingual datasets, specifically the PILAR (Galiano-Jiménez et al., 2024b) Asturian monolingual dataset and the Spanish side of the English-Spanish Wikimedia dataset from OPUS, with filtration procedures detailed in section 3.1.2.

**Development Set:** We evaluate the trained models using the FLORES+[5] Spanish-Asturian development set, referred to as the development set throughout the paper.

**Hardware Specifications:** All experiments were conducted on a single machine with an Intel i9-9900K CPU, 64GB of RAM, and an Nvidia Quadro RTX 6000 (24GB VRAM).

**Software Specifications:** All models and training code were developed using the HuggingFace (HF) Transformers (Wolf et al., 2020) library. For evaluation, we use chrF and BLEU scores from the `evaluate`[6] library of HF. We utilized the work done by Nayak et al. (2023) to obtain the Jensen-Shannon divergence scores.

**Models:** We use NLLB-200-600M (NLLB Team et al., 2022) (we will address it as NLLB-600M throughout the paper), M2M100-418M (Fan et al., 2020) (we will address it as M2M100 throughout the paper), and SMaLL-100 (Mohammadshahi et al., 2022) to conduct experiments. All training and fine-tuning was performed on SMaLL-100 (see section 3.2)

**Training Details:** We use the HF Transformers Trainer API with the AdamW optimizer, a learning rate of $1 \times 10^{-5}$, and a batch size of 16. Gradient accumulation steps of 16 were used to increase the effective batch size to 256. Training is conducted in two steps: first, training the entire model using the filtered Spanish-Asturian CCMatrix dataset (Table 3); then, fine-tuning the best model from the training phase by freezing the encoder layers (Table 4).

**Dataset used in fine-tuning step:** For fine-tuning, we combined multiple datasets as follows

(see Table 4 for results): Dataset $A$ is the filtered Spanish-Asturian Wikimedia; Dataset $B$ includes $A$ + the filtered PILAR crawled monolingual data; Dataset $C$ includes $B$ + the filtered PILAR literary monolingual data; Dataset $D$ includes $C$ + the Spanish monolingual data from English-Spanish Wikimedia. All monolingual data were translated using NLLB-600M.

### 3.1 Data Filtration and Curation

#### 3.1.1 Bilingual Data Filtration

Our Bilingual Data Filtration pipeline consists of two stages: Primary and Secondary.

**Primary Filtration:** We start by removing punctuation and whitespace from the Spanish text, then deduplicate it. Using Idiomata Cognitor (Galiano-Jiménez et al., 2024a), we classify the language of each sentence, removing those with inconsistent predictions. We analyze sentence length ratios and word counts, using the development set as a benchmark to remove anomalies.

| Dataset | Before (M) | After (M) | % drop |
|---|---|---|---|
| Wikimedia - es_ast | 0.04 | 0.03 | 38.99 |
| CCMatrix - es_ast | 5.39 | 2.01 | 62.72 |
| Wikipedia - es_en | 2.80 | 1.26 | 55.19 |
| Wikimedia - es_en | 1.81 | 1.31 | 27.47 |

Table 1: The table presents the sample count (in millions) before and after primary filtration, along with the percentage of samples dropped during this phase.

**Secondary Filtration:** For datasets over 300K sentences, we limit the size to 50K-300K sentences. We use the Jensen-Shannon divergence to refine the data. We sample sets of 2000 sentence pairs (with replacement) to form 1000 sets, remove Spanish stop words and punctuation with the *NLTK library*[7], and calculate word frequency distributions for each sample and the development set. We sort the samples by their divergence scores (Nayak et al., 2023) against the development set, and iteratively merge and deduplicate low-divergence samples until we have the target dataset size (as specified in Table 3). The implementation of the secondary filtration method is detailed in Algorithm 1.

**Algorithm 1:** Secondary Filtration

**Input:** Dataset $D$ with $|D| \geq 500K$, development set $E$, number of batches $N$, batch size $L$, desired size $S$

**Output:** Deduplicated batch set $B_{final}$ of size $S$

```
// Initialize empty batch list
```
$B \leftarrow \{\}$
**for** $i \leftarrow 1$ **to** $N$ **do**
 ```
// Randomly sample L rows with
//    replacement
```
 $B_i \leftarrow RandomSample(D, L)$
 $B \leftarrow B \cup \{B_i\}$
**end**
```
// Initialize empty scores list
```
$scores \leftarrow \{\}$
**for** $each\ batch\ B_i \in B$ **do**
 ```
// Jensen-Shannon Divergence
//    between B_i and E
```
 $JS_i \leftarrow JS\_div(B_i, E)$
 $scores \leftarrow scores \cup \{(B_i, JS_i)\}$
**end**
Sort($scores$) by $JS_i$ in ascending order
```
// Initialize final batch List
```
$B_{final} \leftarrow \{\}$
$current\_size \leftarrow 0$
**while** $current\_size < S$ **do**
 $B_{candidate} \leftarrow scores[i].batch$
 $B_{final} \leftarrow B_{final} \cup \{B_{candidate}\}$
 De-duplicate($B_{final}$)
 $current\_size \leftarrow Length(B_{final})$
**end**
**return** $B_{final}$

### 3.1.2 Monolingual Data Filtration

We extract monolingual data from PILAR crawled and literacy datasets for Asturian, and Wikimedia OPUS datasets for Spanish using the English-Spanish direction. Using *sentence-splitter*[8], we segment the PILAR text into sentences using the Spanish setting, achieving good performance despite the library's lack of support for Asturian. We removed URL links and retained sentences with a word count between four and sixty (the maximum in the development set). Sentences were then classified using the language identifier model Idiomata Cognitor (Galiano-Jiménez et al., 2024a), and those

not identified as Spanish or Asturian were removed.

### 3.2 Model Selection

**Translation Model Selection:** Based on zero-shot performance scores, NLLB-600M was selected as the best model for translating the filtered monolingual dataset, outperforming both M2M-100 and SMaLL-100 in chrF and BLEU scores (Table 2).

**Training Model Selection:** Given the limited GPU resources in our training environment, we selected SMaLL-100 as the model for training and experimentation due to its smaller size and superior performance compared to M2M-100 (Table 2).

## 4 Results and Discussion

In this section, we present the results of our experiments using the proposed methods. The results from our primary data filtration step (Section 3.1.1) demonstrate the row counts of each dataset before and after filtration. Notably, our filtration method had the most significant impact on the CCMatrix (es-ast) and Wikipedia (es-en) datasets, with data reduction percentages exceeding 50%. Further investigation could be conducted to understand the factors contributing to this substantial data drop in these sources and to determine whether this observation is consistent across other language pairs in these datasets.

Table 2 displays the zero-shot performance of three state-of-the-art open-source models: NLLB-600M, M2M100, and SMaLL-100. NLLB-600M was selected as the model for generating translations for monolingual sentences due to its significantly better performance compared to the other two models. Given its smaller size and superior performance compared to M2M100, SMaLL-100 was chosen as the model for training and experimentation. By choosing the SMaLL-100 model, we gained the added advantage of using larger batch sizes due to the model's small size. This proved crucial in our low-compute resource environment.

Table 3 presents the results from the first step of the two-step training regime described in Section 3, applied to the CCMatrix filtered dataset. The data was incrementally filtered based on increasing Jensen-Shannon divergence scores in steps of 50k. We observe that model performance improves up to a subset size of 100K, after which it gradually declines. This observation aligns with the findings of Ranathunga et al. (2024), emphasizing that data

| Model | chrF | BLEU | # of Trainable Params (M) |
|---|---|---|---|
| NLLB-600M | **49.77** | **17.16** | 615.07 |
| M2M100 | 46.27 | 14.71 | 483.91 |
| SMaLL-100 | 48.47 | 14.85 | **332.74** |

Table 2: Scores for the zero-shot performance of the models evaluated on FLORES+ Spanish-Asturian dev set and number of trainable parameters for each model.

quality is more important than the sheer size of the dataset. Selecting smaller, higher-quality datasets not only enhances performance but also offers the additional benefits of reduced training time and lower computational requirements, which in turn minimizes the carbon footprint. As shown in Table 3, the best-performing subset required only 1.48 hours of training and consumed just 0.44 kWh of GPU power.

| Size (k) | chrF | BLEU | Time (hrs) | Power (kWh) |
|---|---|---|---|---|
| 50 | 49.63 | 17.25 | 0.72 | 0.21 |
| 100 | **50.04** | **17.52** | 1.48 | 0.44 |
| 150 | 49.84 | 17.26 | 2.25 | 0.66 |
| 200 | 49.95 | 17.34 | 2.92 | 0.86 |
| 250 | 49.81 | 17.36 | 3.68 | 1.09 |
| 300 | 49.70 | 17.02 | 4.42 | 1.30 |

Table 3: Scores, training durations, and GPU power consumption for training CC-Matrix Spanish-Asturian filtered datasets at intervals of 50K jumps. This is the first step in the training phase.

The results of our second step of model training (fine-tuning), where the encoder layers were frozen, are presented in Table 4. Among the various combinations, Dataset $C$ demonstrated the best performance. This dataset includes Spanish-Asturian Wikimedia data as well as Spanish-Asturian PILAR crawled and literary datasets. Notably, the Spanish side of this dataset was generated using the translation model (NLLB-600M) for the PILAR Asturian monolingual crawled and literary datasets.

Interestingly, the performance drops when using Dataset $D$, which consists of Dataset $C$ combined with Spanish data (Spanish-English Wikimedia) translated into Asturian using the translation model. This observation underscores the importance of high-quality target-side sentences, as the monolingual PILAR dataset comprises carefully curated, high-quality Asturian data. We hypothe-

| | Size (k) | chrF | BLEU | Time (hrs) | Power (kWh) |
|---|---|---|---|---|---|
| A | 24.8 | 51.14 | 17.80 | 0.28 | 0.08 |
| B | 38.5 | 51.38 | 18.02 | 0.60 | 0.18 |
| C | 60.2 | **51.47** | **18.17** | 0.95 | 0.28 |
| D | 84.9 | 50.92 | 17.97 | 1.35 | 0.40 |

Table 4: Dataset name, Size, Scores, Time duration and the GPU Power consumption of fine-tuning the best model from Table 3. A = filtered Spanish-Asturian Wikimedia; B = A + filtered and translated PILAR crawled data; C = B + filtered and translated PILAR literary data; D = C + Spanish-Asturian data from English-Spanish Wikimedia.

size that since Spanish is a high-resource language, the model's encoder has likely been exposed to extensive Spanish data. By freezing the encoder and allowing the model to learn during this fine-tuning step, the model was better able to focus on the target language. Based on these observations, we conclude that when fine-tuning a pre-trained multilingual translation model for a high-resource source language and a low-resource target language, it is essential to leverage high-quality monolingual data for the target language and freeze the encoder to retain the learned knowledge of the source language while making the other layers trainable.

The training time and GPU power consumption for the best datasets from our two-step training procedure (as shown in Table 3 and Table 4) remains well within 3 hours and consumes less than 1 kWh of GPU power. This makes our proposed method highly suitable for low-compute environments.

## 5 Conclusion

We presented a purely statistical-based pipeline for data filtering, demonstrating that simple statistical methods should not be overlooked, particularly for low-resource languages where deep learning-based methods may fail to provide adequate support. Our proposed pipeline achieved competitive performance in a low-compute environment for the constrained task, proving to be both economical, with training times well under 3 hours as well as eco-friendly, with GPU power consumption kept under 1 kWh. This work reinforces the findings of previous studies that emphasize the importance of data quality over quantity. We hope that our methodology will encourage and empower researchers in low-compute environments to contribute to an egalitarian representation of lan-

guages.

## Acknowledgements

## References

Idris Abdulmumin, Michael Beukman, Jesujoba Alabi, Chris Chinenye Emezue, Everlyn Chimoto, Tosin Adewumi, Shamsuddeen Muhammad, Mofetoluwa Adeyemi, Oreen Yousuf, Sahib Singh, and Tajuddeen Gwadabe. 2022. Separating grains from the chaff: Using data filtering to improve multilingual translation for low-resourced African languages. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1001–1014, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Vishrav Chaudhary, Yuqing Tang, Francisco Guzmán, Holger Schwenk, and Philipp Koehn. 2019. Low-resource corpus filtering using multilingual sentence embeddings. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 261–266, Florence, Italy. Association for Computational Linguistics.

Jan Christian Blaise Cruz and Lintang Sutawika. 2022. Samsung research Philippines - datasaur AI's submission for the WMT22 large scale multilingual translation task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1034–1038, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT Sentence Embedding.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024a. Idiomata cognitor.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024b. Pan-iberian language archival resource.

Kevin Heffernan, Onur Çelebi, and Holger Schwenk. 2022. Bitext mining using distilled sentence representations for low-resource languages. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2101–2112, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Huda Khayrallah and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83, Melbourne, Australia. Association for Computational Linguistics.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Jinghui Lu, Maeve Henchion, and Brian Mac Namee. 2020. Diverging divergences: Examining variants of Jensen Shannon divergence for corpus comparison tasks. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6740–6744, Marseille, France. European Language Resources Association.

Nguyen-Hoang Minh-Cong, Nguyen Van Vinh, and Nguyen Le-Minh. 2023. A fast method to filter noisy parallel data WMT2023 shared task on parallel data curation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 359–365, Singapore. Association for Computational Linguistics.

Alireza Mohammadshahi, Vassilina Nikoulina, Alexandre Berard, Caroline Brun, James Henderson,

and Laurent Besacier. 2022. Small-100: Introducing shallow multilingual machine translation model for low-resource languages. *arXiv preprint arXiv:2210.11621*.

Shravan Nayak, Surangika Ranathunga, Sarubi Thillainathan, Rikki Hung, Anthony Rinaldi, Yining Wang, Jonah Mackey, Andrew Ho, and En-Shiun Annie Lee. 2023. Leveraging auxiliary domain parallel data in intermediate task fine-tuning for low-resource translation.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Surangika Ranathunga, Nisansa de Silva, Velayuthan Menan, Aloka Fernando, and Charitha Rathnayake. 2024. Quality does matter: A detailed look at the quality and utility of web-mined parallel corpora. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 860–880, St. Julian's, Malta. Association for Computational Linguistics.

Steinþór Steingrímsson, Hrafn Loftsson, and Andy Way. 2023. Filtering matters: Experiments in filtering training sets for machine translation. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 588–600, Tórshavn, Faroe Islands. University of Tartu Library.

Steinthor Steingrimsson. 2023. A sentence alignment approach to document alignment and multi-faceted filtering for curating parallel sentence pairs from web-crawled data. In *Proceedings of the Eighth Conference on Machine Translation*, pages 366–374, Singapore. Association for Computational Linguistics.

Felipe Sánchez-Martínez, Juan Antonio Pérez-Ortiz, Aarón Galiano-Jiménez, and Antoni Oliver. 2024. Findings of the WMT 2024 shared task on translating into low-resource languages of spain: Blending rule-based and neural systems. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Brian Thompson and Philipp Koehn. 2019. Vecalign: Improved sentence alignment in linear time and space. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1342–1348, Hong Kong, China. Association for Computational Linguistics.

Pavanpankaj Vegi, Sivabhavani J, Biswajit Paul, Prasanna K R, and Chitra Viswanathan. 2022. ANVITA-African: A multilingual neural machine translation system for African languages. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1090–1097, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Wenbo Zhang. 2023. IOL research machine translation systems for WMT23 low-resource Indic language translation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 978–982, Singapore. Association for Computational Linguistics.

# Hybrid Distillation from RBMT and NMT: Helsinki-NLP's Submission to the Shared Task on Translation into Low-Resource Languages of Spain

**Ona de Gibert[1]** and **Mikko Aulamo[1]** and **Yves Scherrer[1,2]** and **Jörg Tiedemann[1]**

[1]University of Helsinki, Dept. of Digital Humanities
[2]University of Oslo, Dept. of Informatics
[1]`firstname.lastname@helsinki.fi`
[2]`firstname.lastname@ifi.uio.no`

## Abstract

The Helsinki-NLP team participated in the 2024 Shared Task on Translation into Low-Resource languages of Spain with four multilingual systems covering all language pairs in the open submission track. The task consists in developing Machine Translation (MT) models to translate from Spanish into Aragonese, Aranese and Asturian. Our models leverage known approaches for multilingual MT; namely, data filtering, fine-tuning, data tagging, and distillation. We use distillation to merge the knowledge from neural and rule-based systems and explore the trade-offs between translation quality and computational efficiency. We demonstrate that our distilled models can achieve competitive results while significantly reducing computational costs. Our best models ranked 4th, 5th, and 2nd in the open submission track for Spanish–Aragonese, Spanish–Aranese, and Spanish–Asturian, respectively. We release our code and data publicly at `https://github.com/Helsinki-NLP/lowres-spain-st`.

## 1 Introduction

In this work, we describe the participation of our team to the Shared Task on Translation into Low-Resource Languages of Spain 2024 (Sánchez-Martínez et al., 2024), the first edition of its kind. The task involves developing Machine Translation (MT) systems for translating from Spanish (*spa*) into three closely related Romance target languages: Aranese (*arn*), Aragonese (*arg*) and Asturian (*ast*). Aranese is a variety of Occitan spoken in the northwestern part of Catalonia; Aragonese is spoken in Aragon, in northwest Spain; and Asturian is spoken in Asturias, in northeast Spain.

Although these minority languages have some form of official status in their respective regions, they are all considered endangered. According to the linguistic taxonomy proposed by Joshi et al. (2020), these languages fall into the category of the "Scraping Bys". This means that, while there is

some available unlabeled data, substantial and coordinated efforts are necessary to raise awareness and gather labeled datasets to improve the prospects of these languages in the future. This task is designed precisely to address these challenges by fostering the development of resources and tools for these under-resourced languages.

In terms of current technological support, some linguistic resources are available for these languages, including online dictionaries and established orthographic standards. Apertium (Forcada et al., 2011) is an open-source Rule-Based MT (RBMT) toolkit initially developed for related languages, that offers substantial coverage for the three target languages. Nevertheless, resources remain notably sparse for data-driven approaches like Neural Machine Translation (NMT). By contributing to this task, we aim to change this picture.

We focus our participation efforts on data collection – by gathering additional data from Wikipedia and online dictionaries –, data augmentation – by producing back-translations (Sennrich et al., 2016) of monolingual data –, and data preparation – by carrying out corpus-targeted cleaning. We also experiment with different data tagging strategies. We submit four multilingual models, which arise from fine-tuning and applying Knowledge Distillation (KD), by leveraging both neural and RBMT outputs, similarly to Aulamo et al. (2021). We evaluate our models both for translation quality and efficiency, resulting in a diverse set of submissions that balance accuracy and speed. Our contributions, including our code and data, are publicly available for further research in our Github repository.

The rest of the paper is organised as follows. Section 2 describes the benchmarking of existing models. Section 3 provides a detailed description of our data collection and preparation efforts. Section 4 describes the submitted models in detail. Section 5 outlines the results and, finally, section 6 concludes our work.

| Model | Params (M) | spa–xxx | | | xxx–spa | | | Speed (s) |
| | | spa–arg | spa–arn | spa–ast | arg–spa | arn–spa | ast–spa | |
|---|---|---|---|---|---|---|---|---|
| Apertium | – | 53.8 | 32.5 | 14.4 | 56.4 | 31.3 | 17 | 99.00 |
| opus-mt/itc-itc | 212 | 42.0 | 9.4 | 15.3 | 70.1 | 32.9 | 22 | 284.24 |
| opus-mt/deu+eng+fra+por+spa-itc* | 222 | **42.6** | **9.6** | **16.9** | – | – | – | 307.06 |
| opus-mt/roa-deu+eng+fra+por+spa** | 222 | – | – | – | **71.1** | 37.5 | 22.4 | 289.09 |
| opus-mt/itc-deu+eng+fra+por+spa | 222 | – | – | – | 70.4 | **37.7** | 22.3 | 245.22 |
| nllb-200-distilled-600M | 600 | – | 8.55 | 13.38 | – | 30.38 | 21.66 | 442.94 |
| nllb-200-distilled-1.3B | 1,300 | – | 8.66 | 13.95 | – | 32.59 | 22.48 | 800.38 |
| nllb-200-1.3B | 1,300 | – | 8.62 | 12.46 | – | 34.35 | 22.85 | 822.73 |
| nllb-200-3.3B | 3,300 | – | 8.75 | 13.38 | – | 35.16 | **23.65** | 914.43 |

Table 1: BLEU scores on the development set for all language pairs in both directions of existing MT models. We also report the average decoding speed in seconds on a single Nvidia V100 GPU. The asterisks (* and **) indicate that we use those models for our work.

## 2 Benchmarking of Existing Models

The first step we took when approaching this task was to benchmark existing models for the target languages. This process enabled us to assess the current landscape of available models, identify those suitable for fine-tuning, and determine which models could be utilized for back-translation.

We evaluate three types of models: OPUS-MT models[1], the smaller NLLB variants (Costa-jussà et al., 2022), and the rule-based Apertium systems. OPUS-MT models are trained with the Tatoeba Translation Challenge dataset[2] (Tiedemann, 2020) and data from the massively parallel Bible corpus (Mayer and Cysouw, 2014) as part of the JHUBC corpus (McCarthy et al., 2020). These are all *transformer-big* (Vaswani et al., 2017) systems. The NLLB models are trained on a diverse collection of multilingual text and come in different sizes. Apertium and OPUS-MT cover all three target languages, whereas NLLB does not support Aragonese.

We evaluate the systems with the provided development set by the organizers and BLEU (Papineni et al., 2002), as implemented in sacreBLEU (Post, 2018). The development set consists on a manually-crafted revision of the 997 sentences from Flores+ (Goyal et al., 2022). Results are shown in Table 1.

We can see how OPUS-MT models, although much smaller in size, outperform all NLLB variants when the target language is the Romance minority language. All NLLB variants perform similarly, independently of their size. We attribute the lower score on the Spanish–Aranese language pair to the

models being trained on Occitan data rather than Aranese. Additionally, the remarkable performance of the Apertium models in translating from Spanish stands out, as they surpass the neural systems, except in the case of Asturian.[3] This demonstrates the effectiveness of rule-based systems in handling closely related languages.

With Spanish as the target language, the NLLB models follow the scaling laws and their score increases along with their size, as would be expected. The OPUS-MT models exhibit comparable performance. As expected, the compact OPUS-MT models are much faster when considering the decoding speed of the different models. OPUS-MT models have been trained using Marian (Junczys-Dowmunt et al., 2018), while the NLLB family was trained using Fairseq (Ott et al., 2019).

Taking this into account, we decide to select two OPUS-MT models from Table 1 for our work: the model marked with * for fine-tuning; and model ** for producing back-translations. Given that the NMT models significantly outperform the Apertium systems for translation into Spanish, the rule-based back-translation strategy employed by Aulamo et al. (2021) did not suit our context.

## 3 Data

The data used to train our NMT systems consists of parallel and monolingual datasets provided by the organizers, as well as additional Wikipedia and dictionary data. We utilize the monolingual datasets by back-translating them to create synthetic parallel

---

[1] https://github.com/Helsinki-NLP/OPUS-MT-train
[2] Version v2023-09-26.

[3] Asturian sentences were professionally translated from English, while Aragonese and Aranese sentences were machine translated from Spanish using Apertium and later post-edited. Hence, the higher score for Apertium and the language pairs involving Aragonese and Aranese.

| | Monolingual | | | | | | Parallel | |
| | PILAR crawled | PILAR literary | PILAR cat–arn | Dictionary | Wikipedia | Wikipedia Discussions | Tatoeba Challenge | TOTAL |
|---|---|---|---|---|---|---|---|---|
| **arg** | | | | | | | | |
| raw | 60,028 | 24,675 | – | – | 255,149 | – | 41,623 | 381,475 |
| langid | 60,028 | 20,241 | – | – | 241,415 | – | 22,354 | 344,038 |
| filtered | 56,103 | 19,328 | – | – | 237,793 | – | 19,479 | **332,703** |
| **arn/oci** | | | | | | | | |
| raw | 7,358 | 229,886 | 85,491 | 14,874 | 616,530 | 14,591 | 744,731 | 1,713,461 |
| **arn** | | | | | | | | |
| langid | 7,358 | 228,512 | 64,141 | 14,874 | 29,627 | 2,429 | 106,248 | 453,189 |
| filtered | 7,243 | 213,960 | 64,141 | 14,874 | 27,160 | 2,249 | 87,189 | **337,801** |
| **oci** | | | | | | | | |
| langid | 0 | 474 | 0 | 0 | 511,713 | 11,415 | 354,202 | 877,804 |
| filtered | 0 | 357 | 0 | 0 | 493,216 | 10,810 | 299,440 | **803,823** |
| **ast** | | | | | | | | |
| raw | 14,776 | 24,093 | – | 82,009 | 2,230,855 | – | 5,511,336 | 7,863,069 |
| langid | 10,538 | 17,112 | – | 82,009 | 1,920,758 | – | 3,705,483 | 5,735,900 |
| filtered | 9,975 | 16,072 | – | 82,009 | 1,862,821 | – | 991,617 | **2,880,485** |

Table 2: Number of sentence pairs in training datasets. The "raw" line shows the sizes before any filtering, the "langid" line shows the number of sentence pairs in the correct language according to Idiomata Cognitor, and the "filtered" line shows the final sizes of the clean datasets. The Aranese and Occitan data are separated into two sets after language identification. The final cleaned data size for each language is shown in bold.

training data. We remove noise from the training data using the Idiomata Cognitor (Galiano-Jiménez et al., 2024a) language identification tool, Opus-Cleaner (Bogoychev et al., 2023) and its visual user interface, and the configurable filtering toolbox OpusFilter (Aulamo et al., 2020).

## 3.1 Data Collection

Table 2 shows the sizes of the datasets used for training. As original parallel data, we use only the Tatoeba Challenge data[4] (Tiedemann, 2020), which contains all data in OPUS (Tiedemann, 2012), deduplicated and shuffled. We also use the crawled and literary PILAR corpora (Galiano-Jiménez et al., 2024b) as monolingual data for all three language pairs. Additionally, we use the Aranese side of the Catalan-Aranese PILAR corpus also as monolingual data.

We also leverage monolingual data that is not provided by the organizers, which puts our models in the open track: Wikipedia and online dictionaries for Aranese[5] and Asturian[6]. From Wikipedia, we obtain the latest dump per language. Moreover, for Occitan, we also make use of OcWikiDisc (Miletic and Scherrer, 2022), a corpus extracted from the talk pages associated with the Occitan

Wikipedia. We assume that Occitan datasets include Aranese data, since it is a variety of Gascon, one of the main dialects of Occitan. For the online dictionaries, we develop our own scraping scripts to gather definitions. The scripts can be found in our Github repository. For the monolingual data, we produce back-translations into Spanish with the openly available OPUS-MT model (marked with ** in Table 1) to produce synthetic parallel data.

## 3.2 Data Cleaning

The first step of our data cleaning pipeline is language identification. We use the Idiomata Cognitor tool (Galiano-Jiménez et al., 2024a) to identify the correct target languages in all data sets. Idiomata Cognitor also allows us to distinguish Aranese from other Occitan varieties. Hence, from this point onwards, we treat Aranese and (non-Aranese) Occitan data separately in order to experiment with different model training strategies as described in Section 4.

Next, we create customized filtering configurations for each corpus (and for each subcorpus in Tatoeba) to apply optimal data cleaning based on the style and domain of the texts. To this end, we use OpusCleaner (Bogoychev et al., 2023), which is a parallel data cleaning tool that allows the user to add and adjust filters and see their effects on a sample of the corpus in real time in a graphical interface. The filters most commonly applied

---

[4]We use the same version as the original OPUS-MT model.
[5]https://www.diccionari.cat/cerca/diccionari-der-aranes
[6]https://diccionariu.alladixital.org/

| Model | Submission | oci tag | arn tag | spa–arg | spa–arn | spa–ast |
|---|---|---|---|---|---|---|
| OPUS-MT | - | - | - | 42.6 | 9.6 | 16.9 |
| A.1 | 1 | »oci« | »oci« | **54.8** | 12.3 | 18.5 |
| A.2 | 1 | – | »oci« | 51.5 | **28.2** | 18.5 |
| B.1 | – | »oci« | »xxx« | 51.6 | 26.1 | 18.5 |
| B.2 | – | »xxx« | »oci« | 51.5 | **26.8** | 18.5 |
| B.3 | – | »oci«»oci« | »oci«»xxx« | **55.2** | 26.0 | 18.5 |
| B.4 | – | »oci«»xxx« | »oci«»oci« | 52.8 | 25.7 | 18.5 |

Table 3: BLEU scores on the development set of fine-tuning the OPUS-MT model with different tagging strategies. We provide the scores of the OPUS-MT model for reference. We report the best checkpoint score per language pair. Model A.2 does not use Occitan data.

to our training sets are: (1) `src_trg_ratio`: The ratio between the number of source and target tokens. (2) `num_mismatch`: The ratio between the number of overlapping and differing numerals. (3) `alpha_ratio`: The ratio between the number of words and non-words, and the ratio between the number of language and non-language characters.

Additionally, some corpora contain unwanted structures, such as HTML tags or transcription content between double square brackets in Wikipedia data, which we remove from the sentences. Finally, we apply OpusFilter to concatenate the different corpora, normalize whitespace characters, remove all sentences shorter than 3 or longer than 150 words and remove all duplicate sentence pairs. Table 2 shows the data sizes for each language pair after applying language identification and corpus cleaning. The final size of our corpus is 4,35M sentences, with 66.14% spa–ast, 7.63% spa–arg, 7.75% spa–arn, and 18.45% spa–oci. All of our data cleaning configuration files can be found in our Github repository.

# 4 Models

In this section, we detail our modeling choices for the four submissions, all of which employ one-to-many multilingual models. Our models leverage fine-tuning and data tagging, and the integration of RBMT with neural models via Sequence-Level KD (Seq-KD) (Kim and Rush, 2016). All models are based on the Transformer architecture (Vaswani et al., 2017) and use the OPUS-MT model, as described earlier, as the initial checkpoint in some form. For tokenization, we use the OPUS-MT model's SentencePiece vocabularies (Kudo and Richardson, 2018), two distinct 32k piece vocabularies: one shared among all source languages (in our case, only Spanish) and another shared among all targets. All models are trained on 4 Nvidia V100 GPUs, except models C.2 and D.2, which are trained on 8 AMD MI250x GPUs. Further configuration details are provided in Appendix A.

## 4.1 Models A: Fine-tuning

As an initial step, we use the openly available OPUS-MT model described in Section 2 and fine-tune it using different data sampling schemes. We train one model with all available training data (model A.1) and another excluding the Occitan data (A.2). The decision to exclude Occitan data was made because both languages share the same language tag, which could potentially confuse the model, since there is much more training data on Occitan than on Aranese. The development set scores are presented in Table 3.

Compared to the original OPUS-MT model, we observe a significant increase in BLEU scores for the spa–arg language pair (+12.2) and for spa–arn (+18.6). However, for Asturian, the increase is more modest (+1.6). Removing the Occitan data results in an increased score of almost +16 BLEU points for the spa–arn language pair. Interestingly, despite having the largest amount of new data for Asturian, the model quickly reaches a performance plateau during training, as shown in Figure 2 in Appendix B. This trend persists throughout our experiments, leading us to conclude that the spa–ast language pair is the most challenging task.

For our **Submission #1**, we ensemble the best $n$ checkpoints per language pair across both A.1 and A.2 models.[7]

---

[7] We perform ensembling using the top 10 best checkpoints for each language pair and submit the ensemble with the highest score on the development set.

| Model | Teacher(s) | Size | Submission | spa–arg | spa–arn | spa–ast |
|---|---|---|---|---|---|---|
| C.1 | A | base | 2 | 53.6 | **28.1** | 18.5 |
| C.2 | A | tiny | – | 51.3 | 25.5 | 18.2 |
| C.3 | B.3 | base | – | 55.4 | 26.6 | 18.5 |
| D.1 | A + RBMT | base | 2 | **54.2** | 27.3 | 18.5 |
| D.2 | A + RBMT | tiny | 3 | **52.8** | **27.1** | 18.2 |
| D.3 | B.3 + RBMT | base | 4 | 56.9 | **30.2** | 18.5 |
| D.3_fixed | B.3 + RBMT | base | – | **57.0** | 26.9 | 18.5 |
| D.4 | RBMT | base | – | 62.4 | 36.8 | 16.9 |

Table 4: Comparison of BLEU scores for our distillation experiments between NMT-only models and NMT+RBMT systems across different language pairs on the development set. We report the best checkpoint score per language pair, except for models D, where we use the same single checkpoint.

## 4.2 Models B: Data tagging

In multilingual systems, it is a common practice to prepend a language tag to the source sentence to indicate the target language. For consistency, we applied uniform tagging across all models. Nevertheless, for Aranese, we experimented with different tagging schemes, given that it is a variety of Gascon, a dialect of Occitan.

Exploring the OPUS-MT vocabulary, we identified an unused tag, »xxx«, which prompted us to experiment with various combinations of the »oci« and »xxx« tags, including the use of double tags. We fine-tune the original OPUS-MT model with all available training data and different tagging schemes. Results of are provided in Table 3.

While the performance of Asturian remained unaffected, using different tags for Aranese and Occitan led to a much higher BLEU on the spa–ara language pair compared to model A.1; due to the effectiveness of the data tagging schemes. Notably, Aragonese appeared to be the most impacted by data tagging, although we are unsure why. On average, the best performing model leverages double tags (B.3). We do not submit any of these, but use B.3 as a teacher in our distillation experiments.

## 4.3 Sequence-Level Distillation

Seq-KD (Kim and Rush, 2016) is a technique where a student model is trained using translations generated by one or more teacher model(s), with the goal of transferring knowledge from a large, powerful teacher model to a smaller, more efficient student model. We experiment with Seq-KD to train fast students.

**Models C: NMT-distilled**

First, we distill student models using the previously fine-tuned *transformer-big* NMT systems as teachers. We leverage Sequence-Level Interpolation (Kim and Rush, 2016), generating 8-best candidate translations for all the training data using the best checkpoint for each language pair. From these, we select the translation with the highest ChrF (Popović, 2015) with the reference to create a distilled dataset, which is then used to train the student model. We use ChrF instead of BLEU as a more fine-grained metric at character level.

To explore the tradeoff between translation quality and speed, we use models A[8] as the teachers and train two students of different sizes. Model C.1 is a *transformer-base* model (67.5M parameters), while model C.2 follows the *tiny* architecture described in Bogoychev et al. (2020), its size is 20.4M parameters (3.3 times smaller). We train model C.2 using the OpusDistillery[9], a pipeline for multilingual Seq-KD of open NMT models. In addition, to investigate the effect of multi-teacher distillation, we distill another *transformer-base* model (C.3) using a single NMT teacher, in this case, model B.3. The development set scores for these student systems are shown in Table 4.

When comparing models C.1 and C.2, it becomes evident that the capacity gap between the teacher and student models significantly impacts student performance. In KD, student models are typically smaller than their teacher counterparts, which can hinder their ability to effectively learn and fit noisy data. This is reflected in the lower scores of the smaller C.2 model across all language pairs compared to C.1. On the other hand, models C.1 and C.3 share the same size, but their training

---

[8]For Occitan, we use the original OPUS-MT model, as our fine-tuned model has a lower score on Occitan, due to the catastrophic forgetting phenomenon (Goodfellow et al., 2013).

[9]https://github.com/Helsinki-NLP/OpusDistillery

Figure 1: Overview of our Seq-KD distillation process to merge NMT and RBMT data. Given a source sentence *(S)*, we produce a hypothesis translation *(H)* with both our RBMT and NMT models. Then, we choose the translation *(H\*)* that has the maximum ChrF with the ground truth *(G)* to create the distilled dataset *(D)*.

strategies differ. Model C.1 is distilled from multiple teachers, while C.3 is distilled from a single teacher. Notably, distilling from a single model in C.3 appears to offer greater stability.

**Models D: Hybrid-distilled**

Since rule-based translation models are remarkably good for the given language pairs, as shown in Table 1, we further experiment with Seq-KD to train student models that benefit from both RBMT and NMT outputs.

In this case, we use two types of teachers: (1) the best checkpoint per language pair of the NMT model(s) (as in the previous section) and (2) the Apertium RBMT models. We forward translate the training data with both teachers. For each source sentence, we select the translation that has the highest ChrF score with the ground truth to create the distilled dataset. Finally, we train a new student model on the distilled dataset. An overview of this process is depicted in Figure 1.

For each of the former models C, we train a comparable hybrid-distilled student using a combination of the NMT and RBMT data. The development set scores for these models (D.1–D.3) are shown in Table 4. The proportions of RBMT data selected for the final distilled dataset are provided in Appendix C.

The inclusion of RBMT data in the distillation process leads to better performance across all language pairs overall. For model D.1, the addition of rule-based distilled data results in a slight decrease for spa–arn, in comparison to C.1. For spa–ast, the performance is identical across all models. It is remarkable to note that model D.3 surpasses the performance of its own teacher with +1.7 BLEU for spa–arg and +4.2 BLEU for spa–arn.

After the submission deadline, we discovered that the NMT distilled dataset for model D.3 had been generated using incorrect language tags for Aranese and Occitan (»oci« instead of »oci«»oci« and »oci«»xxx«, respectively). We provide the corrected results for the D.3 model

(model D.3_fixed in Table 4). Interestingly, the initial D.3 model performed better for Aranese due to the higher proportion of RBMT data for that language (as shown in Appendix C, 15% vs. 2.3%), which favored the RBMT-heavy development set. Motivated by this finding, we trained a student using RBMT-only distilled data after the submission deadline (model D.4 in Table 4), which outperforms all other models except for Asturian. This opens up a new avenue of research, leveraging linguistically informed methods for distillation.

Table 4 demonstrates the effectiveness of using distillation to train a single model that performs well across all three language pairs. Among the language pairs, Aragonese shows the most significant improvement when RBMT data is incorporated, highlighting the particular benefit of combining rule-based and neural translation methods for this language. This aligns with our expectations since as can be seen from Table 1, the spa–arg Apertium model achieves the highest BLEU score.

Out of our distillation experiments, we make three submissions. For **submission #2**, we ensemble the best $n$ checkpoints per language pair across models C.1 and D.2. For **submission #3**, we submit model D.2. In this case, we do not use ensembling because we want to test it for speed. Finally, for **submission #4**, we ensemble the best $n$ checkpoints per language pair from model D.3.

## 5 Results

We make four submissions in the open submission track. The test set corresponds to the 1,012 lines of Flores+ evaluation set. We summarize our submissions' test results in Table 5, as provided by the organizers of the Shared Task. For comparison, we also include the scores of the top-performing competitor, overall and in the open submission track. The official evaluation metrics of the task are BLEU and ChrF. Additionally, we report the average decoding speed and the model sizes.

Our best models ranked 4th, 5th, and 2nd in the open submission track for spa–arg, spa–ara, and

| # | Method | BLEU / ChrF | | | Params (M) | Speed (s) |
|---|---|---|---|---|---|---|
| | | arg | arn | ast | | |
| 1 | Fine-tuning Data Sampling Ensembling | 51.5 / 75.6 | 22.1 / 45.1 | 18.2 / 51.6 | 222.9 | 852.22 |
| 2 | Distillation RBMT+NMT Ensembling | 50.6 / 75.4 | 22.4 / 45.7 | 18.0 / 51.6 | 65.7 | 361.33 |
| 3 | Distillation RBMT+NMT | 49.1 / 75.4 | 21.6 / 45.0 | 17.9 / 51.4 | 20.4 | 4.06 |
| 4 | Distillation Data Tagging RBMT+NMT Ensembling | 52.7 / 75.9 | 24.3 / 46.6 | 18.0 / 51.5 | 67.5 | 891.76 |
| Best (overall) | – | 63.0 / 80.3 | 30.4 / 50.1 | 23.2 / 55.2 | – | – |
| Best (open) | – | 62.7 / 80.0 | 28.8 / 49.4 | 23.2 / 55.2 | – | – |

Table 5: Summary of our submissions. BLEU refers to the score obtained by the best ensemble on the development set; Speed refers to the averaged decoding speed for submission across language pairs on one single AMD MI250x GPU. In addition, we provide the best competitor scores for each target language.

spa–ast, respectively. On average, our best submission for each language pair falls short of the top competitor by 4 BLEU points and 3.8 ChrF points. This narrow margin reflects the competitive nature of this year's task, which saw over 178 submissions.

Our best model is submission #4, followed closely by submissions #1, #2 and, finally #3, in that order. It is noteworthy that our distilled models perform really well compared to their teachers. Submission #2, a distilled model from Submission #1, demonstrates an increase of +0.3 BLEU for spa–arn over its teacher, highlighting the potential of distillation to not only preserve but even enhance translation quality. Moreover, our smallest model, Submission #3, although showing a slight average decrease of –1.1 BLEU compared to its teacher, offers a significant advantage in terms of speed—it is 210 times faster.

## 6 Conclusions

In this work, we have presented our participation in the Shared Task of Translation into Low-Resource Languages of Spain 2024. We have described our data collection and preparation efforts, as well as our four submissions based on multilingual models. We explore fine-tuning of an existing open model with different data tagging schemes and use Seq-

KD to train small efficient student models. Furthermore, to our knowledge, we are the first to leverage RBMT to improve distillation for similarly related languages and prove its effectiveness.

This study opens up new research directions for advancing in low-resource MT by demonstrating the potential of data tagging strategies and hybrid distillation methods, ensuring these languages are both preserved and accessible in the digital age.

## 7 Ethical Considerations

In addition to evaluating the performance of our models in terms of translation quality, it is equally important to consider the computational resources required for their training and deployment. By analyzing the GPU consumption of our experiments, including the time spent and energy consumed for each task, we aim to provide a comprehensive assessment of the efficiency and sustainability of our approaches. This will allow the community to take informed decisions about model selection and optimization in real-world applications, where computational efficiency is often as critical as accuracy. We report the energy consumption of the totality of our experiments in Table 6, which amounts to 508 kWh.

| Task | Model | Time (h) | Energy (kWh) |
|---|---|---|---|
| Back-translation | | 18.9 | 19.5 |
| Fine-tune | A.1 | 35.3 | 37.0 |
| | A.2 | 22.4 | 23.0 |
| | B.1 | 50.5 | 52.2 |
| | B.2 | 27.4 | 28.9 |
| | B.3 | 28.9 | 29.6 |
| | B.4 | 28.0 | 28.8 |
| Forward translation | | 7.9 | 6.6 |
| Train via Seq-KD | C.1 | 64.6 | 66.9 |
| | C.2 | 11.5 | 18.1 |
| | C.3 | 57.1 | 54.2 |
| | D.1 | 53.4 | 55.3 |
| | D.2 | 11.2 | 17.8 |
| | D.3 | 66.6 | 68.6 |
| | D.3_fixed | 55.9 | 57.1 |
| | D.4 | 56.0 | 57.8 |
| Ensembling | | 30.2 | 3.52 |
| Submission | | 1.7 | 0.19 |
| | Total | 627.3 | **625.1** |

Table 6: Energy consumption of our work. We report the time (hours) and energy consumption across the different tasks of our experiments, run on 4 Nvidia V100 GPUs. The training of models D has been run on 8 AMD MI250x GPUs. Ensembling and translations for submission have been run on 1 Nvidia V100 GPUs.

## Acknowledgments

## References

Mikko Aulamo, Sami Virpioja, Yves Scherrer, and Jörg Tiedemann. 2021. Boosting neural machine translation from Finnish to Northern Sámi with rule-based backtranslation. In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 351–356, Reykjavik, Iceland (Online). Linköping University Electronic Press, Sweden.

Mikko Aulamo, Sami Virpioja, and Jörg Tiedemann. 2020. OpusFilter: A configurable parallel corpus filtering toolbox. In *Proceedings of the 58th Annual Meeting of the Association for Computational Lin-*guistics: System Demonstrations, pages 150–156, Online. Association for Computational Linguistics.

Nikolay Bogoychev, Roman Grundkiewicz, Alham Fikri Aji, Maximiliana Behnke, Kenneth Heafield, Sidharth Kashyap, Emmanouil-Ioannis Farsarakis, and Mateusz Chudyk. 2020. Edinburgh's submissions to the 2020 machine translation efficiency task. In *Proceedings of the Fourth Workshop on Neural Generation and Translation*, pages 218–224, Online. Association for Computational Linguistics.

Nikolay Bogoychev, Jelmer van der Linde, Graeme Nail, Barry Haddow, Jaume Zaragoza-Bernabeu, Gema Ramírez-Sánchez, Lukas Weymann, Tudor Nicolae Mateiu, Jindřich Helcl, and Mikko Aulamo. 2023. Opuscleaner and opustrainer, open source toolkits for training machine translation and large language models. *Preprint*, arXiv:2311.14838.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Mikel L Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine translation*, 25:127–144.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024a. Idiomata cognitor.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024b. Pan-iberian language archival resource.

Ian J. Goodfellow, Mehdi Mirza, Xia Da, Aaron C. Courville, and Yoshua Bengio. 2013. An empirical investigation of catastrophic forgeting in gradient-based neural networks. *CoRR*, abs/1312.6211.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 Evaluation Benchmark for Low-Resource and Multilingual Machine Translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann,

Alham Fikri Aji, Nikolay Bogoychev, et al. 2018. Marian: Fast neural machine translation in c++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121.

Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327.

Young Jin Kim, Marcin Junczys-Dowmunt, Hany Hassan, Alham Fikri Aji, Kenneth Heafield, Roman Grundkiewicz, and Nikolay Bogoychev. 2019. From research to production and back: Ludicrously fast neural machine translation. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 280–288, Hong Kong. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Thomas Mayer and Michael Cysouw. 2014. Creating a massively parallel Bible corpus. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3158–3163, Reykjavik, Iceland. European Language Resources Association (ELRA).

Arya D. McCarthy, Rachel Wicks, Dylan Lewis, Aaron Mueller, Winston Wu, Oliver Adams, Garrett Nicolai, Matt Post, and David Yarowsky. 2020. The Johns Hopkins University Bible corpus: 1600+ tongues for typological exploration. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2884–2892, Marseille, France. European Language Resources Association.

Aleksandra Miletic and Yves Scherrer. 2022. OcWikiDisc: a corpus of Wikipedia talk pages in Occitan. In *Proceedings of the Ninth Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 70–79, Gyeongju, Republic of Korea. Association for Computational Linguistics.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the*

*Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Felipe Sánchez-Martínez, Juan Antonio Pérez-Ortiz, Aarón Galiano-Jiménez, and Antoni Oliver. 2024. Findings of the WMT 2024 shared task on translating into low-resource languages of spain: Blending rule-based and neural systems. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218.

Jörg Tiedemann. 2020. The Tatoeba Translation Challenge – realistic data sets for low resource and multilingual MT. In *Proceedings of the Fifth Conference on Machine Translation*, pages 1174–1182, Online. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

## A  Hyperparameters

All models are based on the transformer architecture. They all share the following: the Adam optimizer is used with $\beta_1$=0.9 and $\beta_2$=0.998. The models are trained until convergence with early-stopping on development data after BLEU has stalled 10 times. Next, we specify each model's unique configuration details.

**Models A and B** are *transformer-big* models. They use a 6-layered transformer with 16 heads, 1024 dimensions in the embeddings and 4,096 dimensions in the feed-forward layers.

**Models C.1, C.3, D.1 and D.3** use a 6-layered transformer with 8 heads, 512 dimensions in the embeddings and 2,048 dimensions in the feed-forward layers.

**Models C.2 and D.2** are trained using tiny architecture proposed in Bogoychev et al. (2020). The student model has a transformer encoder with 6 layers and a light-weight RNN based decoder with Simpler Simple Recurrent Unit (SSRU) (Kim et al., 2019) with 2 layers; 8 heads, 256 dimensions in the embeddings and 1,536 dimenstions in the feed-forward layers.

## B  Learning Curves

Figure 2 shows the BLEU score progression over training updates per language pair for model A.2. It shows how the performance for spa–ast quickly reaches a plateau.



Figure 2: BLEU score progression over training updates and epochs for model A.2.

## C  Rule-based MT Data

For the distilled models D, we use a combination of NMT and RBMT teachers to build a distilled dataset. The RBMT teachers are the Apertium models. For each source sentence, we generate a hypothesis translation using both teachers and then compute the ChrF score against the ground truth. We retain the hypothesis with the highest ChrF score for each sentence. Table 7 shows the proportion of sentences originating from RBMT across our experiments.

| Teacher(s) | A.1, A.3 | B.3 | B.3 |
|---|---|---|---|
| Model(s) | C.1, C.2 | D.3 | D.3_fixed |
| Submission | #2 #3 | #4 | - |
| Pair | % | % | % |
| spa–arg | 4.39 | 7.37 | 7.37 |
| spa–arn | 1.32 | 15.32 | 2.33 |
| spa–ast | 3.85 | 3.75 | 3.75 |
| spa–oci | 8.95 | 1.69 | 1.64 |

Table 7: Distribution of distilled data coming from RBMT in sentence count and percentage (%).

# Robustness of Fine-Tuned Models for Machine Translation with Varying Noise Levels: Insights for Asturian, Aragonese and Aranese

[1]**Martin Bär**, [1]**Elisa Forcada Rodríguez** and [2]**María García-Abadillo Velasco**

[1]Erasmus Mundus Master in Language and Communication Technologies (LCT)

[2]Master in Language Analysis and Processing (LAP)

University of the Basque Country (UPV/EHU)

{mbr001, eforcada001, mgarciaabadill001}@ikasle.ehu.eus

## Abstract

We present the LCT-LAP proposal for the shared task on *Translation into Low-Resource Languages of Spain* at WMT24 within the constrained submission category. Our work harnesses encoder-decoder models pretrained on higher-resource Iberian languages to facilitate MT model training for Asturian, Aranese and Aragonese. Furthermore, we explore the robustness of these models when fine-tuned on datasets with varying levels of alignment noise. We fine-tuned a Spanish-Galician model using Asturian data filtered by BLEU score thresholds of 5, 15, 30 and 60, identifying BLEU 15 as the most effective. This threshold was then applied to the Aranese and Aragonese datasets. Our findings indicate that filtering the corpora reduces computational costs and improves performance compared to using nearly raw data or data filtered with language identification. However, it still falls short of the performance achieved by the rule-based system Apertium in Aranese and Aragonese.

## 1 Introduction

Spain is home to a rich linguistic landscape, yet this diversity is accompanied by disparities in terms of speaker numbers and language resources. Languages with co-official status, such as Basque, Catalan and Galician, were previously considered to have limited resources but are now included in numerous popular LLMs. Consequently, research in this field has shifted towards cases where data scarcity is even more pronounced, such as Asturian, Aragonese and Aranese. These languages are the focus of a shared task at the Conference on Machine Translation 2024. The objectives of this task include investigating transferability among Romance languages and determining the most effective methods for utilizing pretrained models in translations between Spanish and low-resource Romance languages.

The methodology employed involved the following steps:

1. Implementing automated methods for curating data. The constrained submission framework enables researchers to utilize corpora that may be notably noisy. Our work aims to propose solutions to this challenge.

2. Creating synthetic data for the monolingual PILAR Galiano-Jiménez et al. (2024b) corpora.

3. Harnessing models trained on other, resource-richer (Iberian) Romance languages with the presumption that this facilitates cross-lingual transfer. The model fine-tuned for Asturian was originally trained on Galician, while the models fine-tuned for Aranese and Aragonese were originally trained on Catalan.

The official metrics for the shared task are BLEU and chrF. The metrics employed in this study are BLEU and chrF++, as they are relatively straightforward to calculate and there is currently no robust neural-based metric for our target languages.

## 2 Background

### 2.1 Spanish Linguistic Landscape

Although the official language in Spain is Spanish, it coexists with other co-official and minority languages. The predominance of Spanish over the other languages and dialects is associated with historical reasons: since the Middle Ages, Spain had undergone a process of Castilianisation, which became very important in the 14th century, when the dominance of the Kingdom of Castile in the centre of the Iberian Peninsula led to the expansion of the use of Castilian. This continued until the 20th century, with the consequent marginalisation of the other vernacular languages (Martínez, 1982). The

co-official languages, Basque, Catalan and Galician, were considered to have limited resources in the past. This picture has changed, as efforts from both research and industry have contributed significantly to integrating them into the field of Language Technology. However, there are also non-official Ibero-Romance languages that are considered as having limited resources:

- Asturian: spoken in Asturias, the northeastern part of Leon, Zamora and the north of Portugal (ARIAS, 2002).

- Aragonese: spoken in the north of the province of Huesca and in the extreme northwest of Zaragoza (Marco Villanueva, 2012).

- Aranese: a variant of Occitan, spoken in the province of Aran (Rey and Canalís, 2006).

In this context, initiatives such as PILAR (Pan-Iberian Language Archival Resource) work to enrich and expand the resource availability of these languages (see Section 3).

## 2.2 Other related works

The interest in low-resource languages has recently increased, leading to a considerable amount of research on the subject (Ranathunga et al., 2023). Several studies on machine translation for low-resource languages can be found, such as the article by Karakanta et al. (2018), which works with non-parallel corpora, or Kumar et al. (2021), which focuses on recasting systems from high-resource languages to low-resource languages.

As far as Iberian languages are concerned, there are other investigations, such as the one published by Oliver et al. (2023), which explores techniques for training NMT systems applied to high- and low-resource Iberian languages or the work by Ko et al. (2021), which adapts high-resource NMT models to translate low-resource languages related to Spanish.

With respect to WMT, since its first edition in 2016, there have been three shared tasks related to the field: In 2020, a task was proposed on unsupervised and very low-resource languages, focusing on Upper Sorbian (Fraser, 2020). The following year, a workshop on multilingual low-resource translation for Indo-European languages was presented, focusing on North Germanic languages such as Icelandic and Romance languages such as Occitan (Libovický and Fraser, 2021). Finally, in 2022,

a task related to unsupervised MT and very low-resource supervised MT was suggested, with Upper and Lower Sorbian languages (Weller-Di Marco and Fraser, 2022).

## 3 Data

This research falls into the constrained submission category, as all data used was obtained from the mentioned sources in the shared task: the Open Parallel Corpora, also known as OPUS (Tiedemann, 2009), and the Pan-Iberian Language Archival Resource, shortened as PILAR (Galiano-Jiménez et al., 2024b).

The Spanish development set is part of the FLORES+ Evaluation Benchmark (NLLB Team et al., 2022). The Asturian, Aragonese, and Aranese counterparts of FLORES+ are published alongside PILAR.

Finally, both the BLEU reference translations for the OPUS data and the synthetic Spanish counterparts for the PILAR data were generated with Apertium (Forcada and Tyers, 2016).

## 3.1 OPUS corpora

OPUS is a public multilingual collection of parallel corpora that gathers open-source documents available on the Internet (Tiedemann and Thottingal, 2020) and supports 744 languages. The constraint submission is limited to all data in OPUS, thereby enabling researchers to create synthetic translations from other languages into Asturian (ast), Aragonese (arg), Aranese (arn), or Spanish (es). However, the data utilized in this work exclusively employs the corpora for the combinations es <> ast/arg/arn.

Given that the collected corpora were not consistently well-aligned, we implemented a filtering pipeline, as detailed in Section 3.2, to produce a smaller but cleaner dataset. The effectiveness of this approach is reflected in the "BLEU 15" column of Table 1.

## 3.2 OPUS Data filtering

Around 8 million aligned sentences were collected from OPUS for the three target languages, although this number was significantly reduced when applying a filtering pipeline. Three main steps were followed to filter out invalid sentences:

- **Basic filtering:** removing unnecessary white spaces, empty lines, and characters not sup-

ported by the file encoding for all target languages.

- **Idiomata Cognitor:** filtering out all sentence pairs whose target language was not labeled as Asturian, Aragonese or Aranese and whose source language was not labeled as Spanish by Idiomata Cognitor (Galiano-Jiménez et al., 2024a), a high-precision classifier trained using Bayesian methods and capable of identifying 10 Romance languages.

- **BLEU threshold filtering**[1]: we first translated the Spanish counterparts of the Asturian/Aragonese/Aranese datasets into the respective target languages using Apertium. Next, we calculated BLEU scores for the original Asturian/Aragonese/Aranese sentences as references and their translations as hypotheses. Then, we filtered the datasets to various BLEU thresholds, assuming that alignments are more likely to be correct if the sentence pairs have high BLEU scores[2]. For Asturian, this was done using four different BLEU thresholds: 60, 30, 15 and 5 BLEU. For Aragonese and Aranese, we only used one threshold.

### 3.3 PILAR Corpora

PILAR is a recently created corpus of texts from different languages spoken in the Iberian Peninsula, including Asturian, Aragonese, Aranese, Balearic and Valencian.

For our purposes, the monolingual data from Asturian, Aragonese and Aranese, and the Aranese counterpart from the Catalan-Aranese parallel corpora was backtranslated into Spanish (see Table 2) using Apertium: backtranslation can be understood as providing monolingual training data with a synthetic sentence source obtained by automatically translating the target sentence into the source language (Sennrich et al., 2015).

---

[1] We used Bleualign as a reference (Sennrich and Volk, 2010). However, we did not calculate the BLEU score between the hypothesis and reference sentences for both languages, nor did we compute the subsequent harmonic mean, given the fact that Apertium web tool does not support translations from Asturian into Spanish. Instead, we limited our calculations to the BLEU score of the original Asturian/Aragonese/Aranese sentences as references and the translation of its Spanish counterpart obtained with Apertium as hypotheses.

[2] BLEU evaluates translations by comparing n-grams between the model output and a reference, favouring those that are closest in terms of word and order. This may favour sentences in both the source and target languages that are easier to translate for Apertium.

| | Asturian | | Aragonese | | Aranese | |
|---|---|---|---|---|---|---|
| | BLEU 15 | Raw | BLEU 15 | Raw | BLEU 15 | Raw |
| GNOME | 18,435 | 68,668 | 2,004 | 5,529 | 0 | 77 |
| KDE4 | 4,515 | 26,023 | - | - | 667 | 49,593 |
| NLLB | 585,683 | 6,470,015 | - | - | 65,797 | 925,448 |
| QED | 125 | 421 | 18 | 222 | 45 | 282 |
| Tatoeba | 58 | 159 | 3 | 13 | 5 | 189 |
| TED2020 | 40 | 116 | - | - | - | - |
| WikiMatrix | - | - | 13,639 | 33,724 | 7,398 | 35,805 |
| wikimedia | 27,776 | 45,506 | 2,908 | 4,457 | 629 | 1,980 |
| XLEnt | 0 | 274,257 | 3 | 16,822 | 0 | 99,472 |
| | **636,632** | 6,884,903 | **18,575** | 60,767 | **74,502** | 1,112,879 |

Table 1: Number of raw sentence pairs obtained from the OPUS repository and the final number of sentences after filtering them with a BLEU score threshold of 15.

| | **Asturian** | **Aragonese** | **Aranese** |
|---|---|---|---|
| crawled | 14,776 | 60,028 | 7,358 |
| literary | 24,093 | 24,675 | 229,886 |
| paragraphs | - | - | 86,568 |
| sentences | - | - | 64,141 |
| Total | 38,869 | 84,703 | 387,953 |

Table 2: Number of monolingual sentences from PILAR that were backtranslated with Apertium.

## 4 Methodology

The methodology of this work involved fine-tuning two pretrained models (see section 4.1) on backtranslated PILAR and filtered OPUS data (see section 3 and section 3.2). The total number of sentences for each language is presented in Table 3.

The experimental setup utilized a Tesla V100-PCIE-32GB GPU running with NVIDIA driver version 535.104.12 and CUDA version 12.2, alongside the HuggingFace Transformers library for model loading and fine-tuning.

All models underwent training for at least 1 epoch (Table 3 shows when each model converged). The best model selection was based on the BLEU score derived from the development set. Additionally, zero-shot translation without fine-tuning was conducted as a baseline for comparing results.

### 4.1 Models

Two models from `Helsinki-NLP` (Tiedemann and Thottingal, 2020) were used for our experiment:

- `opus-mt-es-gl`: a transformer-align model from Spanish into Galician that achieved a BLEU 67.6 and a chr-F score of 80 in the Tatoeba test. Given the close linguistic relationship between Asturian and Galician-Portuguese, we aimed to explore transfer learning when fine-tuning on Asturian data.

| | Model | Data | Sentences | Epochs | Steps | BLEU | chrF++ |
|---|---|---|---|---|---|---|---|
| **AST** | **apertium** | - | - | - | - | 17.1 | **50.69** |
| | es-gl-noft-ast | - | - | - | - | 5.75 | 38.66 |
| | es-gl-ft-basic | basic clean | 6,884,903 | 0.87 | 55k | 17.07 | 49.89 |
| | es-gl-ft-idiomata | idiomata cognitor | 4,521,302 | 1.76 | 36k | 17.32 | 50.24 |
| | es-gl-ft-bleu | bleu_60 | 440,794 | 2.61 | 14k | 17.61 | 50.39 |
| | es-gl-ft-bleu | bleu_30 | 582,883 | 3.95 | 22k | 17.79 | 50.48 |
| | es-gl-ft-bleu | bleu_5 | 743,846 | 3.44 | 25k | 17.84 | 50.57 |
| | es-gl-ft-bleu | bleu_15 | 636,632 | 2.61 | 18k | 17.85 | 50.46 |
| | **es-gl-ft-backtr** | **bleu_15 + PILAR** | **675,501** | **3.22** | **22k** | **17.90** | 50.58 |
| **ARG** | **apertium** | - | - | - | - | **66.05** | **82.23** |
| | es-ca-noft-arg | - | - | - | - | 8.38 | 46.23 |
| | es-ca-ft-arg | idiomata cognitor | 27,335 | 4.67 | 1k | 32.87 | 64.79 |
| | es-ca-ft-arg | basic clean | 60,767 | 2.91 | 1k | 33.17 | 65 |
| | es-ca-ft-arg | bleu_15 | 18,575 | 47.95 | 7k | 41.39 | 70.38 |
| | es-ca-ft-arg | bleu_15 + PILAR | 103,278 | 7.43 | 8k | 41.53 | 70.84 |
| **ARN** | **apertium** | - | - | - | - | **38.02** | **60.01** |
| | es-ca-noft-arn | - | - | - | - | 5.75 | 38.66 |
| | es-ca-ft-arn | idiomata cognitor | 383,575 | 4.67 | 14k | 9.61 | 40.67 |
| | es-ca-ft-arn | basic clean | 1,112,879 | 2.65 | 14k | 9.70 | 40.74 |
| | es-ca-ft-arn | bleu_15 | 74,502 | 6.86 | 9k | 10.19 | 41.88 |
| | es-ca-ft-arn | bleu_15 + PILAR | 462,455 | 0.83 | 8k | 29.04 | 54.85 |

Table 3: BLEU and chrF++ scores on the FLORES+ devset comparing baselines (`apertium` and models with `noft` in their names) and fine-tuned models (`-ft-`) across varying levels of alignment noise. Baselines always occupy the first two rows for each language. Subsequent models are listed in ascending order of BLEU scores. Best performing architectures are highlighted in bold.

- `opus-mt-es-ca`: a transformer-align from Spanish into Catalan with a BLEU score of 68.9 and a chr-F score of 0.832 in the Tatoeba test. We aimed to explore transfer learning when fine-tuning Catalan for Aranese and Aragonese.

## 5 Results

As Table 3 shows, the results of our experiments were compared with two baselines: Apertium, a rule-base system that supports translations in the same languages as those investigated in this work, and the respectively selected model for our experiments with zero-shot translations (i.e. without fine-tuning).

Overall, the results for Aragonese and Aranese show the same trend: the highest performance was achieved by fine-tuning on data filtered with a BLEU threshold of 15, combined with the back-translated PILAR corpora. While the backtrans-

lated data yielded improvements of 18.85 BLEU for Aranese, this improvement was only 0.14 BLEU for Aragonese. Interestingly, fine-tuning on data that had only undergone basic cleaning outperformed our approach of filtering out sentences in other languages. The zero-shot translation approach yielded the lowest results by a significant margin. Despite these efforts, our results still fall short of the baseline Apertium by approximately 9 points in Aranese and nearly 25 points in Aragonese.

Our best result for Asturian is the only one comparable to the baseline Apertium. Our fine-tuned model, which uses a BLEU score threshold of 15 and the PILAR corpora, outperforms the baseline by 0.8 BLEU points. However, it falls short of the baseline by 0.11 chrF++ points.

See the following sections for a more detailed description of each language's results.

## 5.1 Asturian

Our results show that setting a threshold of 15 BLEU for OPUS-aligned corpora yields the best performance in Asturian. It slightly outperforms thresholds of 5 and 30 BLEU and achieves an improvement of almost 0.25 over the cleanest filtered set of corpora with a threshold of 60 BLEU.

Note that an Asturian tokenizer was trained and implemented; however, its performance did not exceed a BLEU score of 17.6 and it was consequently omitted from Table 3. Consequently, no tokenizer was trained for Aranese and Aragonese.

Integrating backtranslated Asturian PILAR results in almost a 1-point BLEU score improvement compared to the slightly preprocessed raw OPUS data (basic clean in Table 3), and a slight improvement of 0.05 compared to the filtered OPUS data with 15 BLEU threshold without PILAR data.

Regarding the baselines, our best method (data filtered with a 15 BLEU threshold and backtranslated PILAR) achieves similar performance as Apertium, with a 0.8 BLEU score improvement and a 0.11 lower chrF++ score. The zero-shot translation results are by far the worst, with scores approximately 12 points below the best results.

## 5.2 Aragonese

As detailed in section 7, only the best filtering threshold for OPUS data in Asturian was also applied to Aragonese.

Our best result is again the result of fine-tuning the bleu_15 + PILAR corpora on a model initially fine-tuned for Spanish-Catalan translation. It outperforms the model finetuned on almost raw data by 8.36 BLEU points. Comparing these results to a model only trained on the bleu_15 data, reveals that using the backtranslated data only yielded an improvement of 0.14 BLEU. However, these results lag behind the Apertium baseline, which obtains scores with approximately 25 points difference in BLEU and around 12 points in chrF++.

The zero-shot baseline model (`es-ca-noft-arg`) achieved a similarly low score as the zero-shot models in the other languages and performed significantly worse than the other approaches. It lags behind the best result from Apertium by approximately 58 BLEU points and around 36 chrF++ points.

## 5.3 Aranese

As detailed in section 7, only the best filtering threshold for OPUS data in Asturian was applied to Aranese.

Showing the same trend as the results for the other two languages, the approach using bleu_15 + PILAR corpora is the most effective. It achieves an improvement of about 20 BLEU points and approximately 14 chrF++ points over the other data sets, which only underwent basic cleaning or language filtering with Idiomata Cognitor. In contrast to our results for Aragonese, the Aranese backtranslated data helped to increase performance tremendously (+18.8 BLEU). As expected, the zero-shot baseline performs the worst, with even greater score disparities.

Despite this significant improvement compared to our other techniques, the BLEU filtering approach fails to outperform the Apertium baseline. Apertium performs significantly better, with approximately 9 points difference in BLEU and over 5 points in chrF++.

## 6 Conclusions

This work was conducted within the constrained category of the shared task *Translation into Low-Resource Languages of Spain* at WMT24. It introduces a pipeline for filtering low-quality alignments in parallel corpora and subsequently fine-tuning translation models to assess the noise robustness of Neural Machine Translation. The paper details the data collection and curation processes for the three target languages selected for this task (Asturian, Aragonese and Aranese), with a particular focus on fine-tuning models for Spanish to Asturian under varying levels of noise and generalizing the results to the other two language pairs.

The initial phase involved curating the OPUS corpora for the Asturian-Spanish pair. This pipeline included **1)** cleaning unsupported characters and blank spaces, **2)** filtering out sentence pairs that were not in Spanish or one of the target languages using Idiomata Cognitor, and **3)** generating translations with Apertium to determine alignment quality of the sentence pairs and establishing four different BLEU thresholds for filtering. After observing that a BLEU threshold of 15 yielded the best performance, we incorporated backtranslated PILAR data into the filtered OPUS corpora. Part of step **3** was omitted for Aranese and Aragonese due to computational constraints.

Despite these filtering approaches resulting in the loss of some significant portions of the available corpora, we observed that the fine-tuned models effectively leveraged prior knowledge from the chosen related languages (Galician and Catalan).

- Our best performing fine-tuned model for Asturian outperformed the baseline Spanish-Galician model by 12.15 BLEU points.

- Our best performing fine-tuned model for Aragonese outperformed the baseline Spanish-Catalan model by 33.15 BLEU points.

- Our best performing fine-tuned model for Aranese outperformed the baseline Spanish-Catalan model by 23.29 BLEU points.

The results for our best fine-tuned Asturian model were relatively strong, achieving competitive scores compared to Apertium. Although the same approach was applied to Aranese and Aragonese, it did not surpass the Apertium baseline by a significant margin.

Overall, we demonstrated that 1) filtering out low-quality translations from a noisy parallel dataset improves fine-tuning results and yields faster training times, and 2) results for Asturian can reach baseline levels with a smaller, cleaner and more computationally efficient corpus, suggesting that the selected models can handle noise only to a certain degree. However, we cannot assert that this approach is effective for Aranese and Aragonese, as the results for these languages fall short of the rule-based baseline.

## 7 Limitations

The scope of this work is mainly limited by computational resources. The HiTZ Basque Center for Language Technology kindly allowed the authors access their resources, but understandably, priority was given to projects more closely related to their main research focus at the time. This led us to 1) implement our own BLEU score filter and dispense with newer, more accurate sentence alignment algorithms, 2) generalize the best BLEU score threshold in Asturian to the other two languages, Aranese and Aragonese.

One potential improvement to our approach would be the application of curriculum learning, where initial fine-tuning is performed on large synthetic data, followed by further fine-tuning on high-quality parallel data.

## 8 Further Work

Future work could address the limitations discussed in Section 7 by 1) exploring the outcomes of fine-tuning a language model on corpora cleaned using not just one, but various sentence alignment algorithms such us Bertalign (Liu and Zhu, 2022) or Vecalign (Thompson and Koehn, 2019), and 2) investigating whether Aranese and Aragonese tolerate different noise thresholds compared to Asturian. Additionally, future research might:

- estimate the amount of KWh required to fine-tune different amounts of corpora,

- examine whether data augmentation through backtranslation of additional OPUS corpora could enhance performance, as this is permitted in the constrained category,

- explore whether a tokenizer trained on a larger corpus and specialized in Asturian, Aranese, and Aragonese could improve results.

## Ethics Statement

The authors of this study adhered to the principles outlined in the ACL Ethics Policy and the ACM Code of Ethics. Our goal is for this research to benefit society by exploring language transferability among three low-resource languages, thus advancing machine translation techniques for underrepresented languages. All data used in this study was sourced from institutional and legal channels, explicitly approved and aligned with the original guidelines of the constrained submission for the shared task on *Translation into Low-Resource Languages of Spain* at WMT24. Throughout the research process, we prioritized transparency and fairness. We conducted honest and reliable evaluations and clearly communicated the limitations of our methods.

## Acknowledgements

# References

X. LL. GARCÍA ARIAS. 2002. Breve reseña sobre la lengua asturiana. *Informe sobre la llingua asturiana*, page 15.

Mikel L. Forcada and Francis M. Tyers. 2016. Apertium: a free/open source platform for machine translation and basic language technology. In *Proceedings of the 19th Annual Conference of the European Association for Machine Translation: Projects/Products*, Riga, Latvia. Baltic Journal of Modern Computing.

Alexander Fraser. 2020. Findings of the wmt 2020 shared tasks in unsupervised mt and very low resource supervised mt. In *Proceedings of the Fifth Conference on Machine Translation*, pages 765–771.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024a. Idiomata cognitor.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024b. Pilar.

Alina Karakanta, Jon Dehdari, and Josef van Genabith. 2018. Neural machine translation for low-resource languages without parallel corpora. *Machine Translation*, 32:167–189.

Wei-Jen Ko, Ahmed El-Kishky, Adithya Renduchintala, Vishrav Chaudhary, Naman Goyal, Francisco Guzmán, Pascale Fung, Philipp Koehn, and Mona Diab. 2021. Adapting high-resource nmt models to translate low-resource related languages without parallel data. *arXiv preprint arXiv:2105.15071*.

Sachin Kumar, Antonios Anastasopoulos, Shuly Wintner, and Yulia Tsvetkov. 2021. Machine translation into low-resource language varieties. *arXiv preprint arXiv:2106.06797*.

Jindřich Libovickỳ and Alexander Fraser. 2021. Findings of the wmt 2021 shared tasks in unsupervised mt and very low resource supervised mt. In *Proceedings of the Sixth Conference on Machine Translation*, pages 726–732.

Lei Liu and Min Zhu. 2022. Bertalign: Improved word embedding-based sentence alignment for Chinese–English parallel corpora of literary texts. *Digital Scholarship in the Humanities*, 38(2):621–634.

Cristian Marco Villanueva. 2012. Lengua aragonesa: Historia y situación actual.

Jesús Neira Martínez. 1982. La desaparición del romance navarro y el proceso de castellanización. *Revista española de lingüística*, 12(2):267–280.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Antoni Oliver, Mercè Vàzquez, Marta Coll-Florit, Sergi Álvarez, Víctor Suárez, Claudi Aventín-Boya, Cristina Valdés, Mar Font, and Alejandro Pardos. 2023. Tan-ibe: Neural machine translation for the romance languages of the iberian peninsula. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 495–496.

Surangika Ranathunga, En-Shiun Annie Lee, Marjana Prifti Skenduli, Ravi Shekhar, Mehreen Alam, and Rishemjit Kaur. 2023. Neural machine translation for low-resource languages: A survey. *ACM Comput. Surv.*, 55(11).

Cecilio Lapresta Rey and Ángel Huguet Canalís. 2006. Identidad colectiva y lengua en contextos pluriculturales y plurilingües. el caso del valle de arán (lleida. españa). *Revista internacional de sociología*, 64(45):83–115.

R Sennrich, B Haddow, and A Birch. 2015. Improving neural machine translation models with monolingual data. arxiv 2015. *arXiv preprint arXiv:1511.06709*.

Rico Sennrich and Martin Volk. 2010. MT-based sentence alignment for OCR-generated parallel texts. In *Proceedings of the 9th Conference of the Association for Machine Translation in the Americas: Research Papers*, Denver, Colorado, USA. Association for Machine Translation in the Americas.

Brian Thompson and Philipp Koehn. 2019. Vecalign: Improved sentence alignment in linear time and space. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1342–1348, Hong Kong, China. Association for Computational Linguistics.

Jörg Tiedemann. 2009. *News from OPUS - A Collection of Multilingual Parallel Corpora with Tools and Interfaces*, volume V, pages 237–248.

Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-MT – building open translation services for the world. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 479–480, Lisboa, Portugal. European Association for Machine Translation.

Marion Weller-Di Marco and Alexander Fraser. 2022. Findings of the wmt 2022 shared tasks in unsupervised mt and very low resource supervised mt. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 801–805.

# Training and Fine-Tuning NMT Models for Low-Resource Languages using Apertium-Based Synthetic Corpora

**Aleix Sant**[1]**, Daniel Bardanca Outeiriño**[2]**, José Ramom Pichel Campos**[2]
**Francesca De Luca Fornaciari**[1]**, Carlos Escolano**[1]**, Javier García Gilabert**[1]
**Pablo Gamallo Otero**[2]**, Audrey Mash**[1]**, Xixian Liao**[1]**, Maite Melero**[1]

[1] Barcelona Supercomputing Center (BSC),
[2] Centro Singular de Investigación en Tecnoloxías Intelixentes (CiTIUS)
Universidade de Santiago de Compostela

{aleix.santsavall, francesca.delucafornaciari, carlos.escolano
javier.garcia1, audrey.mash, xixian.liao, maite.melero}@bsc.es
{daniel.bardanca, jramon.pichel, pablo.gamallo}@usc.gal

## Abstract

In this paper, we present the two strategies employed for the WMT24 Shared Task on Translation into Low-Resource Languages of Spain. We participated in the language pairs of Spanish-to-Aragonese, Spanish-to-Aranese, and Spanish-to-Asturian, developing neural-based translation systems and moving away from rule-based approaches for these language directions. To create these models, two distinct strategies were employed. The first strategy involved a thorough cleaning process and curation of the limited provided data, followed by fine-tuning the multilingual NLLB-200-600M model (Constrained Submission). The other strategy involved training a transformer from scratch using a vast amount of synthetic data (Open Submission). Both approaches relied on generated synthetic data and resulted in high ChrF and BLEU scores. However, given the characteristics of the task, the strategy used in the Constrained Submission resulted in higher scores that surpassed the baselines across the three translation directions, whereas the strategy employed in the Open Submission yielded slightly lower scores than the highest baseline.

## 1 Introduction

This article presents the work done by the ILE-NIA team, which includes researchers from the Barcelona Supercomputing Center and Proxecto Nós (CiTIUS - Universidade de Santiago de Compostela), for the WMT24 Shared Task on Translation into Low-Resource Languages of Spain[1]. Our participation covered three translation directions: Spanish-to-Aragonese, Spanish-to-Aranese, and Spanish-to-Asturian, all of which are Romance languages.

---

[1] https://www2.statmt.org/wmt24/romance-task. html

Aragonese is spoken in several valleys of the Pyrenees in the autonomous community of Aragon. It is one of Europe's smallest language communities, with around 8,500 native speakers and 25,000 total speakers. According to UNESCO, Aragonese is an increasingly endangered language (Moseley, 2010).

Aranese is spoken in Vall d'Aran, in the northwest of Catalonia. It is the native language of this unique region, with approximately 5,090 native speakers. Aranese is a variant of Gascon, one of the main dialects of the Occitan language.

Asturian is the variant of Astur-Leonese spoken in the autonomous community of Asturias, in northern Spain. Currently, around 250,000 people have the ability to understand, speak, read, and write Asturian, representing the 25% of this autonomous community.

For this Shared Task, we participated in two types of submissions. In the **Constrained Submission**, we were allowed to use the specified resources, such as corpora, dictionaries, Apertium-based systems, and documents defining the contemporary orthographic conventions for each language. Regarding the models, we could use publicly available models, provided they did not exceed 1 billion parameters. To meet these requirements, we collected and generated synthetic data from the available resources, built and applied a comprehensive cleaning pipeline to preprocess the data, and fine-tuned three separate NLLB-200-600M (Costa-jussà et al., 2022) on the corresponding Spanish-to-Aragonese, Spanish-to-Aranese, and Spanish-to-Asturian translation directions. BSC researchers conducted the experiments for this submission.

For the **Open Submission**, we were allowed to use any publicly available resources, including corpora and models of any size, as long as the resulting

925

outputs were made available. For this submission, we chose to generate large amounts of synthetic corpora from data released by ILENIA and the Proxecto Nós using Apertium, a rule-based translator (Khanna et al., 2021). We then trained three models based on the transformer architecture (Vaswani et al., 2017) from scratch using OpenNMT-py 3.2[2] (Klein et al., 2018). The researchers from Proxecto Nós were responsible for this second submission.

## 2 Data

### 2.1 Data Collection and Synthetic Creation

#### 2.1.1 Constrained Submission

**OPUS** The organizers created the Aragonese and Aranese FLORES+ dev and devtest sets using Apertium, translating the corresponding Spanish texts from the FLORES-200 multilingual dataset (NLLB Team, 2022) into these languages. We consider this to be a limitation of the Shared Task since the reference test sets are biased towards Apertium-generated data. Nevertheless, in order to achieve the highest possible score on our submission for the Shared Task, we decided to use Apertium to generate synthetic Aragonese and Aranese translations from Spanish monolingual data, instead of directly using the parallel data provided by OPUS (Nygaard and Tiedemann, 2003). Specifically, we used the Spanish side of the es-arg and es-oc parallel corpora from OPUS to generate synthetic data. For Aragonese, we used the GNOME, Ubuntu, Wikimatrix, and Wikimedia corpora, and for Aranese, we used these same corpora in addition to Kde4 and NLLB.

For Asturian, we did not generate any synthetic data from OPUS since the Asturian FLORES+ dev and devtest sets were simply enhanced versions of the original FLORES-200. Instead, we downloaded the following es-ast parallel data: GNOME, Kde4, NLLB, Tatoeba, Ubuntu, and Wikimedia.

**PILAR**[3] We generated synthetic Spanish translations from the monolingual data provided by the organizers in the three respective Romance languages (Galiano-Jiménez et al., 2024b). For Aragonese and Aranese, we used Apertium, while for Asturian, we employed NLLB-200-600M, which fell within the submission's limits. Given the similarity between Aranese and Aragonese to Catalan, we explored whether cascading through Catalan could enhance translation quality. In machine translation, cascading refers to the sequential use of multiple translation systems to improve overall translation accuracy. As demonstrated in Table 1, this method yielded higher scores when translating from Aranese. Consequently, we generated synthetic Spanish translations from Aranese by cascading through Catalan, while for Aragonese and Asturian, we produced the translations directly into Spanish.

| | Aragonese → Spanish | | Aranese → Spanish | |
|---|---|---|---|---|
| | **ChrF** | **BLEU** | **ChrF** | **BLEU** |
| Direct | **80.93** | **66.09** | 68.79 | 45.02 |
| Cascade | 79.41 | 63.00 | **69.67** | **47.29** |

Table 1: Scores obtained with and without cascading through Catalan for FLORES+ dev test.

**Provided PDFs** We extracted monolingual data in Aragonese and Aranese from the provided `ortografia-aragones.pdf` and `DICCIONARI-DER-ARANÉS.pdf` respectively. After the text extraction, we semi-automatically post-processed the data to obtain a structured and clean corpus. Then we generated the corresponding Spanish translations using Apertium following the same method described in the previous paragraph.

**FLORES+** It consists of an extension of FLORES-200 (NLLB Team, 2022), a multilingual English-centric machine translation dataset involving 200 languages, that includes Aragonese, Aranese, and an improved version of Asturian. The FLORES+ dev set for each language served as the validation set during the training phase to optimize our MT engines, while the devtest set was used to evaluate participants' models in the competition. After the final submissions, the devtest set was released, allowing us to obtain scores for the baseline models in this additional set.

#### 2.1.2 Open Submission

For this submission, we used Apertium to generate synthetic data. We created synthetic datasets using a parallel corpus of 30M Galician-Spanish sentence pairs. We only kept the Spanish side of the corpus as the source language data, and then used Apertium to translate it into Aragonese and Asturian, resulting in two 30M sentence parallel corpora (Spanish-Aragonese and Spanish-Asturian).

---

In the case of Aranese, we used a high-quality 30M sentence Spanish-Catalan parallel dataset. We translated the Catalan side into Aranese using Apertium, creating a Spanish-Aranese corpus.

## 2.2 Data Preprocessing and Cleaning

### 2.2.1 Constrained Submission

For this submission, we dedicated substantial effort to cleaning, curating, and normalizing the provided data. We designed a comprehensive cleaning pipeline that processed all the parallel data described in the previous section, resulting in well-structured parallel corpora for the Spanish-Aragonese, Spanish-Aranese, and Spanish-Asturian language pairs.

Following the automatic cleaning, we curated the resulting data to ensure it aligned to the orthographic standards outlined in the task statement and matched the characteristics of the corresponding FLORES+ sets for each language.

**Blank Spaces, Hard- and Soft-Duplicates Removal**    The initial step involved removing any unnecessary blank spaces and exact duplicates within the corpus. Then, NLPDedup[4] was used to remove near duplicates.

**Idiomata Cognitor**[5]    This language identifier (Galiano-Jiménez et al., 2024a), specifically designed for certain Romance languages, was employed to accurately determine the languages of each data pair and exclude pairs with sentences belonging to other languages. This method ensures that the translator model is trained on appropriate data.

**Perl Corpus Cleaner**    We employed Moses (Koehn et al., 2007) preprocessing script `clean-corpus-n.perl` to further clean the parallel corpus. It eliminates sentences containing more than 150 tokens and discards sentence pairs with a length ratio exceeding 3.

**Linguistic Data Normalization**    Since the released corpora included text in various orthographies, and both the FLORES+ dev and devtest sets adhered to the current standards endorsed by their respective language academies, we ensured that our training data conformed to these established

norms through a normalization process. This process was carried out semi-automatically: incorrect patterns in the data were detected and replaced with the correct ones according to the relevant linguistic rules. This normalization was primarily applied to the Aragonese and Aranese monolingual data. For example, in Aragonese, we encountered different types of definite articles used interchangeably, such as "o"/"lo", "a"/"la", "os"/"los", and "as"/"las", as well as various ways to write the word "university", including "unibersidad", "unibersidá", and "universidat". In the case of definite articles, all forms needed to be standardized to "lo"/"la"/"los"/"las", except when following a word ending in 'n', where the forms are "o"/"a"/"os"/"as". For the term "university", the officially accepted word is "universidat".

**Data Curation**    Using the FLORES+ dev set as our reference, we further examined the parallel data to identify misleading translations in the training data. This curation process involved both semi-automatic and manual methods, primarily focusing on the word level. For example, Apertium often leaves unknown words unchanged in the target sentence or produces incorrect translations due to insufficient contextual understanding of the source sentence. These are common behaviors of rule-based translators. We aimed to detect these issues and correct these translation errors.

According to Table 2, a large number of sentence pairs are discarded, mainly due to the high volume of duplicates in the three corpora, with Aranese and, particularly, Asturian exhibiting the highest number of duplicates.

| | Aragonese | Aranese | Asturian |
|---|---|---|---|
| Original | 74,014 | 1,336,229 | 6,603,733 |
| Filtered | 47,521 | 407,397 | 704,933 |

Table 2: Parallel corpus statistics per target language. `Original` refers to all the collected data pairs before going through the pipeline. `Filtered` refers to the number of pairs resulting from the data cleaning pipeline.

**LaBSE scoring**    To evaluate the quality of translations in the parallel datasets, we used a sentence embedding model. Specifically, we employed LaBSE (Feng et al., 2022) to generate embeddings for both source and target sentences and calculated the co-

| | | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | $\sigma$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Aragonese** | **ChrF** | 84.6 | 84.59 | 84.59 | 84.61 | 84.59 | **84.63** | 84.58 | 84.62 | 84.5 | 0.6752 |
| | **BLEU** | 71.05 | 71.05 | 71 | 71.09 | 71.03 | **71.14** | 70.99 | 71.05 | 70.84 | 0.0832 |
| **Aranese** | **ChrF** | 75.95 | **76.04** | 76.02 | 75.99 | 75.97 | 75.97 | 75.91 | 75.84 | 75.32 | 0.2218 |
| | **BLEU** | 55.39 | 55.5 | **55.57** | 55.43 | 55.39 | 55.43 | 55.3 | 55.2 | 54.39 | 0.3535 |
| **Asturian** | **ChrF** | 52.25 | 52.25 | **52.26** | 52.25 | 52.21 | 52.19 | 52.21 | 52.22 | 52.19 | 0.1897 |
| | **BLEU** | **19.35** | 19.34 | 19.33 | 19.32 | 19.28 | 19.23 | 19.26 | 19.27 | 19.24 | 0.1737 |

Table 3: Scores for different LaBSE thresholds on FLORES+ dev set.

sine similarity score between them. Following the approach outlined in (Garcia Gilabert et al., 2024), we obtained scores across various thresholds of cosine similarity to determine the most suitable training dataset for fine-tuning (Table 3). However, the low standard deviation among the results indicates no considerable differences in performance between the sets. We selected a threshold of 0.6 for Aragonese, a threshold of 0.3 for Aranese, and a threshold of 0.1 for Asturian, prioritizing a higher BLEU over a ChrF.

### 2.2.2 Open Submission

No extra preprocessing or cleaning was performed on the synthetic corpora generated with Apertium for this submission. The source data had previously been processed before publication with their own pipeline [6]. This preprossessing includes fixing encoding issues, deduplication, perplexity filtering and language recognition.

## 3 Methodology

### 3.1 Baselines

In Table 4, ChrF and BLEU scores on the FLORES+ dev set for the state-of-the-art models including the directions of interest are shown. Except for the NLLB models, which employ deep learning, all the other evaluated engines are rule-based.

### 3.2 Constrained Submission: Fine-tuning

### 3.2.1 Model

NLLB-200-600M is the smallest model from the NLLB family of multilingual machine translation models. It is a dense transformer model distilled from the pre-trained NLLB-200, a 54.5B sparsely gated mixture-of-experts model, designed to support translations between 202 languages, including many low-resource languages. Therefore, it

incorporates substantial cross-lingua knowledge, making it suitable for further fine-tuning to other languages. It has a vocabulary size of 256k tokens, plus additional tokens for the language tags corresponding to all the languages supported by the model. With respect to our languages of interest, it just handles Asturian. Occitan is also in the list of languages, but not the Aranese variant.

### 3.2.2 Fine-tuning

For this approach, we fully fine-tuned three separate NLLB-200-600M models for Spanish-to-Aranese, Spanish-to-Aragonese, and Spanish-to-Asturian, leveraging the cross-lingua knowledge NLLB possesses.

**Adding New Language Tags** Among the languages of interest, NLLB only supports Asturian natively. To extend NLLB's translation capabilities to translate to Aragonese and to Aranese, we incorporated new tokens referring to their respective language tags (arn_Latn and arg_Latn), since these languages are not present in NLLB[7]. Language tags enable NLLB to identify the source and target languages for translation. Adding new language tags implies extending the embedding matrix with additional embeddings. These new embeddings were initialized using the embeddings of other language tags already supported by NLLB, which were linguistically close to our target language. Specifically, we used the spa_Latn embedding for Aragonese and the oci_Latn embedding for Aranese. Finally, we retrained the embedding matrix during the fine-tuning to enable correct Spanish-to-Aragonese and Spanish-to-Aranese translation. For Asturian, since NLLB already supports translation to this language, our objective was simply to improve the model's performance

---

[7] Aranese is a variant of Occitan, but due to observed differences in the test sets, we treated Aranese as a distinct language with its own language tag.

| | | Spanish → Aragonese | | Spanish → Aranese | | Spanish → Asturian | |
|---|---|---|---|---|---|---|---|
| | | **ChrF** | **BLEU** | **ChrF** | **BLEU** | **ChrF** | **BLEU** |
| Baselines | Apertium | 82 | 65.34 | 72.63 | 48.96 | 50.57 | 16.66 |
| | Traduze | 69.51 | 37.43 | - | - | - | - |
| | Softcatala | 73.97 | 50.21 | 58.61 | 34.43 | - | - |
| | Eslema | - | - | - | - | 50.77 | 17.3 |
| | NLLB-600M | - | - | - | - | 49.72 | 17.23 |
| | NLLB-1.3B | - | - | - | - | 50.04 | 17.44 |
| | NLLB-3.3B | - | - | - | - | 50.15 | 17.96 |
| Constrained Submission | NLLB-600M fine-tuned | 84.63 | 71.14 | 76.04 | 55.5 | 52.26 | 19.33 |
| Open Submission | Transformer from scratch | 81.35 | 63.95 | 71.48 | 45.92 | 50.37 | 16.86 |

Table 4: Evaluations computed on the FLORES+ dev set.

through fine-tuning.

**Treatment of "«" and "»" tokens**   Given that "«" and "»" symbols are not in the NLLB vocabulary and they were present in the reference test sets, we decided to preprocess the training data by replacing "«" and "»" with "<<" and ">>" (as these tokens are in the vocabulary) and then revert this replacement after the model's inference.

### 3.2.3   Inference experiments

The generated translation is restricted to a length of 512 tokens. Testing on FLORES+ dev, we conducted a grid search on the top-performing model obtained from the LaBSE scoring step. We experimented with various beam sizes ($B$) and repetition penalty terms ($\beta$). We tried all combinations between $B = [3, 5, 10]$ and $\beta = [1, 3, 4]$. Nevertheless, no significant differences in performance were observed for these languages, so we ended up using the same hyperparameters employed during training: $B = 5$ and $\beta = 1$. For detailed results, see Appendix.

### 3.2.4   Configurations

In the fine-tuning, we used the AdamW optimizer (Loshchilov and Hutter, 2019) with $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 10^{-6}$, and $\lambda = 0.001$. The learning rate was set to $3 \times 10^{-4}$. We applied an inverse square root scheduler with 15,000 warmup steps.

For Aragonese, the batch size was set to 16, the gradient accumulation steps to 8, and the model was trained for 15 epochs. For Aranese and Asturian, the batch size was set to 8, the gradient ac-

cumulation steps to 4, and the models were trained for 10 epochs. All models were fine-tuned using the Transformers[8] library on H100 GPUs. Every 1,000 training steps, the ChrF score was computed on the FLORES+ dev set, and the model checkpoints were saved when the score improved.

### 3.3   Open Submission: Training

#### 3.3.1   Model

We trained three transformer models from scratch using OpenNMT-py 3.2, each with its own BPE vocabulary. The vocabulary size for each model was set to 20,000 units, based on previous internal research investigating the impact of vocabulary size on BLEU scores (Outeirinho et al., 2024).

#### 3.3.2   Configurations

All three models were trained on a single A100 GPU using the AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.9998$, $\epsilon = 10^{-8}$, and learning rate $5 \times 10^{-4}$. The batch size was set to 2048 sentences, with a maximum length of 150 tokens per sentence. All models were trained for a maximum of 10 epochs.

### 3.4   Evaluation

To evaluate the performance of our models during the development phase, we used the FLORES+ dev set, which contains 997 general domain sentences. The results obtained on this test set for the two developed strategies can be seen in Table 4.

---

[8] https://huggingface.co/docs/transformers/

|  |  | Spanish → Aragonese | | Spanish → Aranese | | Spanish → Asturian | |
|---|---|---|---|---|---|---|---|
|  |  | ChrF | BLEU | ChrF | BLEU | ChrF | BLEU |
| Baselines | Apertium | 79.31 | 61.11 | 49.42 | 28.85 | 50.84 | 16.99 |
|  | Traduze | 67.66 | 35.47 | - | - | - | - |
|  | Softcatala | 71.99 | 47.08 | 48.29 | 26.07 | - | - |
|  | Eslema | - | - | - | - | 50.91 | 17.17 |
|  | NLLB-600M | - | - | - | - | 49.05 | 16.21 |
|  | NLLB-1.3B | - | - | - | - | 49.71 | 16.54 |
|  | NLLB-3.3B | - | - | - | - | 50.03 | 17.09 |
| Constrained Submission | NLLB-600M fine-tuned | 79.88 | 62.32 | 50.05 | 30.12 | 52.14 | 18.43 |
| Open Submission | Transformer from scratch | 78.61 | 59.76 | 48.84 | 27.31 | 50.54 | 16.68 |

Table 5: Evaluations computed on the FLORES+ devtest set.

### 3.4.1 Constrained Submission

Following this strategy, we surpassed the baselines across the three languages, achieving better performance than the state-of-the-art. We observed an increase of +2.63 in ChrF and +5.8 in BLEU for Aragonese (compared to Apertium), +3.41 in ChrF and +6.54 in BLEU for Aranese (compared to Apertium), and +1.49 in ChrF and +1.37 in BLEU for Asturian (compared to Eslema and NLLB-200-3.3B, respectively). These results suggest that both the thorough curation of data and the cross-lingual knowledge possessed by NLLB contributed to these improvements.

### 3.4.2 Open Submission

Leveraging transformer models has led to results that are only slightly behind the best rule-based approaches. However, the significance of these results lies in the fact that they enable the community to access and develop neural models that perform competently in a relatively short time compared to developing a new rule-based system from the ground up. These neural models, particularly transformers, offer new possibilities, such as the ability to learn from limited data and improved scalability, which can help prevent languages with fewer speakers from being marginalized in the online world.

### 4 Results

The FLORES+ devtest set, containing 1,012 sentences, was used to evaluate and rank the participants' models. Once the competition ended, the organizers made the FLORES+ devtest set public.

To further expand the evaluation of our new models, we also obtained scores for the baselines using this test set. Consult Table 5 for all the scores. We see the same trend as with the FLORES+ dev set. Using the fine-tuned version of NLLB-200-600M on the cleaned data, we surpass all the baseline models for the three languages, whereas training the models with OpenNMT-py 3.2 lags behind. Specifically, we enhance the scores by +0.57 in ChrF and +1.21 in BLEU for Aragonese (compared to Apertium), +0.63 in ChrF and +1.27 in BLEU for Aranese (compared to Apertium), and +1.23 in ChrF and +1.26 in BLEU for Asturian (compared to Eslema). At the time of writing this paper, the final ranking scores were not available, so no mention of our final positions in the competition is included in this paper.

### 5 Discussion

Compared to traditional rule-based translation systems, neural models offer greater flexibility, scalability, and adaptability, making them the state-of-the-art in Machine Translation. Hence, our work in developing neural systems for Aragonese, Aranese and Asturian represents an advance in the preservation and promotion of the use of these languages. It also allows the research community to use our models for further advancements in language technology, linguistic research, and the development of more sophisticated and accurate translation systems.

# 6 Conclusions

This paper summarizes the work done by the ILE-NIA team for the Shared Task on Translation into Low-Resource Languages of Spain. By participating in this public competition, we have contributed to the creation and improvement of NMT models for Aragonese, Aranese, and Asturian - three minority languages of Spain. Prior to this task, no NMT models were available for Spanish-to-Aragonese and Spanish-to-Aranese translation.

We presented a Constrained and an Open Submission, each employing different approaches. For the Constrained Submission, adhering to data and model restrictions, we fine-tuned the NLLB-200-600M model, with considerable effort devoted to data cleaning and curation. For the Open Submission, we generated a large amount of synthetic data using Apertium, a rule-based MT system, and used it to train a transformer-based model from scratch.

Results on both FLORES+ dev and devtest sets across the three language directions show that the first strategy achieves better performance and improves translation quality compared to the baselines, whereas the second strategy lags slightly behind the best baseline models.

## Acknowledgements

## References

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024a. Idiomata cognitor.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024b. Pilar.

Javier Garcia Gilabert, Aleix Sant, Carlos Escolano, Francesca De Luca Fornaciari, Audrey Mash, and Maite Melero. 2024. BSC submission to the AmericasNLP 2024 shared task. In *Proceedings of the 4th Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP 2024)*, pages 143–149, Mexico City, Mexico. Association for Computational Linguistics.

Tanmai Khanna, Jonathan North Washington, Francis M. Tyers, Sevilay Bayatli, Daniel G. Swanson, Flammie A. Pirinen, Irene Tang, and Hèctor Alòs i Font. 2021. Recent advances in apertium, a free/open-source rule-based machine translation platform for low-resource languages. *Mach. Transl.*, 35(4):475–502.

Guillaume Klein, Yoon Kim, Yuntian Deng, Vincent Nguyen, Jean Senellart, and Alexander M. Rush. 2018. Opennmt: Neural machine translation toolkit.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

C. Moseley. 2010. *Atlas of the World's Languages in Danger*. Memory of peoples Series. UNESCO.

James Cross Onur Çelebi Maha Elbayad Kenneth Heafield Kevin Heffernan Elahe Kalbassi Janice Lam Daniel Licht Jean Maillard Anna Sun Skyler Wang Guillaume Wenzek Al Youngblood Bapi Akula Loic Barrault Gabriel Mejia Gonzalez Prangthip Hansanti John Hoffman Semarley Jarrett Kaushik Ram Sadagopan Dirk Rowe Shannon Spruit Chau Tran Pierre Andrews Necip Fazil Ayan Shruti Bhosale Sergey Edunov Angela Fan Cynthia Gao Vedanuj Goswami Francisco Guzmán Philipp Koehn Alexandre Mourachko Christophe Ropers Safiyyah Saleem Holger Schwenk Jeff Wang NLLB Team, Marta R. Costa-jussà. 2022. No language left behind: Scaling human-centered machine translation.

Lars Nygaard and Jörg Tiedemann. 2003. Opus—an open source parallel corpus. In *Proceedings of the 13th Nordic Conference on Computational Linguistics*.

Daniel Bardanca Outeirinho, Pablo Gamallo Otero, Iria de Dios-Flores, and José Ramom Pichel Campos. 2024. Exploring the effects of vocabulary size in neural machine translation: Galician as a target language. In *Proceedings of the 16th International Conference on Computational Processing of Portuguese - Vol. 1*, pages 600–604, Santiago de Compostela, Galicia/Spain. Association for Computational Lingustics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

# Appendices

## A   Spanish-to-Aragonese

| | | Beam Width | | |
|---|---|---|---|---|
| | | **3** | **5** | **10** |
| **Rep. Pen.** | **1** | 84.63 | 84.63 | 84.63 |
| | **3** | 84.64 | 84.64 | 84.66 |
| | **4** | 86.63 | 84.64 | 84.63 |

Table 6: ChrF scores obtained using grid search in inference.

| | | Beam Width | | |
|---|---|---|---|---|
| | | **3** | **5** | **10** |
| **Rep. Pen.** | **1** | 71.12 | 71.14 | 71.14 |
| | **3** | 71.15 | 71.16 | 71.16 |
| | **4** | 71.15 | 71.16 | 71.15 |

Table 7: BLEU scores obtained using grid search in inference.

## B   Spanish-to-Aranese

| | | Beam Width | | |
|---|---|---|---|---|
| | | **3** | **5** | **10** |
| **Rep. Pen.** | **1** | 76.04 | 76.04 | 76.04 |
| | **3** | 76.04 | 76.04 | 76.04 |
| | **4** | 76.04 | 76.04 | 76.04 |

Table 8: ChrF scores obtained using grid search in inference.

| | | Beam Width | | |
|---|---|---|---|---|
| | | **3** | **5** | **10** |
| **Rep. Pen.** | **1** | 55.5 | 55.5 | 55.5 |
| | **3** | 55.6 | 55.6 | 55.6 |
| | **4** | 55.6 | 55.6 | 55.6 |

Table 9: BLEU scores obtained using grid search in inference.

## C   Spanish-to-Asturian

| | | Beam Width | | |
|---|---|---|---|---|
| | | **3** | **5** | **10** |
| **Rep. Pen.** | **1** | 52.28 | 52.26 | 52.23 |
| | **3** | 52.26 | 52.24 | 52.22 |
| | **4** | 52.25 | 52.25 | 52.23 |

Table 10: ChrF scores obtained using grid search in inference.

| | | Beam Width | | |
|---|---|---|---|---|
| | | **3** | **5** | **10** |
| **Rep. Pen.** | **1** | 19.34 | 19.33 | 19.21 |
| | **3** | 19.25 | 19.24 | 19.22 |
| | **4** | 19.21 | 19.24 | 19.21 |

Table 11: BLEU scores obtained using grid search in inference.

# Vicomtech@WMT 2024: Shared Task on Translation into Low-Resource Languages of Spain

**David Ponce**[*1,2] and **Harritxu Gete**[*1,2] and **Thierry Etchegoyhen**[1]

[1] Fundación Vicomtech, Basque Research and Technology Alliance (BRTA)
[2] University of the Basque Country UPV/EHU
{adponce,hgete,tetchegoyhen}@vicomtech.org

## Abstract

We describe Vicomtech's participation in the WMT 2024 Shared Task on translation into low-resource languages of Spain. We addressed all three languages of the task, namely Aragonese, Aranese and Asturian, in both constrained and open settings. Our work mainly centred on exploiting different types of corpora via data filtering, selection and combination methods, along with synthetic data generated with translation models based on rules, neural sequence-to-sequence or large language models. We improved or matched the best baselines in all three language pairs and present complementary results on additional test sets.

## 1 Introduction

Despite significant progress in Machine Translation (MT) in recent years, notably with the advent of Neural Machine Translation (NMT) approaches (Bahdanau et al., 2015; Vaswani et al., 2017), translation from and into low-resource languages remains a challenge.

Spain features a large variety of languages beyond Spanish, with varying degrees of MT support. Important quality gains have thus been achieved for the Basque language within the NMT framework (Etchegoyhen et al., 2018), with large public deployments of quality MT systems[1]. For Catalan, a romance language with closer proximity to Spanish, earlier NMT improved over rule-based (RMT) and statistical (SMT) models, although with performance losses on out-of-domain test sets (Costa-jussà, 2017); more recent work on translation between similar languages, that included Catalan-Spanish, showed a predominance of NMT approaches to the task (Akhbardeh et al., 2021).

In addition to the aforementioned languages, there are languages such as Aragonese, Aranese

and Asturian which could be viewed as extremely low-resourced in terms of MT technological support. For most, the main technology is still RBMT, based on the Apertium framework (Forcada and Tyers, 2016). The WMT 2024 shared task on translation into low-resourced languages of Spain addresses translation from Spanish into all three of these languages. In this work, we describe Vicomtech's participation in the shared task, where we submitted entries to both the constrained and open tracks.

In the remainder of this paper, we describe our approaches to improve MT performance for the three selected language pairs. We explored data selection, generation, and combination, comparing the use of different types of data to train end-to-end NMT models as well as fine-tuning pretrained multilingual MT models. In addition to typical parallel data curation, where we filtered the available parallel and comparable data according to sentence similarity, length differences and language identification, we also explored the generation of synthetic data along different lines. We notably compared the use of RBMT systems and large language models (LLM) to generate synthetic parallel datasets from the available monolingual data. The latter approach in particular showcased the potential of LLMs to create back-translations from the selected three low-resource Romance languages into the high-resource Spanish language.

## 2 Methodology

### 2.1 Parallel Data Curation

Despite the limited amount of data available for the languages addressed in this task, several crawled corpora were made available. However, after manually examining sampled of the data, they appeared to feature large amounts of noise, including poor alignments, language identification errors, or sentence pairs with empty information in one of the

---

languages. We therefore performed various types of filtering, described below.

**Language Identification.** We performed language identification with the Idiomata Cognitor tool[2], a Bayesian language identifier specialised on Romance Languages (Galiano-Jiménez et al., 2024a). We filtered all sentence pairs where the identified language on either side mismatched the expected language in the parallel dataset.

**Length Ratio.** We filtered all sentence pairs where the ratio of lengths, in terms of characters, was above a predefined threshold. Unless otherwise specified, we used a default ratio of 3.0. Our goal with this type of filtering was to remove obvious erroneous alignments rather than determine an optimal threshold in terms of length differences.

**Sentence Similarity.** We filtered all sentence pairs whose similarity score was below a predefined threshold. Similarity was computed as the cosine similarity of the sentence embeddings for each sentence pair. After preliminary experiments with different models, we opted for the all-MiniLM-L6-v2 model of the Transformers library[3], as it provided sufficient quality for the considered pairs, while also supporting sufficiently fast processing to run multiple filtering experiments. For each language pair, we assigned similarity scores to the parallel corpora after language and length filtering, manually examined samples of the data and determined a similarity threshold accordingly.

## 2.2 Synthetic Data Creation

For low-resource languages, parallel data are typically scarce and monolingual corpora are a rich source of complementary data. We aimed to explore different approaches to exploit this type of data, generating synthetic data by translating via RBMT systems, NMT models and LLMs (see Frontull and Moser (2024) for a similar approach). Depending on model availability and/or quality, we generated data to be used as either back-translations (BT) (Sennrich et al., 2016) from the low-resource languages into Spanish, or as forward-translations (FT) (Li and Specia, 2019) in the opposite translation direction. In either case, the synthetic data generated from monolingual data were used as parallel data to translate into the low-resource languages. Additionally, we used pivot machine-translation from Catalan to Spanish to complement the Spanish-Aranese parallel datasets, as described below.

**RBMT data (BT + FT).** As back-translations, we translated the available monolingual corpora in Aragonese and Aranese into Spanish with the corresponding Apertium systems.[4] As forward-translations, we generated synthetic data from Spanish into all three low-resource languages, since Apertium covered all three language pairs in that direction. Our goal in both the BT and FT cases was to evaluate the impact of data translated via transformation rules that tend to closely follow the structure of the original Spanish data.

**NMT data (BT).** We generated back-translations into Spanish for all three language pairs with baseline NMT models, either trained from scratch or pretrained and fine-tuned, on the curated parallel data, as described in Section 2.3. Considering the low volumes of clean parallel data and the relatively low quality of the baselines, we discarded the use of forward-translations in this scenario. Back-translations are more robust in this type of scenario, as the target language monolingual data are expected to be correct for the decoder to model and the noise in corresponding synthetic source data can be handled relatively well by NMT models in general. Our aim with NMT-based NMT data was to generate synthetic data of relatively fluid translations that would differ from, and could complement, RBMT translations.

**LLM data (BT).** We also leveraged a general-purpose language model in zero-shot fashion to generate back-translations, querying the model to translate from the low-resource language into Spanish. Our preliminary assessment on the three language pairs was that translation into the low-resource languages could not constitute a reasonable alternative, as most translations from Spanish into either low-resource language were of low quality, irrespective of the size of the selected model. However, in the reverse direction, in all three pairs translation quality was markedly better, indicating that the meaning of the text in the low-resource language could be properly captured by the model, while generating correct output in the high-resource Spanish language.

---

[2]https://github.com/transducens/idiomata_cognitor
[3]https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

[4]https://www.apertium.org/ Note that there was no readily available system for Asturian to Spanish.

**Pivot MT data.** Among the available corpora for the task were data in Catalan-Aranese (see Section 3.1), which could be exploited as well via pivot MT. To this end, we translated the available corpora from Catalan to Spanish with a high-quality in-house NMT model trained on OPUS parallel data (see Appendix C for further details).

## 2.3 Models & Training

**Models.** We trained two main types of models: Transformer-base encoder-decoder models trained from scratch on the available parallel, with or without complementary synthetic data, and a pretrained multilingual model fine-tuned with the same data, namely an NLLB-200-600M model (Costa-jussà et al., 2022). By opting for two parallel approaches, we aimed to evaluate the positive or negative impact of accessing pretrained multilingual knowledge on the task. With either type of model, we trained baseline variants on the curated parallel data, which were used used to generate back-translations, as described in Section 2.2. Both types of models were also used on the combined datasets to train final models, as our main aim was to contrast and compare the use of pre-trained multilingual knowledge vs. focused training on a specific low-resource language pair.

**Tagging.** To train the model variants, we performed several experiments around data tagging, which has been shown to be an efficient approach to training data discrimination, for back-translations (Caswell et al., 2019) or comparable data (Gete and Etchegoyhen, 2022), for instance. We used specific tags, prepended to each training instance, to indicate the type of data at hand, namely <BT> or <FT>. We aimed to investigate in particular whether data tags would be beneficial or detrimental in the case of low amounts of parallel data, combined with larger sets of synthetic data.

## 3 Experimental Setup

### 3.1 Corpora

For the constrained track, we selected the parallel corpora for Asturian, Aranese, and Aragonese from the PILAR collection (Galiano-Jiménez et al., 2024b), the monolingual WikiMedia data for Spanish and Asturian, the parallel data for Spanish-Asturian and Spanish-Occitan (with Occitan related to Aranese) from CCMatrix (Schwenk et al., 2021b) and WikiMatrix (Schwenk et al., 2021a) for Spanish-Aragonese, all of which were downloaded

| Corpus | Lang. | # Sent. | # Filt. | Constr. |
|---|---|---|---|---|
| PILAR | ast | 38.8K | - | ✓ |
| | arg | 84.7K | - | ✓ |
| | arn | 273.2K | - | ✓ |
| | cat-arn | 64.1K | - | ✓ |
| WikiMedia | es | 3.9M | 2.7M | ✓ |
| | ast | 65.7K | 65.6K | ✓ |
| CCMatrix | es-ast | 6.5M | 533.7K | ✓ |
| | es-oci | 925.5K | 55.5K / 8.9K | ✓ |
| WikiMatrix | es-arg | 33.7K | 19.2K / 13.7K | ✓ |
| WikiDump | ast | 3.2M | 2.1M | ✗ |
| | arg | 508.5K | 255.1K | ✗ |

Table 1: Corpora statistics. We indicate the number of initial (# Sent.) and filtered (# Filt.) sentences and corpus use in the constrained track (Constr.). x/y indicates filtering with similarity thresholds of 0.5 (x) and 0.7 (y).

from the OPUS repository (Tiedemann, 2012). We performed language identification to keep only the Aranese sentences from Occitan, and also translated the Catalan portion of the Catalan-Aranese dataset via pivot translation into Spanish. For the open track, we included Asturian and Aragonese monolingual data extracted from WikiDump[5], for additional back-translations.

Excepting the PILAR datasets, which were used as is, we filtered the contents from the parallel corpora using the methods described in Section 2.1. The similarity threshold was set at 0.7 after manually reviewing portions of the data. For Aragonese and Aranese, since significant portions of the datasets were discarded at this threshold, we also created an additional dataset with a 0.5 threshold. For the constrained task, we selected these larger, though noisier, datasets. For the open task, we opted for the smaller, higher-quality datasets, due to the greater availability of data.

As evaluation data, we selected the dev sets available in PILAR, as well as a filtered subset of 3,000 highquality sentence pairs from CCMatrix for Asturian, with a similarity threshold set at 0.9. The latter was created as all models consistently yielded significantly lower scores on the Asturian PILAR dev set, compared to the other language pairs, and Marian models trained with this development set struggled to converge effectively. We report results on the official development set throughout the paper, but discuss additional results on our own development set in Section 5.

Corpora statistics are summarised in Table 1.

---

[5]https://dumps.wikimedia.org/, accessed June 2024

## 3.2 Models

In this section, we describe the translation models we used for the task, including the baselines and the models trained on the selected data described in the previous section. Since we generated synthetic data for both forward (from Spanish) and backward (into Spanish) translation, we present each type of model in turn. Training details, including additional model characteristics and training hyper-parameters are described in Appendix A. All models were evaluated in terms of BLEU and chrF, computed with the sacreBLEU toolkit [6]. Statistical significance was computed via bootstrap resampling (Koehn, 2004) for all results. Best results, for $p < 0.05$, are indicated in bold in all tables.

### 3.2.1 Translation from Spanish

For translation from Spanish, we first assessed the quality of three baseline models not trained on any of the selected data: the rule-based Apertium for each language pair, as a reference MT system for these languages; the multilingual NLLB-200-distilled-600M model, pretrained on a broad range of languages including Asturian and Occitan, as a neural baseline under the constrained track limitation of pretrained models with fewer than one billion parameters; and the Llama3-8B instruction model (AI@Meta, 2024), as an experimental testbed for zero-shot LLM-based translation.

| Lang. | Model | Dev | |
|---|---|---|---|
| | | BLEU | chrF |
| ES→AST | Apertium | **17.1** | **50.7** |
| | NLLB | 14.3 | 44.2 |
| | Llama3 | 15.2 | 48.9 |
| ES→ARG | Apertium | **66.0** | **82.2** |
| | NLLB | 7.9 | 42.1 |
| | Llama3 | 30.4 | 64.5 |
| ES→ARN | Apertium | **38.0** | **60.0** |
| | NLLB | 8.5 | 39.2 |
| | Llama3 | 4.5 | 32.6 |

Table 2: Baseline model results on the development sets for translation from Spanish

Table 2 presents the results for each baseline model in this translation direction, in terms of BLEU (Papineni et al., 2002) and chrF (Popović, 2015). Apertium achieved the highest scores in all three language pairs, demonstrating the value of an RBMT approach for the selected languages. NLLB and Llama3 were notably both outperformed by

large margins on ES-ARN; the former performed equally poorly on ES-ARG but the latter achieved a more reasonable performance of 30.4 BLEU points in this case, still far from the scores obtained by the Apertium baseline. The only language pair where all three models achieved relatively similar low scores was ES-AST, which might be due to the specifics of this development set (see Section 5 for further discussion).

Considering these results, we used Apertium to generate forward synthetic data for all ES→XX translation pairs. To prepare the final models, all related to translation from Spanish in the task, we used two types of approaches: fine-tuning the NLLB model on the selected data and training from scratch a Transformer-base model with 6 encoder layers and 6 decoder layers, trained with the Marian NMT toolkit (Junczys-Dowmunt et al., 2018).

### 3.2.2 Translation into Spanish

| Model | Aranese | Aragonese | Asturian |
|---|---|---|---|
| Apertium | **34.8** | 66.2 | - |
| NLLB | 31.0 | 55.6 | 64.1 |
| Llama3 | 33.1 | 64.3 | 71.5 |
| Marian | **34.0** | **69.6** | **86.5** |

Table 3: BLEU scores for translation into Spanish on the PILAR development sets for Aranese and Aragonese, and on a custom development set for Asturian.

Translation into Spanish was performed to generate synthetic back-translations. For this task, we used the three baseline approaches described in the previous section (except for AST-ES with Apertium, as it is not currently supported) and trained an additional XX→ES Marian model on the selected parallel and forward-translation data.

Table 3 presents the BLEU scores for these models on the task-provided development sets for Aranese and Aragonese, and on our custom development set for Asturian. For Aranese, there was no statistically significant difference between the Marian and Apertium models, both outperforming NLLB and Llama3; in Aragonese, Marian outperformed all other models, with NLLB performing notably worse; in Asturian, it again significantly outperformed both NLLB and Llama3. Considering these results, we selected the Marian model to generate all back-translations. Additionally, since forward-translations were all generated using Apertium, the incorporation of a neural model could add more variety to the synthetic data.

| Lang. | Model | Data | # Sent. | Source | Dev BLEU | Dev chrF | Test BLEU | Test chrF |
|---|---|---|---|---|---|---|---|---|
| ast | Apertium | - | - | - | 17.1 | 50.7 | 17.0 | **50.8** |
| | NLLB | Parallel | 533.7K | CCMatrix | **18.1** | **51.3** | **17.6** | **51.2** |
| | | FT | - | - | | | | |
| | | BT | 638.1K | PILAR+CCMatrix+WikiMedia | | | | |
| arg | Apertium | - | - | - | **66.0** | **82.2** | **61.1** | **79.3** |
| | Marian | Parallel | 19.22K | WikiMatrix | **66.0** | **82.2** | **61.1** | **79.3** |
| | | FT | 2.7M | WikiMedia | | | | |
| | | BT | 103.9K | PILAR+WikiMatrix | | | | |
| arn | Apertium | - | - | - | 38.0 | 60.0 | 28.8 | 49.4 |
| | Marian | Parallel | - | - | **38.7** | **60.3** | **29.8** | **49.8** |
| | | MT | 64.1K | PILAR cat-arn | | | | |
| | | FT | 2.7M | WikiMedia [Tagged] | | | | |
| | | BT | 392.8K | PILAR + CCMatrix + PILARcat-arn | | | | |

Table 4: BLEU and chrF scores for our primary submissions in the constrained track

| Lang. | Model | Data | # Sent. | Source | Dev BLEU | Dev chrF | Test BLEU | Test chrF |
|---|---|---|---|---|---|---|---|---|
| ast | Apertium | - | - | - | 17.1 | 50.7 | 17.0 | 50.8 |
| | NLLB | Parallel | 533.7K | CCMatrix | **18.6** | **51.6** | **18.0** | **51.6** |
| | | FT | - | - | | | | |
| | | BT | 2.7M | WikiDump+PILAR+CCMatrix+WikiMedia | | | | |
| arg | Apertium | - | - | - | **66.0** | **82.2** | **61.1** | **79.3** |
| | Marian | Parallel | 13.7K | WikiMatrix | 65.9 | **82.2** | 61.0 | **79.3** |
| | | FT | 2.7M | WikiMedia | | | | |
| | | BT | 353.5K | WikiDump+PILAR+WikiMatrix | | | | |
| arn | Apertium | - | - | - | **38.0** | **60.0** | **28.8** | **49.4** |
| | Marian | Parallel | 8.9K | CCMatrix (es-oci) | 37.9 | **60.0** | **28.8** | **49.4** |
| | | MT | 64.1K | PILARcat-arn | | | | |
| | | FT | 2.7M | WikiMedia | | | | |
| | | BT | 346.2K | PILAR+CCMatrix+PILARcat-arn | | | | |

Table 5: BLEU and chrF scores for our primary submissions in the open track

A notable result are the relatively high scores of Llama3 zero-shot translation into Spanish, confirming our initial assessments of the potential leveraging this type of LLM to translate from low-resource into high-resource languages. Further variants such as few-shot translation might be worth exploring in this type of scenarios.

# 4 Main Results

The best results for our shared task submissions are summarised in Table 4 and Table 5 for the constrained and open tracks, respectively. We report BLEU and chrF scores on the PILAR development sets and on the task test sets, as reported on the OCELoT website.

## 4.1 Constrained Track

In the constrained setup, our focus was on optimising translation models within the set limits of OPUS data and pretrained models under one billion parameters. For Asturian, the best results were achieved via a fine-tuning of NLLB using both parallel data from CCMatrix and back-translations generated from the PILAR, CCMatrix and Wiki-Media corpora using our custom Marian model.

In the case of Aragonese and Aranese, training Marian models from scratch proved to be the most successful strategy. Given that NLLB was not specifically trained on these languages, this result was not unexpected. For these languages, we also incorporated forward-translations generated using Apertium and back-translations created with our Marian models. For Aranese, the use of parallel

| Lang. | Model | Data | # Sent. | Source | Not tagged | Tagged |
|---|---|---|---|---|---|---|
| ast | Apertium | - | - | - | **17.1** | - |
| | Marian | Parallel | 533.7K | CCMatrix | | |
| | | FT | 2.7M | WikiMedia | **16.9** | **17.4** |
| | | BT | - | - | | |
| arg | Apertium | - | - | - | **66.0** | - |
| | Marian | Parallel | 19.2K | WikiMatrix | | |
| | | FT | 2.7M | WikiMedia | **66.0** | 46.5 |
| | | BT | 103.9K | PILAR+WikiMatrix | | |
| arn | Apertium | - | - | - | 38.0 | - |
| | Marian | Parallel | - | - | | |
| | | MT | 64.1K | PILAR cat-arn | | |
| | | FT | 2.7M | WikiMedia | 37.9 | **38.7** |
| | | BT | 392.8K | PILAR + CCMatrix + PILARcat-arn | | |

Table 6: BLEU scores comparison between models trained with and without tags in the forward-translated data.

data from CCMatrix resulted in lower performance, likely due to the lower quality of these data, which were originally Spanish-Occitan alignments. The inclusion of the pivot translations from Catalan was also beneficial for the Spanish-Aranese pair.

Overall, when comparing our results to the baselines in Table 2, our custom models consistently outperformed the vanilla NLLB across all languages, particularly for Aragonese and Aranese, which were unseen by this model. The models trained for Asturian and Aranese also achieved higher scores than the Apertium baseline. For Aragonese however, our best submission could only match the Apertium baseline scores. This limitation is likely due to the influence of the forward-translations from the rule-based Apertium system, a factor which was not mitigated with data tagging.

### 4.2 Open Track

Our contributions to the open track were twofold: augmenting the training data by incorporating Asturian and Aragonese Wikipedia content, and generating back-translations using Llama3 in a zero-shot setting.

As shown in Table 5, these additions improved the BLEU score for Asturian by 0.5 points compared to the constrained track. However, the results for Aragonese did not benefit from the extra data, showing a slight decrease of 0.1 BLEU. For Aranese, the use of back-translations from Llama3 appeared to be detrimental, resulting in a performance drop of 0.8 BLEU points.

Overall, the open track models yielded mixed results, as the augmented data generated via back-translation and zero-shot LLM translation resulted

in either minor gains or losses. This might be due to the specifics of the development and test sets, in the sense that the augmented data might come from domains of little benefit to improve the translation on these datasets. The results of Section 3.2 are still important in our view, notably the quality of NMT and LLM translations for either direct use or data augmentation.

## 5 Discussion

As previously indicated, given the low performance of all models in Asturian in preliminary experiments, we used a filtered subset of 3,000 sentences from CCMatrix as development set. However, to ensure consistency across languages, we relied on the best-performing model on the original PILAR dev set as the criterion for model selection for submission, leading to the exclusion of models that performed better on our custom dev set.

| Model | Source | Official Dev | Custom Dev |
|---|---|---|---|
| Apertium | - | 17.1 | 79.8 |
| Open Submission | - | **18.6** | 78.4 |
| Marian | CCMatrix | 17.2 | **87.4** |

Table 7: BLEU scores in Spanish-Asturian on the official WMT development set and on our custom development set from CCMatrix.

For reference, Table 7 presents the results of the top-performing model on our dev set: a Marian model trained exclusively on CCMatrix data without any synthetic data. While this model shows lower performance than the one chosen for the official submission and is comparable to the baseline obtained with Apertium, it performed notably bet-

ter on our own development set. Considering the large differences in scores between the PILAR and custom dev sets, it would be interesting to examine in detail the differences between the two datasets in future work.

Among our best submissions to both tracks, only one dataset was tagged, namely the forward-translations based on Wikimedia in ES-ARG. We performed several additional experiments on the use of tags to discriminate between types of data, with the most salient results shown in Table 6. Tags on forward-translations were beneficial for Asturian and Aranese, but for Aragonese their use resulted in a substantial decrease of almost 20 BLEU points on the dev set. This variation might be due to the differing amounts of data available: Asturian and Aranese featured 500K and 364K sentence pairs without tags, respectively, while Aragonese only counted with 119K such pairs. Whereas tags have been shown to be a successful means to discriminate between parallel and other types of data, their use might thus be detrimental when tagged data largely dominate the other types of data.

## 6 Conclusions

We described our submission to the WMT 2024 shared task on translation into low-resource languages of Spain. We followed a multi-pronged approach based on data filtering and augmentation, with multiple types of models trained on different combinations of data with or without tagging. Although we improved over the baselines in general, the gains were minor overall on the development sets provided for the task. Nonetheless, our experiments showed the benefits of training dedicated NMT models, which proved optimal in most cases over fine-tuning pre-trained translation models. We also demonstrated the potential of zero-shot LLM-based translation for translation of the selected low-resource languages into Spanish, an interesting path for future research as standalone translation or as a source of data augmentation.

## Acknowledgments

## References

AI@Meta. 2024. Llama 3 model card.

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, OndÅ™ej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina EspaÃ±a-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (wmt21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88, Online. Association for Computational Linguistics.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Isaac Caswell, Ciprian Chelba, and David Grangier. 2019. Tagged back-translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 53–63, Florence, Italy. Association for Computational Linguistics.

Marta R. Costa-jussà. 2017. Why Catalan-Spanish neural machine translation? analysis, comparison and combination with standard rule and phrase-based technologies. In *Proceedings of the Fourth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial)*, pages 55–62, Valencia, Spain. Association for Computational Linguistics.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Jimmy Ba Diederik P. Kingma. 2015. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations*, San Diego, CA, USA.

Thierry Etchegoyhen, Eva Martínez Garcia, Andoni Azpeitia, Gorka Labaka, Iñaki Alegria, Itziar Cortes Etxabe, Amaia Jauregi Carrera, Igor Ellakuria Santos, Maite Martin, and Eusebi Calonge. 2018. Neural Machine Translation of Basque. In *Proceedings of the 21st Annual Conference of the European Association for Machine Translation*, pages 139–148.

Mikel L. Forcada and Francis M. Tyers. 2016. Apertium: a free/open source platform for machine translation and basic language technology. In *Proceedings of the 19th Annual Conference of the European Association for Machine Translation: Projects/Products*, Riga, Latvia. Baltic Journal of Modern Computing.

Samuel Frontull and Georg Moser. 2024. Rule-based, neural and LLM back-translation: Comparative insights from a variant of Ladin. In *Proceedings of the The Seventh Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2024)*, pages 128–138, Bangkok, Thailand. Association for Computational Linguistics.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024a. Idiomata cognitor.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024b. Pilar.

Harritxu Gete and Thierry Etchegoyhen. 2022. Making the most of comparable corpora in neural machine translation: a case study. *Language Resources and Evaluation*, 56(3):943–971.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.

Zhenhao Li and Lucia Specia. 2019. Improving neural machine translation robustness via data augmentation: Beyond back-translation. In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 328–336, Hong Kong, China. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021a. WikiMatrix: Mining 135M parallel sentences in 1620 language pairs from Wikipedia. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361, Online. Association for Computational Linguistics.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021b. CCMatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500, Online. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*, 30.

## A Training Hyperparameters

The Marian models were transformer-base models. Optimization was performed with Adam (Diederik P. Kingma, 2015), with $\alpha = 0.0003$, $\beta_1 = 0.9$, $\beta_2 = 0.98$, and $\epsilon = 10^{-9}$. We used a working memory of 20GB and automatically chose the largest mini-batch that fit the specified memory. The learning rate was set to increase linearly for the first 16,000 training steps and decrease afterward proportionally to the inverse square root of the corresponding step. The validation data was evaluated every 5000 steps.

For fine-tuning the NLLB model, optimization was performed using Adafactor (Shazeer and Stern, 2018), with a learning rate of 0.0001, a clipping threshold of 1.0, and weight decay set to 0.001. The training used a batch size of 32 and a maximum sequence length of 128 tokens.

Each model was trained on a Nvidia L40 with 48GB of VRAM. Early stopping was applied with a patience of 10 epochs to prevent overfitting.

## B Generation Parameters

For inference with Marian, we set a beam size of 6 and a normalization factor of 0.6.

For the NLLB model, implemented on the transformer library, the maximum input length was configured to 200 tokens, with a beam size of 4.

For Llama3, we set a maximum of 256 new tokens, enabled sampling with a temperature of 0.1, and set top-p to 0.9. We used the following prompt to direct the model to generate translations in the specified target language without additional commentary: "*Traduce a [Español|Aragonés|Aranés|Asturiano] la siguiente frase. No añadas ningún otro comentario.*" .

## C Catalan-Spanish MT Model

We considered two main options to translate Catalan into Spanish, as a means to create additional Aranese-Spanish data via pivot translation: the pretrained multilingual NLLB model or an in-house Marian model trained on parallel corpora from OPUS (namely: dogc, gnome, opensubs, tatoeba, ubuntu, globalvoices, wikimatrix, ted and paracrawl). The latter achieved significantly better results, as shown in Table 8 on a test set of 2,000 sentence pairs randomly sampled from OPUS data.

| Translation Model | BLEU Score |
|---|---|
| NLLB Model | 55.4 |
| Marian Model | **70.7** |

Table 8: BLEU scores for Catalan to Spanish translation.

# SJTU System Description for the WMT24 Low-Resource Languages of Spain Task

**Tianxiang Hu**[1*]  **Pei Zhang**[2*]  **Haoxiang Sun**[1]  **Ruize Gao**[2]  **Jialong Tang**[2]
**Baosong Yang**[2†]  **Rui Wang**[1†]

[1]Shanghai Jiao Tong University, Shanghai, China
[2]Tongyi Lab, Hangzhou, China
{hutianxiang,wangrui12}@sjtu.edu.cn
{xiaoyi.zp,yangbaosong.ybs}@alibaba-inc.com

## Abstract

This paper describes Shanghai Jiao Tong University low-resource languages of Spain translation systems for WMT24 shared task. We participate in the translation task on Spanish → Aragonese, Spanish → Aranese and Spanish → Asturian. Initially, we conduct preliminary experiments to assess the basic translation capabilities of various models and evaluate the impact of fine-tuning with different data types. We then choose to fine-tune the Qwen2-0.5B model using a forward synthesized pseudo-corpus from the Apertium translation system to replicate its fundamental performance. Building on this distillation model, we explore three optimization strategies across the three language directions: (1) Assembling the provided FLORES+ dev sets into a 5-shot format translation training dataset and performing few-shot fine-tuning to enhance model performance. (2) Utilizing the FLORES+ dev sets as training data and applying the Contrastive Preference Optimization (CPO) strategy for further refinement. (3) Retrieving the 20 most similar translation examples from the FLORES+ dev sets using the BM25 algorithm and performing 20-shot translations with the Claude 3.5-sonnet model. After evaluating these strategies, we select the best-performing approach for each language pair as our submission result.

## 1 Introduction

This paper introduces our submissions to the WMT24 Low-Resource Languages of Spain Task. We participate in the competitions for three translation directions: Spanish → Aragonese, Spanish → Aranese, and Spanish → Asturian. For the Spanish → Aragonese and Spanish → Aranese directions, we ultimately submit constrained results, while for Spanish → Asturian, we provide unconstrained (open system) results.

Neural machine translation (NMT) systems have achieved substantial advancements in recent years (Vaswani et al., 2017). However, training neural translation models typically necessitates large-scale parallel corpora (Ranathunga et al., 2021). In many low-resource scenarios, the availability of sufficient parallel data for training is limited, making low-resource translation a critical and valuable research area (Arivazhagan et al., 2019; Wang et al., 2021; Ranathunga et al., 2021). This competition task focuses on translating between Spanish and three other languages: Aragonese, Aranese, and Asturian. Of these, Aragonese and Aranese face particular challenges due to their relatively scarce parallel corpora. While the OPUS[1] website provides a considerable amount of parallel data, the quality of this data remains relatively low.

We initially conduct a preliminary evaluation of translation capabilities using models such as Apertium[2], GPT4 (Achiam et al., 2023), Llama-3 (AI@Meta, 2024), and Qwen2 (Yang et al., 2024) across the three language pairs. Our findings indicate that the Apertium translation system serves as a strong baseline, particularly in terms of BLEU (Papineni et al., 2002; Post, 2018) scores. Subsequently, we explore fine-tuning the Qwen2-0.5B model with various types of synthetic data and data from diverse domains. This exploration reveals that this task presents unique challenges compared to previous low-resource translation tasks. Specifically, forward-translated (Zhang and Zong, 2016) data and data from the OPUS NLLB corpus result in improved performance on dev test sets. We ultimately select the NLLB Spanish corpus from OPUS and perform forward translation using Apertium to generate the corresponding parallel pseudo-corpus. Fine-tuning Qwen2-0.5B with this synthetic data enables us to closely replicate the performance of the Apertium translation system.

---

---

[1]https://opus.nlpl.eu
[2]https://apertium.org

Although simple forward distillation can effectively replicate the performance of the Apertium system, it does not exceed it, and the distilled model does not yield further performance improvements. To enhance the model's effectiveness, high-quality data is crucial. We randomly select a portion of the provided dev test set for additional fine-tuning, with the remaining portion designated as the new dev set. Building on this distilled model, we explore three optimization strategies across three language pairs: (1) We aggregate the provided FLORES+ dev sets into a 5-shot format translation training dataset and perform few-shot fine-tuning (Alves et al., 2023) to further refine the model. (2) We use the FLORES+ dev sets as training data and apply the Contrastive Preference Optimization (CPO) (Xu et al., 2024) strategy to improve model performance. (3) We retrieve the 20 most similar translation examples from the FLORES+ dev sets using the BM25 algorithm and employ the Claude 3.5-sonnet model[3] for 20-shot translations (Agrawal et al., 2022).

## 2 Preliminary Experiment

In this section, we first investigate the basic translation capabilities of various models and identify the Apertium translation system as a particularly strong baseline. We then examine the fine-tuning of the Qwen2-0.5B model using different types of data, which reveals that this task presents unique challenges compared to previous low-resource scenarios. Ultimately, we select the NLLB[4] Spanish corpus from OPUS, forward-translate it using Apertium to create a parallel pseudo-corpus, and fine-tune Qwen2-0.5B with this synthetic data.

**Data** The results presented in this section are derived from experiments conducted on the official FLORES+ dev test sets[5], which come from Pan-Iberian Language Archival Resource (PILAR). The three language pairs under consideration are Spanish → Aragonese (spa-arg), Spanish → Aranese (spa-arn), and Spanish → Asturian (spa-ast), each comprising 997 sentences.

### 2.1 Translation capabilities of different models

We begin by evaluating the BLEU (Papineni et al., 2002; Post, 2018) performance of five models

(Apertium, GPT-4, Llama3-8B, Llama3-70B, and Qwen2-0.5B) on the three language pairs in this task using the FLORES+ dev sets. For the 1-shot scenario, the format used is as follows: "Translate the following sentence from <src lang> into <tgt lang>.\n <src lang>: <src example1>.\n <tgt lang>: <tgt example1>.\n \n Translate the following sentence from <src lang> into <tgt lang>.\n <src lang>: <src sentence>.\n <tgt lang>:". In the 5-shot scenario, this format is extended by providing five examples instead of one. The few-shot examples are randomly sampled from the corresponding language FLORES+ dev sets without repetition.

As shown in Table 1, our results indicate that the Apertium translation system serves as a very strong baseline, significantly outperforming other large models in BLEU scores for the Spanish → Aragonese (spa-arg) and Spanish → Aranese (spa-arn) language pairs. Notably, even the widely used GPT-4 scores considerably lower in BLEU compared to the Apertium system. This superior performance of Apertium may be attributed to the fact that the dev test sets for these two language pairs were derived from Apertium's translations with post-editing. Additionally, we observed that increasing the number of example shots in translation leads to a substantial improvement in performance. This suggests that, for these low-resource languages, providing translation examples enhances the ability of large models to learn and perform the translation task more effectively.

|  | spa-arg | spa-arn | spa-ast |
|---|---|---|---|
| Apertium | 66.0 | 38.0 | 17.1 |
| GPT4 1shot | 35.9 | 16.1 | 18.6 |
| GPT4 5shot | 37.4 | 17.7 | 19.1 |
| Llama3-8B 1shot | 36.3 | 7.8 | 16.6 |
| Llama3-8B 5shot | 41.0 | 10.6 | 18.3 |
| Llama3-70B 1shot | 46.4 | 15.6 | 19.4 |
| Llama3-70B 5shot | 52.4 | 19.9 | 22.4 |
| Qwen2-0.5B 1shot | 22.7 | 4.1 | 8.6 |
| Qwen2-0.5B 5shot | 22.7 | 4.2 | 8.9 |

Table 1: BLEU evaluation of different models on dev test sets for three language pairs. Apertium translation system demonstrates a strong baseline.

### 2.2 Effects of different types of data

To explore the types of data that can be used for fine-tuning the base model, we conduct preliminary experiments focusing exclusively on Aragonese.

As shown in Table 2, we evaluate the impact of different data types on fine-tuning performance. Our findings indicate that forward translation (FT) (Zhang and Zong, 2016) outperforms back translation (Sennrich et al., 2016). This result may be attributed to the fact that the dev test set is derived from Apertium with post-editing, which means that the Aragonese side of the dev test set reflects Apertium's translation style rather than the natural language style of Aragonese. In contrast, back translation targets the authentic Aragonese language style, which does not align with the style of the dev test set, potentially leading to BLEU scores that do not accurately represent the actual translation quality. However, due to the extremely low-resource nature of this language, we have to rely on the official dev test set and BLEU scores for optimization.

Additionally, the table highlights another critical factor affecting performance: the source of the fine-tuning data. Using Spanish monolingual data from the OPUS NLLB corpus[6] provides a noticeable performance advantage over using WMT news[7], Pilar[8] or random samples from OPUS[9]. This suggests that the domain of the dev test set is more closely aligned with the OPUS NLLB corpus, facilitating better adaptation to the dev set for this task. Furthermore, we observe that mixing data from different domains or simultaneously using both BT and FT does not enhance performance, despite increasing the volume of data. In fact, this approach slightly degrades the original performance.

## 2.3 Final Distillation Experiment

Based on the experimental results discussed above, we first perform basic filtering on the NLLB Spanish corpus from OPUS and then randomly sample 1 million sentences. We use the Apertium translation system to translate these 1 million Spanish sentences into the three target languages, creating a parallel pseudo-corpus. We then fine-tune the open-source Qwen2-0.5B model separately for each language using this pseudo-corpus. During training, we fine-tune the model for 1.5 epochs with a batch size of 64, a learning rate of 1e-05, and a weight decay of 0.1. For decoding, we employ beam search with a beam size of 4. As shown

| Data size | Data source | Data type | BLEU |
|---|---|---|---|
| 16k | OPUS | bilingual | 37.7 |
| 16k | OPUS | FT | 61.7 |
| 16k | News | FT | 53.0 |
| 16k | OPUS NLLB | FT | 63.8 |
| 16k | OPUS | BT | 41.1 |
| 16k | Pilar | BT | 34.3 |
| 32k | OPUS | FT+BT | 59.6 |
| 32k | OPUS+News | FT | 59.8 |

Table 2: BLEU evaluation on fine-tuning Qwen2-0.5B using different types of data. Data size refers to the training data size. FT refers to forward translation of Spanish to comprise synthesized parallel data; BT refers to backward translation of Aragonese to comprise synthesized parallel data; News referes to the WMT news.

in Table 3, this approach effectively replicates the baseline performance of the Apertium translation system.

| | spa-arg | spa-arn | spa-ast |
|---|---|---|---|
| Apertium | 66.0 | 38.0 | 17.1 |
| distillation model | 66.0 | 38.0 | 17.0 |

Table 3: BLEU evaluation of the distillation model on dev test sets for three language pairs. We have replicated the baseline capability of Apertium translation system.

## 3 Method

In Section 3, we initially replicate the performance of the strong baseline system Apertium using the Qwen2-0.5B model but are unable to surpass it. We also observe that fine-tuning the model with filtered bilingual data resulted in decreased BLEU scores, likely due to the low quality of available bilingual data. The synthetic pseudo-corpus generated through forward translation reach its performance limits, as further improvements could not be achieved with the distilled model. To address this, we randomly select 700 sentences from the provided dev test set for additional fine-tuning, reserving the remaining 297 sentences as the new dev set. We next explore three optimization strategies to further enhance translation performance for these three language pairs.

## 3.1 Dev 5shot SFT

In Table 1, we observe that providing few-shot examples to large language models improves trans-

[6] https://opus.nlpl.eu/NLLB/corpus/version/NLLB
[7] https://www.statmt.org/wmt11/translation-task.html
[8] https://github.com/transducens/PILAR
[9] https://opus.nlpl.eu

lation performance. However, supervised fine-tuning can reduce some of these few-shot capabilities (Alves et al., 2023). To maintain consistency in the inference format, we structure the fine-tuning data into a 5-shot format during training. For inference, we also use the 5-shot format, with few-shot examples randomly selected from the dev test set. The fine-tuning data consists of 700 sentences from the previously mentioned dev test set.

## 3.2 Dev CPO

Given that the official dev test set is derived from post-edited results, our goal is to assist the model in learning the subtle distinctions between pre-edited and post-edited translations, thereby enhancing its translation capabilities. DPO (Rafailov et al., 2023) is a training strategy focused on optimizing preferences, while CPO (Xu et al., 2024) builds upon DPO by providing further refinements. The following is the formulation of CPO loss:

$$\mathcal{L}(\pi_\theta) = - \mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}} \Big[ \log \sigma \Big( \beta \log \pi_\theta(y_w|x)$$
$$- \beta \log \pi_\theta(y_l|x) \Big) \Big], \quad (1)$$

$$\min_\theta \underbrace{\mathcal{L}(\pi_\theta)}_{\mathcal{L}_{\text{prefer}}} \underbrace{-\mathbb{E}_{(x,y_w)\sim\mathcal{D}}[\log \pi_\theta(y_w|x)]}_{\mathcal{L}_{\text{NLL}}}, \quad (2)$$

where x is source sentence, $y_w$ is preferred translation, $y_l$ is less preferred translation, $\mathcal{D}$ is a dataset of comparisons.

In our approach, we use translations produced by Apertium as negative examples and the corresponding results from the dev test set as positive examples. This CPO training allows the model to learn the nuanced differences between positive and negative instances.

## 3.3 Dev fewshot BM25 with LLM

Previous research suggests that providing similar parallel translation pairs as guidance can improve translation quality with large language models (Agrawal et al., 2022). To leverage this, we use the BM25 algorithm to retrieve several of the most similar translation examples from the dev test set based on the source sentences. These examples are concatenated into the previously described few-shot translation format and positioned before the sentence to be translated. We then employ state-of-the-art LLMs, such as GPT-4 and Claude-3.5, for the translation process.

## 3.4 Post-processing

We observe that translations produced by large language models may encounter issues such as omissions, over-translation, and non-following with instructions (Jiao et al., 2023; Xu et al., 2023). To address these issues, we apply the following rule-based post-processing:

1. For translations generated by the Dev 20-shot BM25 method with LLMs, if the output fails to adhere to instructions, for instance, if it includes phrases such as "I apologize" or "sorry", we perform a retranslation. If the correct translation is not achieved after three attempts, we revert to the translation produced by the Apertium software.

2. Replace any translations where language detection is incorrect with those generated by Apertium software.

3. Replace any translations where the ratio of the length of the source text to the translated text is less than 0.75 or greater than 1.3 with translations generated by Apertium software.

## 4 Experiment

**Data** In this Section, we randomly selected 700 sentences from the provided dev test set for additional fine-tuning, leaving the remaining 297 sentences as the new dev set.

**Experiment Details** For SFT, we fine-tune the distillation model for 5 epochs with a batch size of 8, a learning rate of 1e-05, and a weight decay of 0.1. For decoding, we use beam search with a beam number of 4. For few-shot BM25, we use the BM25 algorithm to select a number of the most similar examples (excluding the sentence itself) from the 997 sentences in the dev set for few-shot translation.

**Results** As illustrated in Table 4, the BLEU scores for the three language pairs across various methods demonstrate noticeable performance improvements over the Distillation model. Specifically, the best performance for Spanish → Aragonese (spa-arg) is achieved with the Distillation model + dev 5-shot SFT, for Spanish → Aranese (spa-arn) with the Distillation model + dev CPO, and for Spanish → Asturian (spa-ast) with Claude 3.5-sonnet + 20-shot BM25.

Furthermore, the dev 5-shot SFT method yields a more consistent performance improvement compared to direct dev SFT. Among the models evaluated, Claude 3.5-sonnet generally outperforms GPT-4-turbo across these three low-resource language pairs, and BM25 retrieval of similar examples significantly boosts translation performance.

|  | spa-arg | spa-arn | spa-ast |
|---|---|---|---|
| Distillation model | 67.4 | 39.5 | 17.0 |
| + dev SFT | 69.3 | 40.5 | 17.3 |
| + dev 5shot SFT | **69.9** | 40.8 | 17.4 |
| + dev CPO | 69.7 | **41.4** | 17.3 |
| GPT4-turbo |  |  |  |
| + 5shot | 40.5 | 32.3 | 20.1 |
| + 5shot BM25 | 44.3 | 33.2 | 20.5 |
| + 20shot BM25 | 47.5 | 33.6 | 21.4 |
| Claude3.5-sonnet |  |  |  |
| + 5shot | 47.8 | 35.2 | 22.9 |
| + 5shot BM25 | 53.6 | 37.4 | 24.2 |
| + 20shot BM25 | 59.9 | 38.1 | **25.2** |

Table 4: BLEU evaluation of different methods on partitioned dev test sets for three language pairs. Our methods all achieve certain performance improvements. For the Aragonese language pair, the best strategy is dev 5-shot SFT. For the Aranese language pair, the optimal strategy is dev CPO. For Asturian language pair, the best approach is using Claude 3.5-sonnet for 20-shot BM25 translation.

## 5 Conclusion

This paper presents the Shanghai Jiao Tong University translation systems for low-resource Spanish languages in the WMT24 shared task. We first create synthetic data through forward distillation using the Apertium translation system, then fine-tune the Qwen2-0.5B model to establish a basic baseline capability. Subsequently, we apply three optimization strategies using the dev test sets: 5-shot format fine-tuning, Contrastive Preference Optimization, and 20-shot translation with BM25 retrieval. Our experiments demonstrate that all three methods lead to performance improvements.

## Acknowledgements

## References

OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, et al. 2023. Gpt-4 technical report.

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2022. In-context examples selection for machine translation. In *Annual Meeting of the Association for Computational Linguistics*.

AI@Meta. 2024. Llama 3 model card.

Duarte M. Alves, Nuno M. Guerreiro, Joao Alves, José P. Pombal, Ricardo Rei, Jos'e G. C. de Souza, Pierre Colombo, and André Martins. 2023. Steering large language models for machine translation with fine-tuning and in-context learning. In *Conference on Empirical Methods in Natural Language Processing*.

Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George F. Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen, and Yonghui Wu. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges. *CoRR*, abs/1907.05019.

Wenxiang Jiao, Jen-tse Huang, Wenxuan Wang, Zhi-wei He, Tian Liang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. ParroT: Translating during chat using large language models tuned with human translation and feedback. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15009–15020, Singapore. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *ArXiv*, abs/2305.18290.

Surangika Ranathunga, En-Shiun Annie Lee, Marjana Prifti Skenduli, Ravi Shekhar, Mehreen Alam, and Rishemjit Kaur. 2021. Neural machine translation for low-resource languages: A survey. *ArXiv preprint*, abs/2106.15115.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pages 86–96.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Rui Wang, Xu Tan, Renqian Luo, Tao Qin, and Tie-Yan Liu. 2021. A survey on low-resource neural machine translation. In *International Joint Conference on Artificial Intelligence*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *ArXiv preprint*, abs/2309.11674.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, et al. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Jiajun Zhang and Chengqing Zong. 2016. Exploiting source-side monolingual data in neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1535–1545.

# Multilingual Transfer and Domain Adaptation for Low-Resource Languages of Spain

**Yuanchang Luo, Zhanglin Wu, Daimeng Wei, Hengchao Shang, Zongyao Li,**
**Jiaxin Guo, Zhiqiang Rao, Shaojun Li, Jinlong Yang,**
**Yuhao Xie, Jiawei Zheng Bin Wei, Hao Yang**
Huawei Translation Service Center, Beijing, China
{luoyuanchang1,wuzhanglin2,weidaimeng,shanghengchao,lizongyao,
guojiaxin1,raozhiqiang,lishaojun18,yangjinlong7,xieyuhao2,
zhengjiawei15,weibin29,yanghao30}@huawei.com

## Abstract

This article introduces the submission status of the Translation into Low-Resource Languages of Spain task at (WMT 2024) by Huawei Translation Service Center (HW-TSC). We participated in three translation tasks: spanish to aragonese (es→arg), spanish to aranese (es→arn), and spanish to asturian (es→ast). For these three translation tasks, we use training strategies such as multilingual transfer, regularized dropout, forward translation and back translation, labse denoising, transduction ensemble learning and other strategies to neural machine translation (NMT) model based on training deep transformer-big architecture. By using these enhancement strategies, our submission achieved a competitive result in the final evaluation.

## 1 Introduction

Neural machine translation (MT) (Lyu et al., 2019; Bahdanau et al., 2014; Gehring et al., 2017) allows translation systems to be trained end-to-end without having to deal with issues like word alignment, translation rules, and complex decoding algorithms that characterize statistical machine translation systems (SMT) (Koehn et al., 2007). Recently, MT technology has evolved towards large language models (LLMs) (Guo et al., 2024). Although neural machine translation has developed rapidly in recent years, it relies heavily on big data - large-scale, high-quality bilingual corpora. Due to the cost and scarcity of real corpora, synthetic data plays an important role in improving translation quality. Existing methods for synthesizing data in NMT focus on leveraging monolingual data during training. Among them, forward translation (Abdulmumin et al., 2021), back translation (Abdulmumin et al., 2021) and data diversity (Nguyen et al., 2020) have been widely used to generate synthetic bilingual corpora. Such synthetic data can be used to improve the performance of NMT

models(Wu et al., 2023b). (Wei et al., 2023) also considers the style of the training data and exploits it to improve performance. Although synthetic data is efficient, synthetic data inevitably contains noise and erroneous translations. Denoising can prevent the training of NMT models from being interfered by noisy synthetic data by introducing high-quality real data as guidance. Another direction to improve the performance of NMT models is to use more efficient training strategies. For example, by mixing similar language data together to train a multi-language pre-training model (Li et al., 2022), due to the shared vocabulary, encoding layer and decoding layer parameters and language similarity, languages with less data can benefit from languages with more data. Regularized dropout (Wu et al., 2021) allows the NMT model to more effectively utilize limited data during the training process. Transduction ensemble learning (Wang et al., 2020) can aggregate the translation capabilities of multiple models into one model.

For the Translation into Low-Resource Languages of Spain task of WMT 2024, we participated in the es→arg, es→arn and es→ast language pair. We use training strategies such as multilanguage pre-training models (Li et al., 2022), regularized dropout (Wu et al., 2021), forward translation (Abdulmumin et al., 2021), back translation (Abdulmumin et al., 2021), Labse denoise (Feng et al., 2020) and transduction ensemble learning (Wang et al., 2020) to train neural machine translation (NMT) models based on deep Transformer architecture.

Next, this article will expand on the details of our translation system in different translation tasks. The structure of the remaining sections is as follows: Section 2 introduces the data scale and data preprocessing process; Section 3 describes the overview of the NMT system; Section 4 gives the parameter settings, data processing results and experimental results; Section 6 gives System conclu-

sions were drawn.

## 2 Dataset

### 2.1 Data Size

In accordance with the requirements of the WMT 2024 outline, on the Translation into Low-Resource Languages of Spain machine translation task, we used the officially provided data to train the NMT system from scratch. Table 1 shows the training data size for each language pair of the bilingual machine translation task. These language pairs include Spanish to Aragonese (es→arg), Spanish to Arabic (es→arn) and Spanish to Asturian (es→ast).

|  | es→arg | es→arn | es→ast |
|---|---|---|---|
| Bilingual | 0.06M | 2.04M | 13.36M |
| Source Monolingual | 0.4M | 8M | 8M |
| Target Monolingual | 0.26M | 6M | 3M |

Table 1: Data size for each bilingual machine translation task

### 2.2 Data Pre-processing

The data pre-processing process is as follows:

- Remove duplicate sentences or sentence pairs.

- Remove invisible characters and xml escape characters.

- Convert full-width symbols to half-width symbols.

- Use fast_align (Dyer et al., 2013) to filter poorly aligned sentence pairs.

- Filter out sentences with more than 80 tokens in bilingual data.

- Remove sentences with duplicate tokens.

- When performing subword segmentation, joint sentencepiece (Kudo and Richardson, 2018) is used for es→arg, es→arn and es→ast translation tasks.

## 3 NMT System

### 3.1 System Overview

Transformer is the state-of-the-art model structure in recent MT evaluations. There are two parts of research to improve this kind: the first part uses wide networks (eg: Transformer-Big (Vaswani, 2017)),



Figure 1: The overall training flow chart of our NMT system on the different translation tasks.

and the other part uses deeper language representations (eg: Deep Transformer (Wang et al., 2019)). For all MT tasks, we combine these two improvements, adopting the Deep Transformer-Big (Wu et al., 2023a) model structure to train the NMT system. Deep Transformer-Big uses pre-layer normalization, features 25-layer encoder, 6-layer decoder, 16-heads self-attention, 1024-dimensional word embedding and 4096-dimensional ffn embedding.

Fig. 1 shows the overall training flow chart of our NMT system on the Translation into Low-Resource Languages of Spain task, we use multilingual transfer (Li et al., 2022), regularization dropout (Wu et al., 2021), forward translation (Abdulmumin et al., 2021), back translation (Abdulmumin et al., 2021), Labse denoise (Feng et al., 2020) and transduction ensemble learning (Wang et al., 2020) and other training strategies are used to train neural machine translation (NMT) models based on deep Transformer-big architecture.

### 3.2 Multilingual Transfer

Recent researches have shown that multilingual models outperform their bilingual counterparts, particularly when the number of languages in the system is limited and those languages are related (Li et al., 2022). This is mainly due to the capability of the model to learn interlingual knowledge (shared semantic representation between languages). Transfer learning using pre-trained multilingual model has shown very promising results

for low resource tasks. In this task, we first select a multilingual system as the base system, then fine-tune the system with low resource language pairs.

Specifically, we add the "<arg>" tag to the Spanish side of the es→arg bilingual data, the "<arn>" tag to the Spanish side of the es→arn bilingual data, and the "<ast>" tag to the Spanish side of the es→ast bilingual data, and sample them. Mix shuf to train a one-to-many pre-training model; sample the es→arg, es→arn and es→ast original bilingual data and then mix shuf to train a many-to-one pre-training model. Then, the one-to-many pre-training model and the many-to-one pre-training model are trained by using the original bilingual data, and three translation models from Spanish to Aragonese, Arabic, and Asturian and three translation models from Aragonese, Arabic, and Asturian to Spanish are obtained.

### 3.3 Regularization Dropout

Dropout (Srivastava et al., 2014) is a widely used technique for regularizing deep neural network training, which is crucial to prevent over-fitting and improve the generalization ability of deep models. Dropout performs implicit ensemble by simply dropping a certain proportion of hidden units from the neural network during training, which may cause an unnegligible inconsistency between training and inference. Regularized Dropout (R-Drop) (Wu et al., 2021) is a simple yet more effective alternative to regularize the training inconsistency induced by dropout. Concretely, in each mini-batch training, each data sample goes through the forward pass twice, and each pass is processed by a different sub model by randomly dropping out some hidden units. R-Drop forces the two distributions for the same data sample outputted by the two sub models to be consistent with each other, through minimizing the bidirectional Kullback-Leibler (KL) divergence (Van Erven and Harremos, 2014) between the two distributions. In this way, the inconsistency between the training and inference stage can be alleviated.

### 3.4 Forward translation and Back translation

Forward translation, also known as self-training (Abdulmumin et al., 2021), is one of the most commonly used data augmentation methods. FT has proven effective for improving NMT performance by augmenting model training with synthetic parallel data. Generally, FT is performed in three steps:

(1) randomly sample a subset from the large-scale source monolingual data; (2) use a "teacher" NMT model to translate the subset data into the target language to construct the synthetic parallel data; (3) combine the synthetic and authentic parallel data to train a "student" NMT model.

Apertium is a free/open-source rule-based architecture for MT that consists of a pipeline of modules performing part-of-speech disambiguation and tagging, lexical transfer, lexical selection, chunk-level or recursive structural transfer, and morphological generation. To make our model better, we use Apertium as a "teacher" model to produce pseudo-corpus.

Back translation (BT) (Abdulmumin et al., 2021) refers to translating the target monolingual data into the source language, and then using the synthetic data to increase the training data size. This method has been proven effective to improve the NMT model performance.

We use the machine translation model obtained by Multilingual Transfer to produce back translation synthetic parallel data, and mix it with forward translation synthetic parallel data and authentic parallel data for training, which can achieve better results than FT or BT.

### 3.5 Labse Denoising

Due to the low quality of our bilingual data, we use LaBSE (Feng et al., 2020) to calculate the semantic similarity of each bilingual sentence pair and exclude bilingual sentence pairs with similarity scores below 0.7 from our training corpus. Use these clean data to better train the model.

### 3.6 Transductive Ensemble Learning

Ensemble learning (Garmash and Monz, 2016), which aggregates multiple diverse models for inference, is a common practice to improve the accuracy of machine learning tasks. However, it has been observed that the conventional ensemble methods only bring marginal improvement for NMT when individual models are strong or there are a large number of individual models. Transductive Ensemble Learning (TEL) (Wang et al., 2020) study how to effectively aggregate multiple NMT models under the transductive setting where the source sentences of the test set are known. TEL uses dev sets finetune a strong model, which boosts strong individual models with significant improvement and benefits a lot from more individual models.

| | BLEU | | | ChrF++ | | |
|---|---|---|---|---|---|---|
| FLORES+ dev sets | es→arg | es→arn | es→ast | es→arg | es→arn | es→ast |
| NMT baseline | 38.5 | 8.5 | 17.3 | 64.6 | 34.3 | 46.6 |
| + FT & BT | 41.7 | 9.5 | 16.9 | 64.8 | 34.9 | 45.5 |
| + Labse denoising | 48 | 10.1 | 17.5 | 72.4 | 38.8 | 47.5 |
| FLORES+ devtest sets | es→arg | es→arn | es→ast | es→arg | es→arn | es→ast |
| + TEL | **63** | **26.3** | **19.8** | **80.3** | **47.9** | **52.2** |

Table 2: BLEU and ChrF++ scores of es→arg, es→arn and es→ast NMT systems

## 4 Experiment

### 4.1 Setup

We use the open-source fairseq (Ott et al., 2019) to train NMT models, and then use SacreBLEU (Post, 2018) and Chrf++ to measure system performance. The main parameters are as follows: each model is trained using 8 V100 GPUs, batch size is 4096, parameter update frequency is 1, and learning rate is 5e-4. The number of warmup steps is 4000, and model is saved every 1000 steps. The architecture we used is described in section 3.1. We adopt dropout, and the rate varies across different training phases. R-Drop (Srivastava et al., 2014) is used in model training, and we set $\lambda$ to 5.

### 4.2 Data processing

| | es→arg | es→arn | es→ast |
|---|---|---|---|
| Bilingual | 0.06M | 2.04M | 13.36M |
| Data Pre-processing | 0.04M | 1.51M | 3.91M |
| Labse Filter | 0.03M | 1.16M | 1.92M |
| Upsampling | 0.56M | 1.74M | 1.92M |

Table 3: Data size for each bilingual machine translation task after data pre-processing

Due to the poor quality of bilingual data in low-resource languages, after the rule cleaning mentioned in section 2.2 and the labse model cleaning mentioned in section 3.2, the amount of data is smaller, and the data amount of es→arg, es→arn and es→ast is quite different. When training one-to-many and many-to-one pre-training models, if the amount of bilingual data for a certain language direction is too small, the translation quality will be extremely poor. Therefore, Following (Conneau and Lample, 2019; Liu et al., 2020) we re-balance the training set by upsampling data from each language $l$ with a ratio:

$$\lambda_l = \frac{1}{p_l} \frac{p_l^{1/T}}{\sum_{l=1}^n p_l^{1/T}} \quad with \quad p_l = \frac{n_l}{\sum_{l=1}^n n_l}$$

where, $T$ is the temperature parameter and we set $T$ to 2. $n_l$ is the number of utterances for language $l$ in the training set. The data amount changes as shown in the following table 3.

### 4.3 Results

Tables 2 shows the evaluation results of es→arg, es→arn and es→ast NMT systems on the brand new FLORES+ dev sets and devtest sets, the results of dev test sets are obtained through OCELoT submission. We use Multilingual Transfer and R-Drop to build a strong baseline, then use FT and BT for data enhancement, and use Labse denoising for more efficient training, and finally use Transductive Ensemble Learning to ensemble multiple models ability.

As can be seen from the table above, after FT & BT and Labse denoising, the translation quality from Spanish to three directions has been improved to varying degrees. This shows that for low-resource scenarios, these two strategies can expand the amount of data and improve the quality of the data. Enhance the translation quality of machine translation models. Among them, the improvement of both strategies in the es→arg direction is higher than that of the other two directions, and the bilingual data of es→arg is also the least. This shows that FT & BT's strategy of expanding the amount of data and labse denoising's strategy of improving data quality are both in situations where the amount of bilingual data is small, The effect is more obvious.

In addition, after Transductive Ensemble Learning, the BLEU value of FLORES+ devtest sets has been greatly improved compared to the FLORES+ dev sets test set. Although it is not the same test set, the BLEU value has improved across latitudes, which shows that The fields of dev sets and devtest sets are very consistent, and Transductive Ensemble Learning, a strategy that utilizes dev sets, can maximize the translation effect of the model on the

test set in the same field.

## 5 Conclusion

This paper presents HW-TSC's submission to the Translation into Low-Resource Languages of Spain task of WMT 2024. For both translation tasks, we use a series of training strategies to train NMT models based on the deep Transformer-big architecture. By using these enhancement strategies, our submission achieves a competitive result in the final evaluation. For example, #607 in the spanish to aragonese constrained submissions, #608 in the spanish to aranese constrained submissions, and #606 in the spanish to asturian constrained submissions.

## References

Idris Abdulmumin, Bashir Shehu Galadanci, and Abubakar Isa. 2021. Enhanced back-translation for low resource neural machine translation using self-training. In *Information and Communication Technology and Applications: Third International Conference, ICTA 2020, Minna, Nigeria, November 24–27, 2020, Revised Selected Papers 3*, pages 355–371. Springer.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Alexis Conneau and Guillaume Lample. 2019. Cross-lingual language model pretraining. *Advances in neural information processing systems*, 32.

Chris Dyer, Victor Chahuneau, and Noah A Smith. 2013. A simple, fast, and effective reparameterization of ibm model 2. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2020. Language-agnostic bert sentence embedding. *arXiv preprint arXiv:2007.01852*.

Ekaterina Garmash and Christof Monz. 2016. Ensemble learning for multi-source neural machine translation. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1409–1418.

Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. In *International conference on machine learning*, pages 1243–1252. PMLR.

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pages 639–649. Association for Computational Linguistics.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions*, pages 177–180. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.

Shaojun Li, Yuanchang Luo, Daimeng Wei, Zongyao Li, Hengchao Shang, Xiaoyu Chen, Zhanglin Wu, Jinlong Yang, Zhiqiang Rao, Zhengzhe Yu, et al. 2022. Hw-tsc systems for wmt22 very low resource supervised mt task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1098–1103.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pretraining for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

He Lyu, Ningyu Sha, Shuyang Qin, Ming Yan, Yuying Xie, and Rongrong Wang. 2019. Advances in neural information processing systems. *Advances in neural information processing systems*, 32.

Xuan-Phi Nguyen, Shafiq Joty, Kui Wu, and Ai Ti Aw. 2020. Data diversification: A simple strategy for neural machine translation. *Advances in Neural Information Processing Systems*, 33:10018–10029.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks

from overfitting. *The journal of machine learning research*, 15(1):1929–1958.

Tim Van Erven and Peter Harremos. 2014. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820.

Ashish Vaswani. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762*.

Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao. 2019. Learning deep transformer models for machine translation. *arXiv preprint arXiv:1906.01787*.

Yiren Wang, Lijun Wu, Yingce Xia, Tao Qin, ChengXiang Zhai, and Tie-Yan Liu. 2020. Transductive ensemble learning for neural machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6291–6298.

Daimeng Wei, Zhanglin Wu, Hengchao Shang, Zongyao Li, Minghan Wang, Jiaxin Guo, Xiaoyu Chen, Zhengzhe Yu, and Hao Yang. 2023. Text style transfer back-translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 7944–7959. Association for Computational Linguistics.

Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, Tie-Yan Liu, et al. 2021. R-drop: Regularized dropout for neural networks. *Advances in Neural Information Processing Systems*, 34:10890–10905.

Zhanglin Wu, Daimeng Wei, Zongyao Li, Zhengzhe Yu, Shaojun Li, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Yuhao Xie, Lizhi Lei, et al. 2023a. Treating general mt shared task as a multi-domain adaptation problem: Hw-tsc's submission to the wmt23 general mt shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 170–174.

Zhanglin Wu, Zhengzhe Yu, Zongyao Li, Daimeng Wei, Yuhao Xie, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Zhiqiang Rao, Shaojun Li, et al. 2023b. Hw-tsc's neural machine translation system for ccmt 2023. In *China Conference on Machine Translation*, pages 13–27. Springer.

# TRIBBLE - TRanslating IBerian languages Based on Limited E-resources

**Igor Kuzmin**[1] **Piotr Przybyła**[1,2] **Euan McGill**[1] **Horacio Saggion**[1]

[1] LaSTUS Lab, TALN Group, Universitat Pompeu Fabra, Barcelona, Spain

[2] Institute of Computer Science, Polish Academy of Sciences, Warsaw, Poland

{igor.kuzmin, piotr.przybyla, euan.mcgill, horacio.saggion}@upf.edu

Figure 1: Language family tree diagram (partial) focusing on the Iberian peninsula

## 1 Introduction

In this short overview paper[1], we describe our system submission for the language pairs Spanish→Aragonese (spa-arg), Spanish→Aranese (spa-arn), and Spanish→Asturian (spa-ast)[2]. We train a unified model for all language pairs in the **constrained** scenario. In addition, we add two language control tokens for Aragonese and Aranese Occitan, as there is already one present for Asturian.

### 1.1 Linguistic background

The Iberian peninsula - which includes the territory of Spain, Portugal and Gibraltar - is a hotspot for linguistic diversity, especially among languages in the Romance family. Spanish, Portuguese and English have official status across these three respective territories.

Basque (a non-Indo-European language) has co-official status in the Spanish Autonomous Communities of the Basque Country and the northern por-

tion of Navarre. In Galicia, Galician is co-official and in the Balaeric Islands, the Valencian Community and Catalonia Catalan/Valencian also enjoys this status.

This status ensures visibility of these languages in the socio-political space as well as a sizeable presence online. Catalan, Basque and Galician are included in many high-performing machine translation (MT) systems (and large language models (LLMs) capable of the task) (Armengol-Estapé et al., 2021) and benchmarks (Federmann et al., 2022).

This is not necessarily the case for the languages which are the focus of this challenge. They are a diverse set of languages, all from different subbranches of the Romance language family. Figure 1 shows their relation to other languages in the Romance family, and to each other, using the wave model (Heggarty et al., 2010) of linguistic evolution. Note the dialect continuum which appears to form between Portuguese → Asturian → Spanish → Aragonese → Catalan and Gascon Occitan.

Figure 2 provides a visual overview of the languages that are translated into from Spanish as part of this challenge. **Aranese**, a dialect of Gascon Occitan, also has co-official status in Catalonia but provision is only made in the Aran Valley for its use.

**Aragonese** and **Asturian** are spoken by larger numbers of people, but mostly as either second language learners or legacy speakers such as the elderly. It is for this reason that these languages all fall under the category of "Endangered" languages according to Ethnologue (Eberhard et al., 2024). All three languages are, however, considered "Vital" in terms of Digital Language Support (Simons et al., 2022). This is the second highest category behind "Thriving", meaning that there are extant corpora and resources available. However, this does not necessarily mean that there is decent quality technology such as MT available for these lan-

---

[1]Igor Kuzmin and Euan McGill are corresponding authors

[2]Submission IDs #622, #623, and #624 respectively

Figure 2: The languages involved in the WMT shared task and some demographic information

guages.

## 1.2 Extant technology for these languages

There is a recent increased push towards including languages with a small digital presence in language technology (Bapna et al., 2022), and effort has been made already to cover the languages of this challenge, including efforts to generate clean corpora from multilingual content from the internet (González and Álvarez, 2023; Ruder et al., 2023).

The first rule-based system to involve translation into and between the present languages is the open source Apertium (Forcada et al., 2011). Other systems and improvements have been built on top of this service such as Softcatalà (Ivars-Ribes and Sánchez-Cartagena, 2011) which focuses on translation into and out of Catalan, and a neural MT translator (NMT) between Spanish→Aragonese (Cortés et al., 2012).

In addition, the Spanish government-funded TAN-IBE project (Oliver et al., 2023) - of which this challenge is a part - seeks to apply modern techniques across NMT and LLM-based approaches to improve this low-resource MT task.

## 2 System description

We take the distilled NLLB-200 model (Costa-jussà et al., 2022) with 600M parameters and extend special tokens with 2 tokens that denote target languages (arn_Latn, arg_Latn) because Asturian was already presented in NLLB-200 model. After we initialized the weights of the new tokens using weights from existing tokens in the vocabulary. We used oci_Latn (Occitan) for arn_Latn (Aranese)

and spa_Latn (Spanish) for arg_Latn (Aragonese) because this languages are from the corresponding language family.

## 2.1 Training and data filtering

To create our corpus, we sampled OPUS[3] and PILAR[4] FLORES+ (revised pairs), which contain Catalan->Aranese (from PILAR), Spanish->Aranese, Spanish->Occitan, Spanish->Asturian and Spanish->Aragonese directions. We used Apertium (Khanna et al., 2021) to translate Catalan to Spanish, but we kept both source languages in our training set. Additionally, for the Occitan target language, we used idiomata cognitor (Galiano-Jiménez et al., 2024) to keep only corresponding target languages. We applied the adapted MOSES Punctuation Normalizer provided by Meta Research group under the stopes library[5] for all language pairs because NLLB was trained on preprocessed texts. Further data filtering followed the NLLB paper (Costa-jussà et al., 2022). We used fastText[6] to delete all pairs with English examples. After that, we computed length ratios and kept all sentences where the length was from five to 1050 characters, with a max length ratio lower than 0.9 and a unique ratio higher than 0.125. Finally, we de-duplicated all translation language pairs, keeping a maximum of two source duplicates and three target duplicates. Additionally we kept all pairs where distance score was in [0.6;1.0]. The result distribution of the source and target languages in

---

[3] https://opus.nlpl.eu/
[4] https://github.com/transducens/PILAR
[5] https://github.com/facebookresearch/stopes/blob/main/stopes/pipelines/monolingual
[6] https://fasttext.cc/

Figure 3: Distribution of language pairs from processed dataset.

our result corpora is captured at the Figure 3.

For the rest of the language pairs, we excluded all samples where the target language did not match the language predicted by idiomata cognitor.

## 2.2 Data augmentation

We adapt the model by training on a special regime of data augmentation with both monolingual and bilingual training data for the language pairs in this challenge.

The OPUS data were filtered in order to discard the spurious sentence pairs. We do that by performing translation of the Spanish sentence to the appropriate target language using Apertium and comparing the translation to the sentence present in the corpus. We assume that certain differences are possible due to imperfect performance of Apertium and natural variability of language, but the two variants should preserve some resemblance. To quantify that, we compute the Levenshtein (Levenshtein, 1966) edit distance ($d(s_1, s_2)$) between the two strings ($s_1$, $s_2$) and transform it into a similarity score defined as:

$$sim(s_1, s_2) = 1.0 - \frac{d(s_1, s_2)}{max(|s_1|, |s_2|)}$$

Based on manual analysis of the scores, we assume the similarity score of minimum 0.6 to be sufficient for the sentence pair to be used. Otherwise, it is discarded.

## 2.3 Fine tuning

The NLLB-200 model with 600M parameters, distilled from a 54B parameter Mixture-of-Experts model, demonstrated superior performance compared to the baseline version. Building on this

foundation, we implemented a series of adaptation steps described above to further enhance the model's capabilities on a new target languages. In this sections, we detail our training methodology and the specific hyperparameters employed to optimize the model's performance across diverse linguistic tasks. The fine-tuning process was done with one T4 GPU using Hugging Face Transformers (Wolf et al., 2020) library with the following hyperparameters presented at the Table 1. Our result model is available at the Hugging Face repository[7].

| Hyperparameter | Value |
|---|---|
| Learning Rate | 1e-4 |
| Weight Decay | 1e-3 |
| Train Batch Size | 4 |
| Eval Batch Size | 4 |
| Training Epochs | 2 |
| Optimizer | Adafactor |
| Clip Threshold | 1.0 |
| Warmup Steps | 10% of total steps |

Table 1: Hyperparameters for NLLB-200 Fine-tuning.

## 3 Results

Our results for the translation task from the Spanish language test set[8] to the target languages, as evaluated through OCELoT[9] submission system are reasonably positive, with respective BLEU and chrF+ scores of 49.2 and 73.6 for spa-arg, 17.9 and 15.5 for spa-arn, and 23.9 and 46.1 for spa-ast.

In terms of comparing the current approach with previous approaches such as Apertium and its successors, many of these studies only report word error rate whereas we used BLEU and chrF+. In those studies where BLEU is reported, it is known that BLEU favours SMT and NMT systems over rule-based ones. Moreover, this challenge introduces the present test set - so there is no previous work on the same data for direct comparison.

We find that this method of training is relatively efficient, with energy usage of 2.93kWh and emissions of approximately 1.81kg of $CO_2$[10].

---

[7]https://huggingface.co/igorktech/tribble-600m
[8]https://github.com/transducens/wmt2024-romance-tests
[9]https://ocelot-west-europe.azurewebsites.net/leaderboard/4
[10]https://wandb.ai/igorktech01/wmt24-tribble/runs/5z9r7tjt

## References

Jordi Armengol-Estapé, Casimiro Pio Carrino, Carlos Rodriguez-Penagos, Ona de Gibert Bonet, Carme Armentano-Oller, Aitor Gonzalez-Agirre, Maite Melero, and Marta Villegas. 2021. Are multilingual models the best choice for moderately under-resourced languages? A comprehensive assessment for Catalan. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4933–4946, Online. Association for Computational Linguistics.

Ankur Bapna, Isaac Caswell, Julia Kreutzer, Orhan Firat, Daan van Esch, Aditya Siddhant, Mengmeng Niu, Pallavi Baljekar, Xavier Garcia, Wolfgang Macherey, Theresa Breiner, Vera Axelrod, Jason Riesa, Yuan Cao, Mia Xu Chen, Klaus Macherey, Maxim Krikun, Pidong Wang, Alexander Gutkin, Apurva Shah, Yanping Huang, Zhifeng Chen, Yonghui Wu, and Macduff Hughes. 2022. Building machine translation systems for the next thousand languages. *Preprint*, arXiv:2205.03983.

Juan Pablo Martínez Cortés, Jim O'Regan, and Francis M Tyers. 2012. Free/open source shallow-transfer based machine translation for spanish and aragonese. In *LREC*, pages 2153–2157.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.

David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2024. Ethnologue: Languages of the World. Twenty-seventh edition.

Christian Federmann, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online. Association for Computational Linguistics.

Mikel L Forcada, Mireia Ginestí-Rosell, Jacob Nordfalk, Jim O'Regan, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Gema Ramírez-Sánchez, and Francis M Tyers. 2011. Apertium: a free/open-source platform for rule-based machine translation. *Machine translation*, 25:127–144.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, and Juan Antonio Pérez-Ortiz. 2024. Idiomata cognitor.

Antoni Oliver González and Sergi Álvarez. 2023. Filtering and rescoring the CCMatrix corpus for neural machine translation training. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 39–45, Tampere, Finland. European Association for Machine Translation.

Paul Heggarty, Warren Maguire, and April McMahon. 2010. Splits or waves? trees or webs? how divergence measures and network analysis can unravel language histories. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365(1559):3829–3843.

Xavier Ivars-Ribes and Victor M. Sánchez-Cartagena. 2011. A widely used machine translation service and its migration to a free/open-source solution: the case of softcatalà. In *Proceedings of the Second International Workshop on Free/Open-Source Rule-Based Machine Translation*, pages 61–68, Barcelona, Spain.

Tanmai Khanna, Jonathan N Washington, Francis M Tyers, Sevilay Bayatlı, Daniel G Swanson, Tommi A Pirinen, Irene Tang, and Hèctor Alòs i Font. 2021. Recent advances in Apertium, a free/open-source rule-based machine translation platform for low-resource languages. *Machine Translation*, 35(4):475–502.

Vladimir Iosifovich Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10:707–710.

Antoni Oliver, Mercè Vàzquez, Marta Coll-Florit, Sergi Álvarez, Víctor Suárez, Claudi Aventín-Boya, Cristina Valdés, Mar Font, and Alejandro Pardos. 2023. TAN-IBE: Neural machine translation for the romance languages of the Iberian peninsula. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages

495–496, Tampere, Finland. European Association for Machine Translation.

Sebastian Ruder, Jonathan Clark, Alexander Gutkin, Mihir Kale, Min Ma, Massimo Nicosia, Shruti Rijhwani, Parker Riley, Jean-Michel Sarr, Xinyi Wang, John Wieting, Nitish Gupta, Anna Katanova, Christo Kirov, Dana Dickinson, Brian Roark, Bidisha Samanta, Connie Tao, David Adelani, Vera Axelrod, Isaac Caswell, Colin Cherry, Dan Garrette, Reeve Ingle, Melvin Johnson, Dmitry Panteleev, and Partha Talukdar. 2023. XTREME-UP: A user-centric scarce-data benchmark for under-represented languages. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1856–1884, Singapore. Association for Computational Linguistics.

Gary F. Simons, Abbey L. L. Thomas, and Chad K. K. White. 2022. Assessing digital language support on a global scale. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4299–4305, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

# CloudSheep System for WMT24 Discourse-Level Literary Translation

**Lisa Liu, Ryan Liu, Angela Tsai, Jingbo Shang**
University of California, San Diego
{lil043, ryl001, cjt002, jshang}@ucsd.edu

## Abstract

This paper describes the CloudSheep translation system for WMT24 Discourse-Level Literary Translation shared task. We participated in the Chinese-English direction on the unconstrained track. Our approach to the task used a pipeline of different tools in order to maximize the translation accuracy and flow of the text by combining the strengths of each tool. In particular, our focus was to translate names consistently and idioms correctly. To achieve consistent names throughout a text, a custom name dictionary was generated for each text, containing person and place names, along with their translations. A common honorific dictionary was applied for consistency with titles, especially in historical or cultivation novels. The names were found and translated with GPT 3.5-turbo. To achieve accurate and concise translations of idioms, which are often translated literally and verbosely, we integrated the CC-CEDICT library to provide official definitions. Then, we used GPT-4 to pick the best dictionary definition that fit the context and rephrase it to fit grammatically within a sentence. For the translation of non-name and non-idiom terms, we used Google Translate. We compared our approach's performance with Google Translate as a baseline using BLEU, chrF, and COMET, as well as A/B testing.

## 1 Introduction

Machine translation techniques customized for webnovels have been researched more during the past few years (Wang et al., 2023). With the widespread availability of commercial and open-source large language models, it has become easier to fine tune existing models for a specific kind of data. Many of the top translation solutions to last year's task approach the problem of webnovel translation from the fine tuning perspective, experimenting with combining and tuning different machine learning models to find the best method for translation (Lopez et al., 2023; An et al., 2023).

When scored by human annotators, each of last year's machine translation systems, without exception, had more errors in the categories of Accuracy and Fluency compared to the other categories of Style, Terminology, Localization, and Other (Wang et al., 2023). This may indicate that inconsistency and inaccuracy are still issues that need more attention.

With a background in reading and translating webnovels as human translators, specifically in the Chinese to English direction, we wanted to approach the machine translation problem from the human readability perspective. As a reader, one of the biggest qualities of a translation is consistency. When a character is referred to as A in one sentence and referred to as B in the next, it is very hard to follow the translation, even if the writing style and vocabulary choices are immaculate. On the other hand, even if the character is wrongly referred to as B the whole time, the consistency allows the reader to follow the translation and the events. At the time of our background research, the most up-to-date version of DeepL, a popular machine translation tool in the webnovel translation community, still had name translation inconsistencies even within the same sentence, as shown in Figure 1.

Another important aspect of a good translation from a reader's perspective is correctly translating Chinese phrases with an English equivalent that matches it in tone. As a translator, that means that we often aim to convey the figurative meaning, rather than the literal meaning. This is especially common for idioms, or "chengyu" in Chinese. These phrases often originated from ancient texts, and their meaning often comes from the myth, story, or historical event they were derived from, rather than the actual characters. Due to this, a literal translation fails to convey the meaning, and often is too formal for the modern settings where they are used as a casual part of speech. For example, the phrase "脑子进水" literally translates to

Figure 1: Name inconsistencies within DeepL translation for a single passage.

"water entered the brain," but the meaning is "lost one's mind" or "gone crazy." Making the appropriate choice between them depends on the sentence's tone and context.

With these two aspects in mind, our translation system aims to target inconsistencies in name translation and inaccuracies in idiom translation. We accomplished the former through generating a dictionary of the names found in the Chinese raws along with their English translations. We accomplished the latter through finding the figurative meaning of idioms from an open-source dictionary and using GPT-4 to rephrase the best definition to fit the sentence.

## 2 Data and Tools

We primarily used the GuoFeng Webnovel Corpus provided by the organizers (WMT23, 2023) (Wang et al., 2024). The data we used for self-evaluation came from the test data in last year's dataset, because of the relatively short lengths of the texts provided per novel and the reference English translation provided as well. We also looked for short excerpts of novels through publicly available translations (found through NovelUpdates) and their original Chinese texts (found through JJWXC) to

test our system's ability to translate idioms and names.

We also used public blog posts to compile a dictionary of honorific translations, in order to maintain consistent translations across novels and texts. We used open-source dictionaries like CC-CEDICT to obtain the most accurate translation for idioms. Finally, we used prompt engineering and GPT models to tie together the different translation tools we used to create a comprehensive translation.

## 3 Translation Pipeline

### 3.1 Text Segmentation (Name Translation)

We wanted to find a way to reliably build a name dictionary that would get a majority of the names without incurring too much cost. The first place submission in last year's task's unconstrained task, DUTNLP (Zhao et al., 2023), used Jieba, a segmentation tool for Chinese. Text segmentation is the process of dividing text into meaningful words or phrases. Different segmentation granularities can significantly impact translation performance, especially for languages like Chinese (Zhao et al., 2013). In Chinese, spaces are not used to separate words, which can be made up of multiple characters, making good text segmentation very important for determining which words are present in a sentence.

We tested Jieba in our own system, aiming to use its ability to identify proper nouns to form a basic dictionary of names in the text. Specifically, we filtered for phrases tagged "nr" (person name) and "ns" (place name) (Jieba, 2020). Unfortunately, Jieba had a high false positive rate, and often split up phrases or names, which made it unsuitable to form our name dictionary. For example, if the name contained a common noun that could be part of many phrases, replacing that part of the name with the English meaning would be very unhelpful and create a weird-sounding name. However, although Jieba was not suitable for identifying proper nouns, it was still useful for determining a phrase's part of speech.

A name dictionary's main purpose is to translate names consistently, and is more useful when it contains names that appear often. If a character or place appears only once within the story, readers do not need a consistent translation across mentions to recognize it. Additionally, it is likely to be insignificant to the story, so even if the translation is not the best, it is unlikely to affect enjoyment much.

Names commonly occurring in the text are likely to be re-occurring characters, such as the main protagonist or important supporting characters.

We decided to try to feed a percentage of lines to GPT-3.5 for name identification. GPT-3.5 was good at identifying names, rarely returning false positives. However, we didn't need to get every single name from a text, just the re-occurring ones. This meant that GPT was sifting through a large number of duplicates, and incurring extra cost through the API.

We manually identified the names within the sample dataset from last year's task, and for each text, we calculated the total number of unique names, the total percentage of characters within the text that belonged to a name, and the average number of lines that would contain a name. We found that the character percentages ranged from 5% to 10%, and the line percentages ranged from 11% to 22%.

By only giving GPT a certain percentage of lines that were randomly selected from the text, we introduced an element of chance into our pipeline that meant GPT may not be able to see all the names from the text it is given. We selected 20% as a number on the higher end of the range we found, so it was likelier that GPT would be given a majority of the names.

In order to pick lines more likely to contain names, we used Jieba to identify the number of nouns within a line. We theorized that it was unlikely for lines to contain no names, so if Jieba didn't identify any nouns at all within a sentence, it likely missed a name, which may be a combination of characters that are verbs or adjectives on their own. We first ran Jieba's segmentation on the text, and then selected only from a pool of longer sentences without any nouns identified.

We also theorized that for characters such as these, their introductions are more likely to be concentrated within the beginning or middle of the text, rather than the end. As we only need to get one occurrence of each name, we decided to weigh sentences earlier as more likely to be selected. We picked 15% of the lines from the first $\frac{3}{4}$ of the text, and 5% of the lines from the last $\frac{1}{4}$ of the text.[1]

---

[1] "Lines" in the text file are sometimes multiple sentences in the Chinese raws; so if Jieba identifies 0 nouns in a "line", that can equate to 0 nouns in a paragraph.



Figure 2: Flow chart describing name translation process.

## 3.2 Honorifics (Name Translation)

We compiled a list of honorifics, ranging from common honorifics such as "哥哥" (brother) to martial art novel honorifics such as "师爷" (grandmaster) (Mountain, 2017) to historical novel honorifics such as "公公" (eunuch) (Wyhcwe, 2022). We acknowledge that the translations of such terms can sometimes vary across different translations, but we wanted to make a standard translation across all of our translations.

The names identified by Jieba and GPT-3.5 in the previous subsection include these honorifics, so by first applying the honorific translation and only asking GPT to translate the remaining characters left behind as the name, it can standardize the name translations and also ensure the honorifics aren't translated as pinyin directly. For example, our code would go through these steps to translate a name: 韩少爷→ Young Master 韩→ Young Master Han. We used the prompt: "Translate this name to English: [name]. Only list the English name."

## 3.3 CC-CEDICT (Idiom Translation)

We used the Chinese-English dictionary, CC-CEDICT. It is an free online dictionary that is regularly updated through crowdsourcing, and every contribution is verified regularly and added to

Figure 3: Flow chart describing idiom replacement process.



Figure 4: Flow chart describing final rephrasing process.

the database (CC-CEDICT, 2020). Due to continual updates by the owners, CC-CEDICT is a good choice for getting the most updated figurative meanings of idioms, slang, and other culturally specific terms. For example, the phrase "脑子进水" from the introduction section has the CC-CEDICT entry "to have lost one's mind crazy soft in the head."

We searched through all entries labelled as "idioms" within the CC-CEDICT dictionary. If any such idioms were found within a line of the original Chinese text, the Chinese idiom would be replaced directly with its corresponding CC-CEDICT English entry. Other Chinese text in the line not identified as idioms would not be translated at this step. The raw dictionary replacements did not account for grammatical context surrounding the idioms, and some entries contained more than one English translation phrase per Chinese idiom. Furthermore, in their raw formatting these entries were surrounded by brackets and contained a text flag "(idiom)". We kept the raw replacement formatting as-is, which we then processed further after translating the rest of the text.

### 3.4 Overall Translation

We were left with text that was primarily still in the original Chinese, but with names and idioms programmatically replaced with English translations. We experimented with two different translation engines, DeepL and Google Translate, to translate the remainder of the text. These engines are the two most mentioned translation engines amongst online webnovel forums before ChatGPT. The en-

gines translated the remaining Chinese text without any modification to the already-present English idiom replacements, thus requiring a step to smooth out the sentences containing idioms.

### 3.5 GPT Rephrasing (Idiom Translation)

We decided to use GPT-4 and LangChain (LangChain, 2024) to replace every line that contained a raw idiom definition, as identified by their surrounding brackets and accompanying "(idiom)" flag, with a grammatically correct rephrasing. Langchain is an open-source framework that makes it easier to develop using GPT's API. We found that GPT-4 was better than GPT-3.5 at rephrasing only the sentences with idiom definitions within a given line. Because any output from this section would be inserted directly into the text as the final step, we decided to switch to GPT-4 for this step for better quality. To minimize the API cost in the rephrasing phase, we recorded the line numbers for the lines modified in the CC-CEDICT step, and only gave those lines to Langchain. Only about 2%-11% of lines across the samples we encountered contained idioms, so GPT-4 was only used on a small percentage of the text.

We used the prompt: "Please pick the idiom definition that best fits the context for the following sentences and rephrase only the part of the sentence with the idiom grammatically. Only output the new translation. Don't change the sentences without idioms. Favor the more concise meaning and find an English equivalent if possible." We added many instructions to our prompt as a result of experimentation; not asking for the "more concise meaning" or "English equivalent" often resulted in translations that were complicated amalgamations of every definition provided by the dictionary entry; not asking for "don't change the sentences without idioms" often resulted in sentences without idioms being changed and other content given being cut out.

Once the GPT-4 rephrasing was complete, the text translation was considered to be finished.

## 3.6 Evaluation

We used three metrics for automatically evaluating machine-translated text: BLEU, chrF, and COMET. BLEU evaluates word-level n-grams, calculating the precision between the machine translation and the reference, weighted by a brevity penalty (Papineni et al., 2002). ChrF evaluates character-level n-grams, scoring the overlap of short sequences of characters between the machine translation and the reference (Popović, 2015). COMET is a fine-tuned neural framework that takes in sentence embeddings from the source text, translation, and reference (Rei et al., 2020). We used these because last year's conference proceedings summary paper used them for the automatic evaluation (Wang et al., 2023).

A shortcoming of automatic metrics such as BLEU is that they lack the ability to evaluate based upon semantics, instead favoring direct word-to-word matches between a translation and reference (Callison-Burch et al., 2006). This means a translation that achieves high grammatical quality but uses different words than a provided reference could potentially score poorly. As such, we also surveyed human readers to compare the quality of our system's translations. Participants were given 4 separate translations of a text sample ranging from 200-300 English words, each generated using a different method: one generated by our translation system using Google Translate ("pipeline Google Translate"), one generated by our translation system using DeepL ("pipeline DeepL"), one generated using only Google Translate ("pure Google Translate"),



Figure 5: random sample 1, video games (20%)



Figure 6: random sample 2, science fiction (23%)

and one using only DeepL ("pure DeepL"). These translations were given in a random order, and participants were not informed of which translation came from which source. After reading the translations, participants were asked to rank them from best to worst based on how readable they found the translations. This process was repeated over several different samples, and the rankings were recorded for each sample.

## 4 Results

We used the automated metrics to evaluate the results of the four techniques: pure Google Translate, pipeline Google Translate, pure DeepL, pipeline DeepL. To decide between Google Translate and DeepL for our final submission, we decided to compare the pure Google Translate and pure DeepL results. In this paper, we show the results for three random samples selected across the dataset, shown in Figures 5, 6, and 7. The genre of the sample is labelled, along with their distribution percentage in the training set (Wang et al., 2023). Google Translate and DeepL performed about the same for the first two samples, but Google Translate was signifi-

Figure 7: random sample 3, martial arts (2%)

cantly better than DeepL in the third sample, which was a classical martial arts novel. Although martial arts novels only make up a small percentage of the dataset, because idioms originate from classical Chinese literature, we decided to employ the translation pipeline with Google Translate ("pipeline Google Translate") for our final conference submission.

In our A/B testing, across all the samples, we found that participants ranked the pipeline Google Translate output the highest most, and the pure DeepL the lowest, as shown in Table 1. However, the distribution was mostly even, and about half the time, participants reported that the difference between translations was slight, which could be due to the limitations mentioned in the following limitations section.

| Technique | 1st | 2nd | 3rd | 4th |
|---|---|---|---|---|
| pure Google (2) | 2 | 2 | 1 | 2 |
| pipeline Google (1) | 3 | 1 | 3 | 0 |
| pure DeepL (4) | 1 | 3 | 0 | 3 |
| pipeline DeepL (3) | 2 | 0 | 4 | 1 |

Table 1: Times each technique was ranked 1st, 2nd, 3rd, or 4th across 7 samples. Ties were allowed.

## 5 Conclusion

We created a machine translation system that creates consistent translations for names and accurate translations for idioms, both of which enhance human readability despite making up a small ratio of the overall text. Even though our pipeline did not see any major improvements in the automated evaluation metrics, the positive reception among human survey participants points to the potential value that our process provides.

## 6 Limitations

When providing lines of text for ChatGPT to identify names, we randomly selected a certain percentage of lines to use in order to reduce API usage costs. Though the selected lines were weighted based on factors such as whether or not Jieba found any proper nouns in a line, there is nonetheless a slight element of randomness that is introduced during our process. One limitation that could be further explored is how consistently our pipeline performs over multiple runs on the same input.

Another limitation in our results lay in our use of human evaluators. Participants were asked to rank translations that used our system against translations that did not. Though they were not informed which translations did or did not use our system, they also were not given any specific metrics to quantify their decisions. Participants also sometimes reported that the passages provided were too long to quickly judge the difference, and that reading four passages in a row that described the same content made it hard to evaluate the difference without an earnest effort to study the differences within the text. In the future, our team could work on developing a more robust approach to the human side of evaluations that addresses these limitations.

## References

Li An, Linghao Jin, and Xuezhe Ma. 2023. Max-isi system at wmt23 discourse-level literary translation task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 282–286.

Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. Re-evaluating the role of bleu in machine translation research. In *11th conference of the european chapter of the association for computational linguistics*, pages 249–256.

CC-CEDICT. 2020. Chinese-English Dictionary. https://cc-cedict.org/wiki/. [Online; accessed 28-July-2024].

Jieba. 2020. Jieba. https://github.com/fxsjy/jieba. [Online; accessed 30-July-2024].

LangChain. 2024. LangChain. https://www.langchain.com/langchain. [Online; accessed 30-July-2024].

Fabien Lopez, Gabriela González-Sáez, Damien Hansen, Mariam Nakhlé, Behnoosh Namdarzadeh, Marco Dinarelli, Emmanuelle Esperança-Rodier, Sui He, Sadaf Mohseni, Caroline Rossi, et al. 2023. The make-nmtviz system description for the wmt23 literary task.

Immortal Mountain. 2017. WuXia Terms of Address. https://immortalmountain.wordpress.com/glossary/terms-of-address. [Online; accessed 28-July-2024].

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Maja Popović. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. *arXiv preprint arXiv:2009.09025*.

Longyue Wang, Siyou Liu, Minghao Wu, Wenxiang Jiao, Xing Wang, Jiahao Xu, Zhaopeng Tu, Liting Zhou, Yan Gu, Weiyu Chen, Philipp Koehn, Andy Way, and Yulin Yuan. 2024. Findings of the wmt 2024 shared task on discourse-level literary translation. proceedings of the ninth conference on machine translation (wmt).

Longyue Wang, Zhaopeng Tu, Yan Gu, Siyou Liu, Dian Yu, Qingsong Ma, Chenyang Lyu, Liting Zhou, Chao-Hong Liu, Yufeng Ma, et al. 2023. Findings of the wmt 2023 shared task on discourse-level literary translation: A fresh orb in the cosmos of llms. *arXiv preprint arXiv:2311.03127*.

WMT23. 2023. WMT23 Task Link. http://www2.statmt.org/wmt23/literary-translation-task.html. [Online; accessed 23-July-2024].

Wyhcwe. 2022. Historical Terms of Address. https://dreamsofjianghu.ca/%e5%85%ab%e5%ae%9d%e5%a6%86-eight-treasures-trousseau/glossary/. [Online; accessed 28-July-2024].

Anqi Zhao, Kaiyu Huang, Hao Yu, and Degen Huang. 2023. DUTNLP system for the WMT2023 discourse-level literary translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 296–301.

Hai Zhao, Masao Utiyama, Eiichiro Sumita, and Bao-Liang Lu. 2013. An empirical study on word segmentation for chinese machine translation. In *International Conference on Intelligent Text Processing and Computational Linguistics*, pages 248–263. Springer.

# SJTU LoveFiction's System for WMT24 Discourse-Level Literary Translation

**Haoxiang Sun[1]*  Tianxiang Hu[1]  Ruize Gao[2]  Jialong Tang[2]  Pei Zhang[2]**
**Baosong Yang[2]  Rui Wang[1]**
[1]Shanghai Jiao Tong University, Shanghai, China
[2]Alibaba Group, Hangzhou, China
{sunny_sjtu,hutianxiang,wangrui12}@sjtu.edu.cn
{xiaoyi.zp,yangbaosong.ybs}@alibaba-inc.com

## Abstract

This paper describes Shanghai Jiao Tong University (SJTU LoveFiction) Discourse-Level Literary Translation systems for the WMT24 shared task. We participate in the literary translation task on Chinese → English, Chinese → German and Chinese → Russian with unconstrained tack. Our system is based on Qwen2-72B(Yang et al., 2024), Claude3.5(Anthropic, 2023) and GPT-4o(OpenAI, 2024) with novel techniques that improve literary translation performance on the target language pairs. (1) Chunk-based SFT and inference: we put several sentences together to form a chunk and try different chunksize during SFT and inference. (2) Merge multi-model translations by agents: we design a Translation Editor Agent based on GPT-4o to generate a better new translation by referencing the source text and merge 3 candidate translations generated by Qwen2-72B, Claude-3.5 and GPT-4o. (3) Terminology Intervention: to ensure terminology consistency, a Term Proofreader Agent, based on GPT-4o, is utilized to extract term pairs from source texts and translations. For each Chinese term, we decide its optimal translation and request the Term Proofreader to modify the translation generated by Translation Editor Agent. In model evaluation:(1) We employ d-BLEU for single model evaluation. (2) We design a Client Agent based on Claude-3.5 to assess the win-tie rate between two translations for cross-model evaluation.

## 1 Introduction

Despite great advancements in machine translation (MT) these years(Artetxe et al., 2017; Wang et al., 2022), achieving high-quality translations for literary texts remains a formidable task, primarily due to the complexities involved in maintaining coherence, consistency, and cultural context across larger text spans (Voita et al., 2019; Lopes et al., 2020).

This paper describes SJTU LoveFiction's submission to WMT24 Discourse-Level Literary Translation. We participate in all 3 language pairs (Chinese → English, Chinese → German and Chinese → Russian) with unconstrained tack.

Our system builds upon Qwen2-72B, Claude-3.5 and GPT-4o models with various practical techniques. We adopt a chunk-based strategy, grouping several sentences into a chunk during supervised fine-tuning (SFT) and inference phase.

Multi-agent structure demonstrates strong performance in discourse-level machine transaltion(Wu et al., 2024). To enhance translation quality, we develop a Translation Editor Agent based on GPT-4o. This agent references the source text and merges muliti-model translations to produce a refined output. While different models may generate varied translations for the same Chinese term, we also implement a Term Proofreader Agent powered by GPT-4o. This agent extracts term pairs from source text and corresponding translations. For each Chinese term, the optimal translation is determined manually, then the term proofreader applies these optimal terms to the merged translations.

In terms of evaluation, we use d-BLEU to assess the performance of a single model under different experimental settings. For cross-model evaluation, we design a Client Agent based on Claude-3.5. This agent references the Chinese source text to evaluate and rank the translations produced by different models by accuracy, fluency, and the preservation of stylistic elements.

This paper is structured as follows: Section 2 describes our data pre-processing strategies, followed by the details of our method in Section 3. Section 4 presents the experimental results and analysis, then we draw conclusions in Section 5.

---

*Work done during internship at Alibaba Group

## 2 Data Processing

We perform Supervised Fine-Tuning (SFT) on the GuoFeng Webnovel Corpus (Wang et al., 2023) . Handling the noise within the dataset is crucial as it can significantly impact LLM's translation performance. We adopt a series of rigorous data filtering strategies.

### 2.1 Chinese-English Data Filtering

1. **Remove lines without Chinese-English pairs**: Delete any line that contains only a single Chinese or English sentence.

2. **Eliminate garbled text, emojis, foreign language characters, and emoticons**: These elements can degrade model performance. We use Unicode range identification and regular expressions for precise removal.

3. **Delete lines containing only punctuation marks**: Such lines typically lack linguistic value and retaining them would introduce noise, thereby impairing model training.

4. **Standardize punctuation**: Convert all Chinese punctuation to English punctuation to enhance model consistency and coherence in translation results.

### 2.2 Chinese-German and Chinese-Russian Data Filtering

Chinese-German/Russian data has the following features.

1. **Chapter-Level Alignment Only**: The alignment is maintained only at chapter level. Within chapters, paragraph or sentence level alignment is not achieved.

2. **Chapter Containment Differences**: In the Chinese files, each file contains a single chapter. In contrast, the German and Russian files may contain multiple chapters per file.

The following filtering strategies are employed:

1. **Remove Unaligned Chapter Pairs**: Delete Chinese-German/Russian file pairs that are not aligned at the chapter level.

2. **Eliminate garbled text, emojis, foreign language characters, and emoticons.**

3. **Remove Chapters Exceeding 8k Tokens**: LLMs struggle with long passages, thus chapters exceeding 8k tokens are excluded.

## 3 Method

In this section, we describe our method and provide a comprehensive explanation of the key components.

### 3.1 System Overview

We depict the overview of our system in Figure 1, which can be divided into four steps:



Figure 1: System Overview

1. **Chunk Division**: To maintain contextual information, we combine several sentences into a single chunk.

2. **Supervised Fine-Tuning & 1-shot Inference for Multi-Model Translator**: We SFT the Qwen2-72B model on Guofeng Webnovel Corpus. Afterwards, we use the fine-tuned Qwen2-72B, Claude-3.5, and GPT-4o to perform 1-shot inference on the test set, generating translation results.

3. **Translation Merging**: We employ a Translation Editor Agent based on GPT-4o to merge the translation outputs of the three models.

4. **Terminology Intervention**: We utilize a Term Proofreader Agent based on GPT-4o to extract term glossaries from source texts and translations. We select the optimal term pairs manually and ask the term proofreader to apply them to the merged translation as the final output.

### 3.2 Chunk Division

The lack of contextual information in sentence-level data poses a significant challenge for achieving high-quality translation results. Combining multiple sentences within each chapter into chunks can alleviate this problem(Zhao et al., 2023). During the SFT phase, we experiment with various

chunksizes of 5, 10, and 20 sentences to determine the optimal size for training. In the inference phase, we further extend our experiments to chunksizes of 1, 5, 10, 20, 40, and 80 sentences. This strategy aims to provide the model with more contextual information, thereby improving the translation quality.

### 3.3 Supervised Fine-Tuning (SFT) & 1-shot Inference

In order to find the best setting for Qwen2-72B, we SFT Qwen2-7B on the Guofeng Webnovel Corpus and conduct inference on the in-domain dev set. Given consistent distribution between the two datasets, this approach will reveal the best setting for LLM to learn the knowledge embedded in Guofeng Webnovel Corpus. d-BLEU scores under different settings are shown in Figure 2.



Figure 2: d-BLEU for Qwen2-7B on In-domain Dev Set

Although training with 5-sentence chunks and inferring with 1-sentence chunks yields the highest d-BLEU score of 29.35, we prefer 10-sentence chunks for training and 5-sentence chunks for inference. This configuration, with a d-BLEU score of 29.17, maintains nearly equivalent performance while preserving contextual information during inference.

As we aim to capture more context, the model must handle longer inputs. However, LLM's ability to handle long inputs is inherently limited. It's essential to acknowledge that we need to strike a balance, i.e. **maintaining sufficient contextual information without exceeding the model's capacity for processing long inputs.**

In our inference experiments with Claude-3.5 and GPT-4o, we employed 1-shot inference, a form of few-shot learning. Few-shot learning aims to enable models to generalize from a limited number of examples(Brown et al., 2020).We determine the best inference chunksize according to the following results.



Figure 3: d-BLEU for Claude-3.5 & GPT-4o on OOD Dev Set

Both models perform best with 5-sentence chunk size, achieving d-BLEU scores of 19.98 and 21.15, respectively. We choose 5-sentence chunksize as the inference setting for Claude-3.5 and GPT-4o.

### 3.4 Translation Merging

After obtaining multi-model translations, we randomly select a chapter for manual verification and observe that different models exhibit distinct strengths in their translations for the same chunk. **To leverage the advantages of all three translations, we employ a Translation Editor agent based on GPT-4o, which is prompted to merge the three candidate translations into an improved version.** Workflow of the Translation Editor is as follows.



Figure 4: Workflow of The Translation Editor

1. **Quality Assessment.** Assess the quality of different translation referencing the source text. After this step, the agent knows the relatively better part in each translation.

2. **Translation Merging.** Put these parts together to form the merged translation.

This process allows the Translation Editor agent to integrate the best elements (highlighted in red in Figure 4) of the three candidate translations, generating a superior translation.

### 3.5 Terminology Intervention

While the Translation Editor agent generates improved results by blending three candidate translations, **different models may produce different translations for the same Chinese term, leading to consistency issues.** To address this, we develop a **Term Proofreader Agent**. Workflow of the agent is as follows.



Figure 5: Workflow of The Term Proofreader

1. **Term Extraction.** The terminology proofreader agent begins by extracting term pairs from the Translation Editor's output, referencing its Chinese source. Glossaries are obtained after this step.

2. **Manual Determination.** For each Chinese term in the glossaries, we manually determine the optimal translation. This step involves reviewing the context and ensuring that the chosen translation accurately reflects the meaning and nuance of the original term.

3. **Term Application.** Once the optimal translations are determined, the terminology proofreader agent applies these optimal translations to Translation Editor's output.

### 3.6 Evaluation

#### 3.6.1 Single Model Evaluation

We calculate d-BLEU scores between our translations and reference texts to evaluate single model performance and determine the optimal experimental settings (i.e. training & inference chunksize).

d-BLEU measures N-gram matching, reflecting the similarity between two distributions.The distribution of the in-domain dev set and the train set are consistent. Thus d-BLEU can assess the model's learning of train set during SFT stage, enabling us to select the optimal SFT setting by d-BLEU. On the other hand, distribution of the ood dev set is inconsistent with the train set. d-BLEU can assess the model's fitting to the ood dev set distribution. Thus we can select the optimal inference setting by d-BLEU.

#### 3.6.2 Cross-model Evaluation

For cross-model evaluation, we find that human-preferred translation can have low d-BLEU score. This discrepancy arises because d-BLEU relies solely on N-gram matching and is unable to capture deeper semantic information. For human-preferred translation, there can be significant lexical differences from the reference translations, even though the semantic content is accurately conveyed. d-BLEU is ineffective in evaluating such cases.

Previous works reveal that LLM-Evaluators can achieve high consistency with human expert on system-level evaluation(Kocmi and Federmann, 2023; Moosa et al., 2024). We build a Client Agent based on Claude-3.5, which considers accuracy, fluency, and the preservation of stylistic elements.



Figure 6: Workflow of The Client

#### 3.6.3 Human Evaluation

We employ 3 language experts to do fine-grained evaluation. They are requested to perform Linguistic Quality Rating (LQR) by the following standard in Table 1.

## 4 Results

We present the effect of our method in this section.

| Score | Quality Description |
|-------|---------------------|
| 1 | Incomprehensible or incorrect. |
| 2 | Severe errors, hard to understand. |
| 3 | Some errors, but understandable. |
| 4 | Mostly correct, minor errors. |
| 5 | Completely correct and fluent. |

Table 1: LQR Scoring Standards

### 4.1 Supervised Fine-Tuning (SFT) & Inference

We train Qwen2-72B on GuoFeng Webnovel Corpus with 10-sentence chunksize. The following table shows the d-BLEU scores for various inference chunksizes on both in-domain and out-of-domain dev sets.

| Inference Chunksize | In-domain | OOD |
|---------------------|-----------|-----|
| 1 | 27.32 | 22.51 |
| 5 | 27.05 | 24.64 |
| 10 | 26.05 | **24.74** |
| 20 | 27.74 | 24.65 |
| 40 | **28.11** | 24.25 |
| 80 | 20.49 | 24.04 |

Table 2: d-BLEU of Qwen2-72B on In-domain & OOD Dev Set

Qwen2-72B achieves best performance under 40-sentence inference chunksize on in-domain dev set while the best performance on OOD dev set is achieved with the 10-sentence chunksize. This indicates that although Qwen2-72B has a stronger capability for handling long texts, out-of-domain data distribution still poses difficulties for translation.

Results for Claude-3.5 and GPT-4o on OOD dev set is in Figure 3.

### 4.2 Translation Merging

We randomly selected 200 chunks from the final test set to evaluate the performance of individual models and our translation merging strategy.

GPT-4o ranks 1st place in single model performance while our translation merging strategy surpasses every single model, indicating that better translation is generated by the Translation Editor Agent.

| Model | LQR3 | LQR4 | LQR5 |
|-------|------|------|------|
| **Translation Merging** | **65%** | **44%** | **24%** |
| **GPT-4o** | 60% | 30% | 9% |
| **Claude-3.5** | 54% | 33% | 12% |
| **Qwen2-72b** | 42% | 24% | 3% |

Table 3: LQR Scores for Different Models

We also employed the Client Agent to compare GPT-4o's results and the merged translations. Table 4 presents the win-tie rate relative to GPT-4o.

| Metric | Rate |
|--------|------|
| Win | 41% |
| Tie | 21% |
| Lose | 38% |
| Net Win Rate | **3%** |

Table 4: Win-tie Rate Compared to GPT-4o

The LLM evaluator also acknowledges that our translation merging strategy brings a slight improvement.

### 4.3 Terminology Intervention

We employ the Term Proofreader Agent to extract term pairs from the entire test set. The following table presents the results before and after the terminology intervention.

| | Before | After |
|---|--------|-------|
| Chinese Terms | 806 | 806 |
| English Translations | 3012 | 902 |
| Average Correspondence | 3.73 | 1.12 |

Table 5: Term Correspondence Before and After Intervention

Before the intervention, 806 unique Chinese terms correspond to 3012 English translations, with an average of 3.73 English translations per Chinese term, indicating high variability and inconsistency.

After the intervention, the number of English translations is reduced to 902. This significant reduction demonstrates that the Term Proofreader Agent effectively standardized the terminology, ensuring consistent translations for each Chinese term.

# 5 Conclusion

Through chunk splitting, multi-model translation merging, and terminology intervention, our system demonstrates strong performance in the WMT24 Discourse-Level Literary Translation task. The translation merging strategy surpasses all individual models in LQR scores. Terminology intervention significantly improves terminology consistency, reducing the average correspondence from 3.73 translations to 1.12. Future work will focus on further optimizing these techniques and exploring new strategies to enhance translation quality, especially in handling long texts and preserving literary styles.

# References

Anthropic. 2023. Claude 3.5.

Mikel Artetxe, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. 2017. Unsupervised neural machine translation. *arXiv preprint arXiv:1710.11041*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.

António Lopes, M. Amin Farajian, Rachel Bawden, Michael Zhang, and André F. T. Martins. 2020. Document-level neural MT: A systematic comparison. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 225–234, Lisboa, Portugal. European Association for Machine Translation.

Ibraheem Muhammad Moosa, Rui Zhang, and Wenpeng Yin. 2024. Mt-ranker: Reference-free machine translation evaluation by inter-system ranking. In *The Twelfth International Conference on Learning Representations*.

OpenAI. 2024. Gpt-4o.

Elena Voita, Rico Sennrich, and Ivan Titov. 2019. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1198–1212, Florence, Italy. Association for Computational Linguistics.

Haifeng Wang, Hua Wu, Zhongjun He, Liang Huang, and Kenneth Ward Church. 2022. Progress in machine translation. *Engineering*, 18:143–153.

Longyue Wang, Zhaopeng Tu, Yan Gu, Siyou Liu, Dian Yu, Qingsong Ma, Chenyang Lyu, Liting Zhou, Chao-Hong Liu, Yufeng Ma, et al. 2023. Findings of the wmt 2023 shared task on discourse-level literary translation: A fresh orb in the cosmos of llms. In *Proceedings of the Eighth Conference on Machine Translation*, pages 55–67.

Minghao Wu, Yulin Yuan, Gholamreza Haffari, and Longyue Wang. 2024. (perhaps) beyond human translation: Harnessing multi-agent collaboration for translating ultra-long literary texts. *Preprint*, arXiv:2405.11804.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Anqi Zhao, Kaiyu Huang, Hao Yu, and Degen Huang. 2023. Dutnlp system for the wmt2023 discourse-level literary translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 296–301, Singapore. Association for Computational Linguistics.

# Context-aware and Style-related Incremental Decoding framework for Discourse-Level Literary Translation

**Yuanchang Luo, Jiaxin Guo, Daimeng Wei, Hengchao Shang, Zongyao Li,**
**Zhanglin Wu , Zhiqiang Rao, Shaojun Li, Jinlong Yang, Hao Yang**
Huawei Translation Service Center, Beijing, China
{luoyuanchang1,guojiaxin1,weidaimeng,shanghengchao,lizongyao,
wuzhanglin2,raozhiqiang,lishaojun18,yangjinlong7,yanghao30}@huawei.com

## Abstract

This report outlines our approach for the WMT24 Discourse-Level Literary Translation Task, focusing on the Chinese-English language pair in the Constrained Track. Translating literary texts poses significant challenges due to the nuanced meanings, idiomatic expressions, and intricate narrative structures inherent in such works. To address these challenges, we leveraged the Chinese-Llama2 model, specifically enhanced for this task through a combination of Continual Pre-training (CPT) and Supervised Fine-Tuning (SFT). Our methodology includes a novel Incremental Decoding framework, which ensures that each sentence is translated with consideration of its broader context, maintaining coherence and consistency throughout the text. This approach allows the model to capture long-range dependencies and stylistic elements, producing translations that faithfully preserve the original literary quality. Our experiments demonstrate significant improvements in both sentence-level and document-level BLEU scores, underscoring the effectiveness of our proposed framework in addressing the complexities of document-level literary translation.

## 1 Introduction

Machine Translation (MT) (Brown et al., 2020; Chowdhery et al., 2023; Touvron et al., 2023) has become an essential tool in breaking language barriers, enabling the automatic translation of text from one language to another. While significant advancements (Vaswani et al., 2017; Sennrich et al., 2016; Wei et al., 2023; Gu et al., 2018; Ghazvininejad et al., 2019; Wang et al., 2021; Guo et al., 2021; Yu et al., 2021) have been made in MT for various text genres, translating literary texts remains a formidable challenge. Literary texts are rich in complex linguistic phenomena, such as nuanced meanings, idiomatic expressions, and intricate narrative structures. Unlike technical or news-related texts, literary works demand a deeper understanding of context, tone, and style, making them particularly challenging for MT systems. This difficulty is compounded by the scarcity of high-quality parallel datasets in the literary domain, limiting the ability of MT models to learn from extensive, diverse examples.

Document-level translation (Sun et al., 2020; Du et al., 2024; Wu et al., 2024) introduces another layer of complexity to MT, especially when dealing with longer texts such as novels. Unlike sentence-level translation, where context is limited to a single sentence, document-level translation requires the model to consider the broader discourse context to maintain coherence and consistency throughout the entire text. This is particularly crucial in literary translation, where the narrative thread, character development, and thematic elements must be preserved across sentences and paragraphs. Long-range dependencies, where information introduced early in a text influences later parts, pose a significant challenge for MT systems, which often struggle to retain and apply such context effectively over extended texts.

In this system report, we describe our participation in the WMT24 Discourse-Level Literary Translation Task, focusing on the Chinese-English language pair under the Constrained Track. Our approach leverages the Chinese-Llama2 model, specifically designed for this task, through a combination of Continual Pre-training (CPT) and Supervised Fine-Tuning (SFT). This methodology allows us to refine the model's understanding of literary texts while adapting it to the specific nuances of Chinese-English translation. Additionally, we employ an Incremental Decoding framework, which enables the model to translate documents sentence by sentence, ensuring that each translation is informed by the broader context. This approach is designed to tackle the challenges of document-level literary translation, aiming to produce translations

973

Figure 1: The overall of our approach.

that are not only accurate but also faithful to the original text's literary quality.

## 2 Background: TP3

Machine Translation (MT) is the automated process of converting text from one language to another using computational methods. Traditionally, MT relies on encoder-decoder models, where the encoder processes the source language and the decoder generates the translation, often requiring large bilingual datasets and data augmentation to improve performance. Recently, Large Language Models (LLMs) like GPT have become prominent in MT, enabling translation through zero-shot or few-shot learning by conditioning on a source sentence (Jiao et al., 2023; Zeng et al., 2023; Chen et al., 2023; Xu et al., 2023; Yang et al., 2023; Zhang et al., 2023). These models can also be fine-tuned with high-quality bilingual data and tailored instructions to enhance translation accuracy and robustness, offering new possibilities for MT with limited resources.

**TP3** Guo et al. (2024) propose a novel training paradigm, consisting of Three-Stages Translation Pipeline (TP3), to boost the translation capabilities of LLMs. The training paradigm includes:

Stage 1: Continual Pre-training using Extensive Monolingual Data. This stage aims to expand the multilingual generation capabilities of LLMs. While it is inherently related to machine translation tasks, it is not essential.

Stage 2: Continual Pre-training with Interlinear Text Format Documents. They construct interlinear text format from sentence-aligned bilingual parallel data and utilize them for continual pre-training of LLMs. Experimental results demonstrate the critical importance of this stage, resulting in a significant improvement in translation quality, particularly for English-Other translations.

Stage 3: Leveraging Source-Language Consistent Instruction for Supervised Fine-Tuning. In this stage, they discover that setting instructions consistent with the source language benefits the supervised fine-tuning process.

## 3 Methods

### 3.1 TP3 for Discourse-Level Literary

We introduce the TP3 training paradigm into the literary translation task, with the entire training process illustrated in Figure 1.

**Stage 1: Continual Pre-training using Chinese and English Monolingual Literary Data** In this stage, we adapt a general-purpose large language model (LLM) into a specialized Literary LLM by using monolingual literary data in both Chinese and English. While existing LLMs like Llama perform well in English-centric tasks, their capabilities in other languages, especially in literary contexts, are often limited. To improve this, we employ continual pre-training with extensive monolingual literary texts, enhancing the model's understanding of nuanced language, stylistic elements, and narrative structures. This step is critical for enabling the model to generate more coherent and contextually appropriate translations.

For this task, continual pre-training is essential, transforming a general LLM into one tailored for

Figure 2: The overall of our incremental decoding framework.

literary translation. We treat each novel as a distinct training unit, combining sentences within each chapter into paragraphs to capture long-range dependencies and context. This approach is vital for maintaining consistency and preserving the literary quality of translations. By focusing on both Chinese and English literary data, the model gains a balanced understanding of the stylistic and structural intricacies in both languages.

**Stage 2: Continual Pre-training with Aligned Chinese-English Interlinear Text Format Literary Documents** In Stage 2, we enhance the model's cross-lingual translation capabilities by using aligned Chinese-English interlinear text format literary documents, building on the foundation established in Stage 1. The interlinear text format, where each source sentence is directly aligned with its translation at the word or phrase level, is essential for enabling the model to understand and map the syntactic and semantic structures between Chinese and English, which is crucial for producing high-quality translations. We implement a continual pre-training approach using LoRA (Low-Rank Adaptation of Large Language Models) (Hu et al., 2021) to efficiently adapt the model with these interlinear text documents.

Initially, the model was trained on general sentence-aligned parallel data to establish a strong cross-lingual alignment foundation. Subsequently, we performed incremental pre-training

with literary-specific interlinear data. By focusing on literary documents, we ensure the model becomes finely attuned not only to general cross-lingual translation but also to the unique stylistic and structural nuances of literary texts. This approach enables the model to capture the intricate relationships between Chinese and English in a literary context, significantly improving translation quality and fidelity.

**Stage 3: Supervised Fine-Tuning with Context-aware and Style-related Instructions** In the final stage of our approach, we conduct supervised fine-tuning using context-aware and style-related instructions, specifically tailored to address the challenges of semantic coherence and stylistic consistency in literary translation. Unlike the traditional approach of using Source-Language Consistent Instruction, which emphasizes alignment with the source language, our method focuses on ensuring that the translated output maintains a consistent narrative flow and adheres to the stylistic nuances of the original text. This adjustment is crucial for literary translation, where preserving the author's voice and the overall tone of the work is just as important as achieving accurate translation.

The fine-tuning process leverages the LoRA to refine specific parameters of the model efficiently. By applying LoRA, we can update the model with low-rank adaptations, which helps in preventing overfitting while ensuring that the model adapts ef-

fectively to the task-specific requirements. This targeted fine-tuning allows the model to better capture the long-range dependencies and stylistic elements that are essential for producing translations that are not only accurate but also faithful to the literary qualities of the source text.

### 3.2 Incremental Decoding framework

In traditional machine translation, sentences are often translated independently of one another, leading to issues with semantic coherence and stylistic consistency when viewed from a broader, document-level perspective. To address these challenges, we propose an Incremental Decoding framework that considers the translation of each sentence as part of a continuous process, taking into account the translations of previous sentences. This method ensures that the translated text maintains a cohesive flow and consistent style throughout the entire document.

The Incremental Decoding framework incorporates two key components: Context-aware information and Style-related information. Context-aware information involves using the translations of the previous n sentences as historical context when translating the current sentence. This helps maintain continuity in the narrative and ensures that the translation aligns with the broader context established in earlier sentences.

Style-related information further refines this process by incorporating translations of sentences that are similar to the current sentence in terms of content and style. These sentences are selected based on sentence and keyword similarity, ensuring that the translation reflects the stylistic nuances present in the original text. By integrating both context-aware and style-related information, the Incremental Decoding framework produces translations that are not only accurate but also harmonious in tone and structure, closely mirroring the original literary work.

### 4 Experiments

### 4.1 Datasets and Evaluation Metrics

We utilized data from the general MT shared task and the GuoFeng Webnovel Corpus. The GuoFeng Webnovel Corpus was employed in Stages 1, 2, and 3, while the general MT data was used exclusively in Stage 2. Detailed statistics of the data are presented in Table 2.

For the evaluation metrics, we utilized Sacre-BLEU (Papineni et al., 2002) to assess system performance. Given that the test set was segmented into sentence-level units, we conducted evaluations using both s-BLEU (sentence-level BLEU) and d-BLEU (document-level BLEU) scores to provide a comprehensive analysis of the translation quality.

### 4.2 Experiment Settings

In our experiments, we used Chinese-LLaMA2 (Cui et al., 2023) as the foundation model. Chinese-LLaMA2 is an enhanced and optimized version of Llama-2, specifically designed for Chinese language understanding and instruction comprehension. This model includes a larger Chinese vocabulary and benefits from incremental pretraining on a large-scale Chinese dataset, which significantly improves its semantic understanding capabilities.

For both the Continual Pre-training and Supervised Fine-Tuning stages, we adhered to the hyperparameters utilized in the Chinese-LLaMA2 project. During Stage 2, the model was trained for 1 epoch, while in Stage 3, the training was extended to 3 epochs to ensure more refined adjustments.

All experiments were conducted using 8 Nvidia GPUs, each with 64GB of memory, and employed DeepSpeed (Rasley et al., 2020) ZeRO 2 for model parallelization, which allowed for efficient handling of the large-scale model and dataset.

### 4.3 Compared Baselines

- **General Sent-Trans**: In this baseline, we directly create sentence-level translation instruction data and use it to perform Supervised Fine-Tuning on the Chinese-LLaMA2 model. This approach focuses on training the model with general sentence-level translation tasks without any specialized pre-training.

- **Literary Sent-Trans**: This baseline builds on the previous stages, as outlined in Stage 1 and Stage 2. We first subject the Chinese-LLaMA2 model to Continual Pre-training using monolingual and bilingual literary data. Following this pre-training, the model undergoes Supervised Fine-Tuning using the same sentence-level translation instruction data as in the General Sent-Trans baseline. This approach is designed to adapt the model to the literary domain before fine-tuning it with general sentence-level instructions.

976

| | Valid 1 | | Valid 2 | | Test 1 | | Test 2 | |
|---|---|---|---|---|---|---|---|---|
| | s-BLEU | d-BLEU | s-BLEU | d-BLEU | s-BLEU | d-BLEU | s-BLEU | d-BLEU |
| General Sent-Trans | 16.81 | 24.1 | 10.74 | 17.39 | 17.97 | 25.87 | 13.32 | 20.37 |
| Literary Sent-Trans | 23.35 | 30.51 | 14.64 | 21.81 | 20.91 | 28.51 | 18.02 | 25.38 |
| Literary Doc-Trans | 23.78 | 31.85 | 14.94 | 22.12 | 20.97 | 29.43 | 18.28 | 25.62 |

Table 1: **The overall results.**

| Data Source | Data Size |
|---|---|
| General MT | 25M |
| GuoFeng Webnovel Corpus | 1.9M |

Table 2: Data Statistics.

- **Literary Sent-Trans**: **This represents our final proposed approach**. After the Continual Pre-training conducted in Stage 1 and Stage 2, we further train the model using the Supervised Fine-Tuning method from Stage 3, which incorporates Context-aware and Style-related Instructions. This method aims to enhance the model's ability to maintain semantic coherence and stylistic consistency across sentences in literary document translation.

## 4.4 Results

The comparison between Literary Sent-Trans and General Sent-Trans reveals significant improvements in both s-BLEU and d-BLEU scores across various test sets, indicating that Stage 1 and Stage 2 effectively incorporated literary knowledge into the model. Furthermore, when comparing Literary Doc-Trans with Literary Sent-Trans, we observe additional gains in both s-BLEU and d-BLEU metrics, demonstrating the effectiveness of Stage 3's Context-aware and Style-related Instructions. These results collectively highlight the incremental benefits of each stage in enhancing the model's performance in literary translation. The detailed results are presented in Table 1.

## 5 Conclusion

In this work, we addressed the complex task of literary translation within the WMT24 Discourse-Level Literary Translation Task, focusing on the Chinese-English language pair. By leveraging the Chinese-Llama2 model, enhanced through Continual Pre-training and Supervised Fine-Tuning, we successfully adapted the model to capture the unique nuances of literary texts. Our Incremental Decoding framework further ensured that each sentence was translated with awareness of its broader context, resulting in more coherent and stylistically consistent translations. The improvements observed in both sentence-level and document-level BLEU scores validate the effectiveness of our approach. These results highlight the potential of combining advanced language models with specialized training strategies to tackle the intricacies of literary translation, paving the way for further research in this challenging domain.

## References

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. *CoRR*, abs/2005.14165.

Yijie Chen, Yijin Liu, Fandong Meng, Yufeng Chen, Jinan Xu, and Jie Zhou. 2023. Improving translation faithfulness of large language models via augmenting instructions.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov,

and Noah Fiedel. 2023. Palm: Scaling language modeling with pathways. *J. Mach. Learn. Res.*, 24:240:1–240:113.

Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.

Zefeng Du, Wenxiang Jiao, Longyue Wang, Chenyang Lyu, Jianhui Pang, Leyang Cui, Kaiqiang Song, Derek F Wong, Shuming Shi, and Zhaopeng Tu. 2024. On extrapolation of long-text translation with large language models. In *Findings of the Association for Computational Linguistics*.

Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. 2019. Mask-predict: Parallel decoding of conditional masked language models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 6111–6120. Association for Computational Linguistics.

Jiatao Gu, James Bradbury, Caiming Xiong, Victor O. K. Li, and Richard Socher. 2018. Non-autoregressive neural machine translation. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.

Jiaxin Guo, Minghan Wang, Daimeng Wei, Hengchao Shang, Yuxia Wang, Zongyao Li, Zhengzhe Yu, Zhanglin Wu, Yimeng Chen, Chang Su, Min Zhang, Lizhi Lei, Shimin Tao, and Hao Yang. 2021. Self-distillation mixup training for non-autoregressive neural machine translation. *CoRR*, abs/2112.11640.

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pages 639–649. Association for Computational Linguistics.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.

Wenxiang Jiao, Jen tse Huang, Wenxuan Wang, Zhiwei He, Tian Liang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Parrot: Translating during chat using large language models tuned with human translation and feedback.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pages 311–318. ACL.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*, pages 3505–3506. ACM.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data.

Zewei Sun, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Lei Li. 2020. Capturing longer context for document-level neural machine translation: A multi-resolutional approach. *Cornell University - arXiv,Cornell University - arXiv*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Minghan Wang, Jiaxin Guo, Yuxia Wang, Yimeng Chen, Chang Su, Daimeng Wei, Min Zhang, Shimin Tao, and Hao Yang. 2021. HI-CMLM: improve CMLM with hybrid decoder input. In *Proceedings of the 14th International Conference on Natural Language Generation, INLG 2021, Aberdeen, Scotland, UK, 20-24 September, 2021*, pages 167–171. Association for Computational Linguistics.

Daimeng Wei, Zhanglin Wu, Hengchao Shang, Zongyao Li, Minghan Wang, Jiaxin Guo, Xiaoyu Chen, Zhengzhe Yu, and Hao Yang. 2023. Text style transfer back-translation.

Minghao Wu, Yulin Yuan, Gholamreza Haffari, and Longyue Wang. 2024. (perhaps) beyond human translation: Harnessing multi-agent collaboration for translating ultra-long literary texts. *arXiv preprint arXiv:2405.11804*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models.

Wen Yang, Chong Li, Jiajun Zhang, and Chengqing Zong. 2023. Bigtranslate: Augmenting large language models with multilingual translation capability over 100 languages.

Zhengzhe Yu, Jiaxin Guo, Minghan Wang, Daimeng Wei, Hengchao Shang, Zongyao Li, Zhanglin Wu, Yuxia Wang, Yimeng Chen, Chang Su, Min Zhang, Lizhi Lei, Shimin Tao, and Hao Yang. 2021. Joint-training on symbiosis networks for deep nueral machine translation models. *CoRR*, abs/2112.11642.

Jiali Zeng, Fandong Meng, Yongjing Yin, and Jie Zhou. 2023. Tim: Teaching large language models to translate with comparison.

Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhengrui Ma, Yan Zhou, Langlin Huang, Mengyu Bu, Shangtong Gui, Yunji Chen, Xilin Chen, and Yang Feng. 2023. Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models.

# NovelTrans: System for WMT24 Discourse-Level Literary Translation

**Yuchen Liu[1], Yutong Yao[1], Runzhe Zhan[1], Yuchu Lin[2], Derek F. Wong[1*]**

[1]NLP2CT Lab, Department of Computer and Information Science, University of Macau
nlp2ct.{yuchen, yutong, runzhe}@gmail.com; derekfw@um.edu.mo
[2]DeepTranx, Zhuhai, China
yuchulin@deeptran.com

## Abstract

This paper describes our submission system, NovelTrans, from NLP2CT and DeepTranx for the WMT24 Discourse-Level Literary Translation Task in Chinese-English, Chinese-German, and Chinese-Russian language pairs under unconstrained conditions. For our primary system, three translations are done by GPT4o using three different settings of additional information and a terminology table generated by online models. The final result is composed of sentences that have the highest xCOMET score compared with the corresponding sentences in other results. Our system achieved an xCOMET score of 79.14 which is higher than performing a direct chapter-level translation on our dataset.

## 1 Introduction

In the rapidly evolving field of natural language processing (NLP), discourse-level literary machine translation remains a challenging task. It involves not only complex semantic phenomena but also long-term dependency, rare or new terminologies, and cultural background (Pang et al., 2024; Liu et al., 2023). These factors pose a high requirement for the translation model. Training or fine-tuning such a model is extremely costly. To address this, pretrained large language models (LLMs) and training-free methods like in-context learning (Brown et al., 2020) are widely used. Up to now, significant advancements have been made in sentence-level machine translation using training-free methods. These methods, such as TEaR (Feng et al., 2024), DUAL-REFLECT (Chen et al., 2024), Multi-Aspect Prompting and Selection (He et al., 2024), and Multi-Agent Debate (Liang et al., 2024), have proven effective. However, few studies have been conducted on the document level.

This paper presents our submission to the WMT24 Discourse-Level Literary Translation

shared task. We utilize online commercial general-purpose LLMs, DeepSeek (DeepSeek-AI et al., 2024) and GPT4o (OpenAI et al., 2024), to perform the translation with the help of techniques including *Document-level Multi-Aspect Prompting and Selection* (*d-MAPS*), *LLM-generated terminology table* and *dynamic retrieval of in-context learning examples using Reranked BM25* (R-BM25; Agrawal et al. 2023). We also explore the potential of *post-correction of punctuation errors* in LLMs' translation results. Using the above method, NovelTrans achieves an xCOMET score of 79.14, 0.68 points higher than the GPT4o baseline. Moreover, the consistency of rare or unseen terminologies has significantly improved and the number of mistranslated or awkwardly translated phrases is greatly reduced. The remaining part of this paper is structured as follows. Section 2 contains an overview of our pipelines and detailed descriptions of each procedure in the pipelines. Experiments and results analysis of our method are given in Section 3. Finally, the conclusion is presented in Section 4.

## 2 System Overview

### 2.1 Pipeline

For our pipeline, we implemented three variants which were named Primary, Contrastive-1, and Contrastive-2. The Primary system has a pipeline shown in Figure 1. For each input document, we first generate a terminology table and then replace all terminologies in the document with their corresponding translations, ensuring the consistency of terminology translation throughout the document. Then the document is split into chapters using regular expressions. Each chapter is divided into 20-line segments. Each segment is translated using GPT4o, with MAPS and R-BM25 enhancing the translation quality. The translated text will then proceed to the post-correction stage, where the GPT4o model will detect and resolve punctuation errors. For the

---
*Corresponding Author.

Figure 1: The translation flowchart of our NovelTrans system where post-correction is not included.

Contrastive-2 system, the MAPS uses a different way to determine the quality of translation and will be discussed in Section 2.2. The Contrastive-1 system is the same as the primary system except for the removal of the post-correction stage. As the API service for GPT4o we used contains a content filter, if a segment's translation is filtered by the content filter, the process will be handled using the DeepSeek API.

## 2.2 Document-level Multi-Aspect Prompting and Selection

Multi-Aspect Prompting and Selection (MAPS) is a powerful prompting strategy that can help a model understand the complicated relationships in discourse-level corpus better. Inspired by the MAPS, we chose to transfer MAPS to the document-level (d-MAPS). Considering both resource limitations and characteristics of web novels, we implemented d-MAPS as follows. We first acquire explanations for colloquialisms and the segment summary through the cooperation of DeepSeek and GPT4o. Then, three different translations are produced by GPT4o: one with explanations, one with the summary, and one without any extra information. Afterward, the COMET-22-kiwi reference-free translation quality evaluation model (Rei et al., 2022) is applied to obtain the quality score of each sentence in these three results. To select the final translation result, we employ two different strategies. In the Primary and Contrastive-

1 system, the final result is composed of sentences that have the highest xCOMET score compared to the corresponding sentences in other translations. In Contrastive-2, the final translation is determined by choosing the result with the highest average xCOMET score.

## 2.3 LLM-generated Terminology Table

In the traditional novel translation pipeline, it is crucial to set up a terminology table before the translation to unify the translations of those rare terms throughout the corpus. To generate the terminology table, we use the DeepSeek API which has better knowledge of Chinese cultural backgrounds to retrieve proper nouns and then translate these words into the target language considering their context. With the terminology table acquired, we then replace all the terms in the source corpus with their corresponding translations to ensure consistency. The consistency mentioned above refers to the uniformity of special terminology translation.

## 2.4 Re-ranked BM25

Re-ranked BM25 (R-BM25; Agrawal et al. 2023) is an in-context example retriever that can ensure both sample quality and retrieving speed. After 100 sentences are retrieved by a normal BM25 retriever, a score will be computed for each sentence using the following formula, in which S and Q denote the source and retrieved sentence's n-grams separately.

$$R_n = \frac{\sum_{\text{ngram} \in S \cap Q} \text{Count}_{\text{matched}}(\text{ngram})}{\sum_{\text{ngram} \in S} \text{Count}_S(\text{ngram})} \quad (1)$$

$$\text{Score} = \exp\left(\frac{1}{n}\sum_n \log(R_n)\right) \quad (2)$$

Then these sentences are re-ranked using these scores to solve the problems that BM25 favors rare words (Robertson and Zaragoza, 2009). To form the sentence pool for the R-BM25 to search, we utilize the GuoFeng Webnovel Corpus[1] (Wang et al., 2023) which has three subsets named TRAIN, VALID1, and VALID2. By combining all three subsets, we formed a large dataset and then filtered out sentences with low xCOMET scores. During the experiment, VALID2 is not included because our valid set is sampled from VALID2. To generate the in-context learning examples for a particular segment, we retrieve three samples for each sentence

---

[1] http://www2.statmt.org/wmt23/literary-translation-task.html

| | Zh-En | | Zh-Ru | | Zh-De | |
|---|---|---|---|---|---|---|
| | xCOMET | d-BLEU | xCOMET | d-BLEU | xCOMET | d-BLEU |
| DeepSeek | 76.58 | 18.03 | - | - | - | - |
| GPT3.5-Turbo-16k | 77.33 | 17.92 | - | - | - | - |
| GPT4o baseline | 78.46 | **18.85** | 83.74 | **26.51** | 80.69 | 38.33 |
| NovelTrans (Ours) | **79.14** | 18.69 | **84.42** | 26.44 | **80.85** | 39.78 |

Table 1: Experiment result compared with other models. Results listed here expect NovelTrans are all generated by direct chapter-level translation. xCOMET scores in this and tables below are all computed using XCOMET-XL.

| Method | xCOMET | BLEU | d-BLEU |
|---|---|---|---|
| GPT4o baseline | 78.46 | 20.17 | 18.85 |
| NovelTrans | **79.14** | 19.94 | 18.69 |
| *w/o ICL* | 78.85 | 19.67 | 18.63 |
| *w/o ICL & Terminology Table* | 78.71 | 20.60 | 18.63 |
| *w/o ICL, Terminology Table & d-MAPS* | 78.68 | **20.80** | **18.97** |

Table 2: Ablation study of our proposed pipeline. ICL examples are selected by R-BM25 score. Terminology table represents the terminology table obtained by the cooperation of GPT4o and DeepSeek. The GPT4o baseline is generated by directly translating the text at the chapter level.

in that segment using R-BM25 and then randomly sample eight sentences to form the final in-context learning example. It is tested that choosing eight examples will result in the best performance boost.

## 2.5 Post-Correction of Translation

After reviewing the translation results, we observed that punctuation errors, such as comma splices, appeared at a high frequency due to the inappropriate use of punctuation in the source corpus. To solve this, we employed a post-processing method that uses GPT4o to correct punctuation errors at the sentence level. Given the sentence above and below the target sentence, we asked the model to check and resolve punctuation errors. This method resulted in a better version of the target sentence.

## 3 Experiments

### 3.1 Experiments Setup

The datasets we used are GuoFeng Webnovel Corpus V1 and V2. V1 contains a Chinese-English parallel corpus while V2 contains Chinese-German and Chinese-Russian nonparallel corpus. For the Chinese-English direction, we performed experiments on 10 chapters in VALID2 of the dataset. These chapters are taken from different books to avoid bias. For Chinese-German and Chinese-Russian direction, we chose 4 chapters from different books and aligned them separately using

GPT4o API before experimenting. The GPT4o API we used is provided by OpenAI. The DeepSeek API is provided by DeepSeek Open Platform[2]. Since the BLEU score faces the problem of inaccuracy in evaluating Zero Pronoun Translation tasks (Zhan et al., 2023; Xu et al., 2023), we focused more on the COMET score. To be better aligned with the human evaluation, we chose to use **XCOMET-XL** (Guerreiro et al., 2023) to compute the xCOMET score. BLEU and d-BLEU scores are all computed by **SacreBleu** (Post, 2018). To compute d-BLEU, we join all sentences in the document together and treat them as a single sentence since it is the method used to compute the d-BLEU score in the previous year's WMT literary translation task (Wang et al., 2023).

### 3.2 Results

Table 1 shows the comparison between our system and other online models in Chinese-English, Chinese-German, and Chinese-Russian translation direction. The result shows that our system achieves a higher xCOMET score in exchange for the d-BLUE performance.

### 3.3 Ablation Study

We conduct ablation study on Chinese-English direction. The result, provided in Table 2, shows that

---

[2]https://platform.deepseek.com/

| Source | GPT4o Baseline | NovelTrans |
|---|---|---|
| 走，全部跟我走，去破坏对方的 (rival) 世界级传送阵. | Go, all of you come with me to destroy the *other side's* (Wrong) world-class teleportation array. | Let's go, everyone follows me to destroy the *enemy's* (Correct) world-class teleportation array. |
| 这四个字，是郑州城人类最后的绝唱 (the last song of mankind in the city of Zhengzhou). | These four words were the *last human song of Zhengzhou* (Bad Phrase Translation). | These four words were the *last elegy of humanity in Zhengzhou city* (Correct). |

Table 3: Case study where examples are taken from different pipeline methods.

| Source | Without Correction | With Correction |
|---|---|---|
| "别紧张，自己人。" | "Don't be nervous, I'm one of you." | "Don't be nervous; I'm one of you." |
| 他们打开背后的涡旋引擎跳了下去 | They activated the vortex engine on their backs, jumping down | They activated the vortex engine on their backs before jumping down. |

Table 4: Comparison of translation results with or without post-translation correction.

| Position | Source | Without Term Table | With Term Table |
|---|---|---|---|
| Near the start of a chapter | 若非此刻在天渡船上,可能已经大打出手. | If they weren't on the *Tian Du ship*, he might have already started a fight. | If they weren't on the *Heavenly Ferry*, he might have already started a fight. |
| Near the end of the same chapter | 不多时,天渡船抵达对岸. | Before long, the *Heaven Crossing Boat* (Inconsistent) reached the other side. | Before long, the *Heavenly Ferry* (Consistent) reached the opposite bank. |

Table 5: Comparison of translation results with or without LLM-generated terminology table.

removal of component in our system will result in a performance drop on xCOMET.

### 3.4 Analysis

Table 3 shows two examples taken from our experiment. In the first example, the direct translation of GPT4o uses an ambiguous phrase, "other side", which can mean both an enemy and a geographically opposite side. However, with the context, we can easily determine that the "other side" here conveys only the meaning of "rival". In the second example, the Chinese word "绝唱" which means the best art piece an artist has ever made is misused as "last song before their death" in the source sentence. Our system understood what the author wanted to convey and chose a suitable word, "elegy", rather than doing a literal translation. These examples show that, compared with the baseline, our method has a stronger understanding of the context and Chinese cultural background. Table 4 demonstrates the effect of post-correction. The GPT4o model can detect and correct punctuation errors, especially comma splices that occur at high frequency, in various ways. Table 5 shows an example of inconsistency in the translation of special terms and our method can greatly reduce this type of problem.

### 4 Conclusion

We successfully deployed a discourse-level translation pipeline using online language models and adapted several sentence-level techniques for discourse-level translation. Our system achieved a higher xCOMET score than direct translation using GPT-4o. However, our research has some limitations. Adapting MAPS to discourse-level translation may disrupt long-term dependencies, indicating a need for further investigation in this

area. Additionally, our method utilizes significantly more tokens than direct translation, necessitating further discussion on how to reduce token usage.

## Acknowledgement

## References

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. In-context examples selection for machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8857–8873, Toronto, Canada. Association for Computational Linguistics.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

Andong Chen, Lianzhang Lou, Kehai Chen, Xuefeng Bai, Yang Xiang, Muyun Yang, Tiejun Zhao, and Min Zhang. 2024. DUAL-REFLECT: Enhancing large language models for reflective translation through dual learning feedback mechanisms. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 693–704, Bangkok, Thailand. Association for Computational Linguistics.

DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Liu, Xin Xie, Xingkai Yu, Xinnan Song, Xinyi Zhou, Xinyu Yang, Xuan Lu, Xuecheng Su, Y. Wu, Y. K. Li, Y. X. Wei, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Zheng, Yichao Zhang, Yiliang Xiong, Yilong Zhao, Ying He, Ying Tang, Yishi Piao, Yixin Dong, Yixuan Tan, Yiyuan Liu, Yongji Wang, Yongqiang Guo, Yuchen Zhu, Yuduan Wang, Yuheng Zou, Yukun Zha, Yunxian Ma, Yuting Yan, Yuxiang You, Yuxuan Liu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhewen Hao, Zhihong Shao, Zhiniu Wen, Zhipeng Xu, Zhongyu Zhang, Zhuoshu Li, Zihan Wang, Zihui Gu, Zilin Li, and Ziwei Xie. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model.

Zhaopeng Feng, Yan Zhang, Hao Li, Bei Wu, Jiayu Liao, Wenqiang Liu, Jun Lang, Yang Feng, Jian Wu, and Zuozhu Liu. 2024. Tear: Improving llm-based machine translation with systematic self-refinement.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection.

Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2024. Exploring human-like translation strategy with large language models. *Transactions of the Association for Computational Linguistics*, 12:229–246.

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Rui Wang, Yujiu Yang, Zhaopeng Tu, and Shuming Shi. 2024. Encouraging divergent thinking in large language models through multi-agent debate.

Xuebo Liu, Yutong Wang, Derek F. Wong, Runzhe Zhan, Liangxuan Yu, and Min Zhang. 2023. Revisiting commonsense reasoning in machine translation: Training, evaluation and challenge. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15536–15550, Toronto, Canada. Association for Computational Linguistics.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report.

Jianhui Pang, Fanghua Ye, Longyue Wang, Dian Yu, Derek F Wong, Shuming Shi, and Zhaopeng Tu. 2024. Salute the classic: Revisiting challenges of machine translation in the age of large language models. *arXiv preprint arXiv:2401.08350*.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389.

Longyue Wang, Zhaopeng Tu, Yan Gu, Siyou Liu, Dian Yu, Qingsong Ma, Chenyang Lyu, Liting Zhou, Chao-Hong Liu, Yufeng Ma, Weiyu Chen, Yvette Graham, Bonnie Webber, Philipp Koehn, Andy Way, Yulin Yuan, and Shuming Shi. 2023. Findings of the WMT 2023 shared task on discourse-level literary translation: A fresh orb in the cosmos of LLMs. In *Proceedings of the Eighth Conference on Machine Translation*, pages 55–67, Singapore. Association for Computational Linguistics.

Mingzhou Xu, Longyue Wang, Siyou Liu, Derek F Wong, Shuming Shi, and Zhaopeng Tu. 2023. A benchmark dataset and evaluation methodology for

chinese zero pronoun translation. *Language Resources and Evaluation*, 57(3):1263–1293.

Runzhe Zhan, Xuebo Liu, Derek F. Wong, Cuilian Zhang, Lidia S. Chao, and Min Zhang. 2023. Test-time adaptation for machine translation evaluation by uncertainty minimization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 807–820, Toronto, Canada. Association for Computational Linguistics.

# LinChance×NTU for Unconstrained WMT2024 Literary Translation

**Kechen Li[1], Yaotian Tao[1], Hongyi Huang[1], Tianbo Ji[2]**

[1]Jiangsu Linchance Technology Co., Ltd. (LinChance)
[2]School of Transportation and Civil Engineering, Nantong University

{likechen,taoyaotian,huanghongyi}@linchance.com, jitianbo@ntu.edu.cn

## Abstract

The rapid growth of deep learning has spurred significant advancements across industries, particularly in machine translation through large language models (LLMs). However, translating literary still presents challenges, including cross-cultural nuances, complex language structures, metaphorical expressions, and cultural differences. To address these issues, this study utilizes the Llama and Phi models using both LoRA and full-parameter techniques, alongside a prompt-based translation system. Full-parameter tuning of the Llama-3-Chinese-8B-Instruct model was unsuccessful due to memory constraints. In terms of the WMT task, the fully fine-tuned Phi 3 model was selected for submission due to its more natural and fluent translations. Nonetheless, results showed that LoRA and the prompt-based system significantly improved the Llama3 model's performance, surpassing other models in BLEU and ROUGE evaluations.

## 1 Introduction

In recent years, the development of deep learning has spread across various industries (Ji et al., 2024), and the impact of large language models (LLMs) on these industries has been particularly significant (Lyu et al., 2023). Despite the fact that many challenges in machine translation (MT) have been overcome (Wang et al., 2023), literary translation still encounters cross-cultural issues, including such as processing complex languages, understanding metaphorical expressions, and addressing cultural differences (Lyu et al., 2020). Meanwhile, choosing the right model has become a key topic in terms of neural network-based (NN-based) MT (Xia, 2020), as models based on the source language typically have advantages in handling tasks regarding that language. Two main models are involved in this research: Llama3-Chinese-8B-

Instruct[1] and Phi-3-mini-128k-instruct-Chinese[2]. The former fine-tuned the Llama3 model (Dubey et al., 2024) with 5 million instruction data points from the community, which significantly enhances its performance in Chinese-language tasks with a better ability of understanding Chinese contexts. The latter is with less than half the size (3.8B parameters) of the Llama3 8B version, which can surpass the performance of Llama3 with less computational resources.

LoRA (Sundaram et al., 2019) is a lightweight fine-tuning technique mainly used for efficiently training large models. Compared to traditional full-parameter fine-tuning, LoRA decomposes the trained parameter matrices into low-rank forms, resulting in a reducing number of parameters and less computational cost of training. It is especially appropriate for fine-tuning LLMs with constrained resources while maintaining high performance. In this research, we utilize Llama-Factory[3] (Zheng et al., 2024), an optimized framework designed specifically for fine-tuning LLMs like Llama. It supports various advanced training techniques, including mixed precision training and gradient accumulation, to improve training efficiency and reduce computational resource requirements. By integrating lightweight methods like LoRA, it can achieve efficient and stable model fine-tuning in resource-constrained environments, helping to quick adjustment and deployment of LLMs.

In general, our contributions can be summarized as follows:

- We conduct a comprehensive experiment of two major large language models, Llama-3-Chinese-8B-Instruct and Phi-3-mini-128k-instruct-Chinese, for the task of WMT2024

---

[1]https://huggingface.co/hfl/Llama-3-chinese-8b-instruct
[2]https://huggingface.co/shareAI/Phi-3-mini-128k-instruct-Chinese
[3]https://github.com/hiyouga/Llama-Factory

Literary Translation. Our results demonstrate that both models perform excellently in handling Chinese tasks, especially when facing cross-cultural challenges in literary translation.

- We applied the LoRA technique to efficiently fine-tune the Llama model, significantly reducing computational costs while maintaining high translation quality. We additionally optimized the fine-tuning process by leveraging the Llama-Factory framework. Our experimental results demonstrate that the combination of LoRA and Llama-Factory can effectively support the adaptation and deployment of large-scale models in resource-constrained environments.

- We investigate the strengths and weaknesses of each model, particularly in terms of fluency and diversity (captured by ROUGE) as well as accuracy (captured by BLEU).

## 2 Related Work

**Large Language Models (LLMs) for Machine Translation** The application of Large Language Models (LLMs) in machine translation (MT) has seen significant advancements, particularly in general domain translation (Wang et al., 2023). Pre-trained models such as Llama and Phi3 have been increasingly employed for tasks requiring semantic understanding across languages. Studies have highlighted how instruction-tuned LLMs can improve translation quality by adapting to the syntactic structures and cultural nuances of target languages. This is particularly relevant for our work as we evaluate the Llama-3 Chinese model (Cui et al., 2023), which is fine-tuned for literary translation.

**Challenges in Literary Translation** While MT systems have progressed in many domains (Du et al., 2024), the translation of literary texts remains particularly challenging due to the need to capture nuanced expressions, idioms, and stylistic elements. Literary translation (Jones, 2019) is often considered the "last frontier" for MT. Prior work has explored how traditional sentence-based MT systems struggle with long, complex passages found in literary texts (Aliguliyev, 2009). This is in line with the focus of our experiments, which attempt to handle these unique challenges using Llama-3 Chinese models and fine-tuning techniques.

**Fine-Tuning Techniques for MT** In order to ad-dress the computational limitations and language understanding challenges associated with large models, various fine-tuning approaches have been proposed (Nicholas and Bhatia, 2023). Recent studies on Low-Rank Adaptation (LoRA) have demonstrated that memory-efficient fine-tuning methods allow for high-quality performance on GPUs with limited memory. LoRA's success in reducing memory consumption while maintaining model accuracy has been a key technique in our experiments, particularly with the Llama-3 Chinese model.

**Evaluation of Translation Quality** The evaluation of literary translations poses its own challenges, as traditional metrics like BLEU may not fully capture the nuances of a good translation. (Pang et al., 2024) Newer approaches, such as Monolingual Human Preference (MHP) and Bilingual LLM Preference (BLP), have been proposed to better assess translation quality in a literary context. (Wu et al., 2024) Our experiments draw on these evaluation techniques, comparing model outputs through both automated metrics and human preference assessments to gauge the effectiveness of different fine-tuning strategies.

## 3 Experiment

### 3.1 Evaluation Metrics

To achieve accurate evaluation of MT result (Chang et al., 2024), two prevailing evaluation metrics are utilized: BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004). BLEU is an automated evaluation metric based on n-gram matching, primarily used to measure the similarity between machine translation outputs and reference translations. By calculating the overlap of n-grams of different lengths between the model's output and the reference translation, BLEU can reflect the accuracy of the translation to a certain extent. ROUGE is a widely used automatic text evaluation metric mainly used to compare the similarity between generated text and reference text. ROUGE-L, in particular, is based on the Longest Common Subsequence (LCS) (Bergroth et al., 2000) and measures the similarity between generated text and reference text in terms of length matching and word order. Compared to BLEU, ROUGE captures both text diversity and fluency.

### 3.2 Prompt Engineering

The prompt design focuses on the task of translating Guofeng (traditional Chinese-style) novels,

aiming to ensure that the translated text faithfully conveys the literary and cultural nuances of the original, while maintaining translation efficiency and accuracy. By establishing clear guidelines, the prompt emphasizes fidelity to the original text, concise output, and quality control to ensure that the translation remains fluent while preserving the original style and tone. The prompt is designed following the CoT (Chain of Thoughts) framework (Wei et al., 2022), with the specific approach outlined below:

First, the prompt introduces automatic language detection and translation features, enabling efficient Chinese-to-English translation of Guofeng novels to enhance processing speed and coherence. Secondly, fidelity to the original is critical, requiring the preservation of the original tone, style, and expression. Special attention is given to details such as pronouns, with a focus on word-for-word translation to avoid distorting the literary essence due to cultural or linguistic differences. The prompt further emphasizes objectivity in translation, avoiding any omissions or commentary, ensuring the completeness and authenticity of the translated text.

Additionally, the translated text must be concise, with no added annotations, ensuring that the classical charm and cultural context of the novel are naturally conveyed, enhancing the reader's experience. To ensure quality, the prompt requires thorough review and correction of the translated output, avoiding mistranslations or omissions, and ensuring that the text aligns with the target language's fluency and conventions. The prompt is task-oriented, providing only the final, revised translation, avoiding irrelevant information or excessive explanation, which improves processing efficiency and suits large-scale Guofeng novel translation projects.

### 3.3 Experiment 1: Training with Llama-3-Chinese-8B-Instruct

An initial attempt was made to fully train the Llama-3-Chinese-8B-Instruct model on the dataset. However, the process failed due to insufficient memory. The model's large size and the memory requirements exceeded the capabilities of the available hardware, necessitating use a smaller model or shift to a more memory-efficient fine-tuning method.

### 3.4 Experiment 2: Fine-Tuning Llama3 8B Model with LoRA

Given the memory constraints, the Llama3 8B model was fine-tuned using the LoRA technique

on a dataset of 10,000 samples. The fine-tuning was performed with several key hyperparameters to optimize model performance and manage computational resources effectively. The learning rate was set to 1e-5 (Jin et al., 2023), using a cosine learning rate scheduler to gradually reduce the learning rate and improve training efficiency (Kim et al., 2021). A per-device train batch size of 2 was chosen to balance between memory usage and model update frequency, with gradient accumulation over 16 steps to simulate a larger batch size without requiring additional GPU memory. The training process was conducted over 10 epochs to ensure sufficient learning from the data, utilizing a maximum of 10,000 samples. Additionally, the model was trained with mixed precision (fp16) (Le Gallo et al., 2018) to reduce memory usage and accelerate computation. The evaluation strategy was set to evaluate the model performance at regular steps to monitor its progress closely. The results of this fine-tuning experiment are summarized in Table 1.

Table 1: Results of Fine-Tuning Llama3 8B Model with LoRA

| Metric | Value | Description |
|---|---|---|
| BLEU-4 | 55.28 | A metric which evaluates the quality of a candidate by computing the n-gram ($n = 4$) precision with references. |
| ROUGE-1 | 60.18 | A variation of ROUGE where 1 means unigrams. |
| ROUGE-2 | 37.40 | A variation of ROUGE where 2 means bigrams. |
| ROUGE-L | 55.91 | A variation of ROUGE where L means longest common subsequences (LCS). |
| Runtime | 3m8s | Total runtime |
| Sample/s | 1.594 | Samples processed per second |
| Step/s | 1.594 | Training steps per second |

**Analysis:** The results from this experiment were quite promising, with high BLEU and ROUGE scores. The LoRA technique allowed the model to be fine-tuned without running into memory issues, demonstrating that it is an effective method for working with large models on limited hardware. The high ROUGE scores suggest that the model

was able to generate translations that were both accurate and fluent.

### 3.5 Experiment 3: Full Fine-Tuning with Phi Chinese Model

The Phi Chinese model was fully fine-tuned on a smaller dataset of 2,000 samples.

Table 2: Results of Full Fine-Tuning with Phi Chinese Model

| Metric | Value |
|--------|-------|
| BLEU-4 | 50.93 |
| ROUGE-1 | 51.90 |
| ROUGE-2 | 27.19 |
| ROUGE-L | 46.41 |
| Runtime | 16m34s |
| Sample/s | 0.503 |
| Step/s | 0.503 |

**Analysis:** The performance of the Phi Chinese model, while adequate, was noticeably lower than that of the Llama3 8B model fine-tuned with LoRA. The lower BLEU and ROUGE scores could be attributed to the smaller model size and the limited dataset, which may not have provided enough data for the model to generalize well.

### 3.6 Experiment 4: Full Fine-Tuning with Phi3 Chinese 3.5B Model

The Phi3 Chinese 3.5B model was fully fine-tuned on a dataset of 1,500,000 samples. The fine-tuning process was carefully configured with a set of key hyperparameters to optimize the model's performance while efficiently managing computational resources. We set the learning rate to 1e-5, which is low enough to ensure stable training and prevent the model from overshooting optimal weights but sufficient to allow for meaningful updates to the model parameters. A batch size of 128 was chosen to strike a balance between training speed and memory constraints. To further accommodate large batch sizes, a gradient accumulation step of 16 was used, effectively increasing the batch size without exceeding GPU memory limits. The model was trained using mixed-precision floating-point, allowing for faster computation and reduced memory usage, which is crucial when dealing with large-scale models. We set the number of epochs to 3.0 to provide sufficient training cycles while minimizing the risk of overfitting. A temperature parameter of



Figure 1: Loss Over Time for the Phi-3-mini-128k-instruct-Chinese (3.8B) model.

0.4 was employed during the generation phase to control the randomness and diversity of the model's output, balancing between creativity and coherence. The results of this fine-tuning experiment are summarized in Table 3.

Table 3: Results of Full Fine-Tuning with Phi3 Chinese 3.5B Model

| Metric | Value |
|--------|-------|
| BLEU-4 | 49.14 |
| ROUGE-1 | 49.80 |
| ROUGE-2 | 25.09 |
| ROUGE-L | 45.24 |
| Runtime | 2m36s |
| Sample/s | 1.278 |
| Step/s | 1.278 |

**Training and Evaluation Loss:** To evaluate the fine-tuning process of the Phi-3-mini-128k-instruct-Chinese model, we tracked the training and evaluation loss over time, as illustrated in Figure 1. The model, which consists of 3.8 billion parameters, was fine-tuned on a diverse dataset using LoRA and full-parameter tuning techniques. Both the training and evaluation losses were monitored to assess model convergence and stability during the fine-tuning process.

**Loss Over Time:** As shown in Figure 1, the initial training loss starts relatively high, around 1.9, and decreases sharply during the early stages of training. By the end of the first epoch, the loss drops to approximately 1.3, indicating that the model quickly learns to generalize to the underlying patterns in the training data. The evaluation loss follows a similar trend, closely mirroring the training loss, which suggests that the model gener-

alizes well without overfitting during the training process. By the second epoch, the loss stabilizes around 1.2 for both training and evaluation, demonstrating the model's ability to maintain consistent performance throughout the training process. The convergence of the loss indicates that the model is reaching its optimal capacity under the current fine-tuning setup.

**Observations and Insights:** The relatively close alignment of training and evaluation losses suggests that the fine-tuning process successfully mitigated the risk of overfitting, which is often a concern when dealing with large models and smaller, task-specific datasets. Moreover, the overall reduction in loss suggests that the Phi-3-mini-128k-instruct-Chinese model was able to effectively capture the nuances of the Chinese language and the intricate nature of literary translation tasks, as intended in this study.

## 4   Conclusion and Future Work

In this paper, we conduct experiments to provide valuable insights into the performance of different models and fine-tuning techniques for the WMT2024 Literary Translation Task. The Llama3 8B model, when fine-tuned using the LoRA technique, demonstrated the best performance, highlighting the importance of memory-efficient training methods in dealing with large models on limited hardware. The results from the Phi and Phi3 models suggest that model size alone may not guarantee the better performance, and the choices of fine-tuning method and dataset size are critical factors in achieving high-quality translations. In the future, we plan to investigate the performance of even larger models (e.g., Llama3 70B) to explore the trade-offs between model size, computational resources, and translation quality. In addition, since the metrics we used may correlate negatively with human judgements (Ji et al., 2022), developing task-specific evaluation metrics would be valuable for the accurate assessment of model performance.

## References

Ramiz M Aliguliyev. 2009. A new sentence similarity measure and sentence based extractive technique for automatic text summarization. *Expert Systems with Applications*, 36(4):7764–7772.

Lasse Bergroth, Harri Hakonen, and Timo Raita. 2000. A survey of longest common subsequence algorithms. In *Proceedings Seventh International Symposium on String Processing and Information Retrieval. SPIRE 2000*, pages 39–48. IEEE.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45.

Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.

Zefeng Du, Wenxiang Jiao, Longyue Wang, Chenyang Lyu, Jianhui Pang, Leyang Cui, Kaiqiang Song, Derek F Wong, Shuming Shi, and Zhaopeng Tu. 2024. On extrapolation of long-text translation with large language models. In *Findings of the Association for Computational Linguistics*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Tianbo Ji, Yvette Graham, Gareth Jones, Chenyang Lyu, and Qun Liu. 2022. Achieving reliable human assessment of open-domain dialogue systems. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6416–6437, Dublin, Ireland. Association for Computational Linguistics.

Tianbo Ji, Kechen Li, Quanwei Sun, and Zexia Duan. 2024. Urban transport emission prediction analysis through machine learning and deep learning techniques. *Transportation Research Part D: Transport and Environment*, 135:104389.

Hongpeng Jin, Wenqi Wei, Xuyu Wang, Wenbin Zhang, and Yanzhao Wu. 2023. Rethinking learning rate tuning in the era of large language models. In *2023 IEEE 5th International Conference on Cognitive Machine Intelligence (CogMI)*, pages 112–121. IEEE.

Francis R Jones. 2019. Literary translation. In *Routledge encyclopedia of translation studies*, pages 294–299. Routledge.

Chiheon Kim, Saehoon Kim, Jongmin Kim, Donghoon Lee, and Sungwoong Kim. 2021. Automated learning rate scheduler for large-batch training. *arXiv preprint arXiv:2107.05855*.

Manuel Le Gallo, Abu Sebastian, Roland Mathis, Matteo Manica, Heiner Giefers, Tomas Tuma, Costas Bekas, Alessandro Curioni, and Evangelos Eleftheriou. 2018. Mixed-precision in-memory computing. *Nature Electronics*, 1(4):246–253.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Chenyang Lyu, Tianbo Ji, and Yvette Graham. 2020. Incorporating context and knowledge for better sentiment analysis of narrative text. In *Text2Story@ ECIR*, pages 39–45.

Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*.

Gabriel Nicholas and Aliya Bhatia. 2023. Lost in translation: large language models in non-english content analysis. *arXiv preprint arXiv:2306.07377*.

Jianhui Pang, Fanghua Ye, Derek F Wong, and Longyue Wang. 2024. Anchor-based large language models. In *Findings of the Association for Computational Linguistics*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Jothi Prasanna Shanmuga Sundaram, Wan Du, and Zhiwei Zhao. 2019. A survey on lora networking: Research problems, current solutions, and open issues. *IEEE Communications Surveys & Tutorials*, 22(1):371–388.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. *arXiv preprint arXiv:2304.02210*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Minghao Wu, Yulin Yuan, Gholamreza Haffari, and Longyue Wang. 2024. (perhaps) beyond human translation: Harnessing multi-agent collaboration for translating ultra-long literary texts. *arXiv preprint arXiv:2405.11804*.

Ying Xia. 2020. Research on statistical machine translation model based on deep neural network. *Computing*, 102(3):643–661.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

# Improving Context Usage for Translating Bilingual Customer Support Chat with Large Language Models

**José Pombal**[1,2,3*], **Sweta Agrawal**[2*], **André F.T. Martins**[1,2,3,4]

[1]Unbabel [2]Instituto de Telecomunicações
[3]Instituto Superior Técnico, Universidade de Lisboa [4]ELLIS Unit Lisbon

jose.pombal@unbabel.com,swetaagrawal20@gmail.com

## Abstract

This paper describes Unbabel+IT's submission to the Chat Shared Task held at the Workshop of Machine Translation 2024. The task focuses on translating customer support chats between agents and customers communicating in different languages. We present two strategies for adapting state-of-the-art language models to better utilize contextual information when translating such conversations. Our training strategy involves finetuning the model on chat datasets with context-augmented instructions, resulting in a specialized model, TOWERCHAT. For inference, we propose a novel quality-aware decoding approach that leverages a context-aware metric, CONTEXTCOMET, to select the optimal translation from a pool of candidates. We evaluate our proposed approach on the official shared task datasets for ten language pairs, showing that our submission consistently outperforms baselines on all and competing systems on 8 out of 10 language pairs across multiple automated metrics. Remarkably, TOWERCHAT outperforms our contrastive submission based on the much larger TOWER-v2-70B model while being 10× smaller. According to human evaluation, our system outperforms all other systems and baselines across all language pairs. These results underscore the importance of context-aware training and inference in handling complex bilingual dialogues.

## 1 Introduction

The focus of this year's chat translation (Chat MT) shared task is the translation of conversations in customer service applications. This task differs from classical MT in that the interactions are bilingual and the texts are often more dynamic, contextualized, and informal than the structured content typically found in news or Wikipedia articles. In such scenarios, leveraging conversation context could potentially help avoid cases of lexical inconsistency and incoherence (Läubli et al., 2018; Toral

---

*Equal contribution.



Figure 1: A WMT24 sample conversation (some turns omitted) with reference English translations. Without context, TOWERCHAT mistranslates "montagem de elencos" to "casting". With context, it correctly translates the source, understanding the customer is talking about a squad building challenge ("dme").

et al., 2018). However, previous editions of the Chat MT shared task have shown that standard MT models are still incapable of doing so (Farajian et al., 2020; Farinha et al., 2022).

Large Language Models (LLMs) present a promising avenue to address this issue. Not only are they becoming the state-of-the-art solution for multilingual machine translation (Zhang et al., 2023; Wei et al., 2023; Alves et al., 2023; Reinauer et al., 2023; Zhu et al., 2024), but they are also known to handle context adeptly (Karpinska and

Iyyer, 2023; Wang et al., 2023; He et al., 2024). Despite their potential, the application of LLMs in understanding and translating bilingual chat conversations remains underexplored. We aim to bridge this gap by investigating how translation LLMs can be adapted for the Chat MT task and how they can effectively leverage conversational context to produce better translations.

Using TOWER LLM (Alves et al., 2024), a strong LLM specialized for MT and related tasks, we show that an LLM not fine-tuned for the chat domain struggles to leverage context for disambiguation, often resulting in translations that are worse than those produced without context. We thus propose two solutions to improve context usage for translation LLMs. First, we build a translation model tailored for Chat MT – TOWERCHAT – finetuned on a carefully constructed context-augmented dataset. Second, to further improve the usage of contextual information during inference, we take a novel approach of performing quality-aware decoding (Fernandes et al., 2022, QAD) with a context-aware MT evaluation metric, CONTEXTCOMET (Agrawal et al., 2024). QAD approaches select one best hypothesis from a pool of candidates using an MT metric, and have been shown to improve translation quality (Freitag et al., 2022; Fernandes et al., 2022; Farinhas et al., 2023).

This serves as our primary submission to the WMT24 Chat MT shared task, along with two contrastive ones – TOWERCHAT without QAD, and TOWER-V2-70B. The TOWER-V2-70B model is the strongest version of TOWER, which was developed for the General MT shared task.[1] The translations obtained from our approach consistently achieve the best scores across all language pairs tested, as measured by both automatic MT metrics (neural and lexical) and lexical cohesion metrics (MUDA accuracy) and human evaluation, beating strong baselines that disregard the context of conversations. Furthermore, TOWERCHAT without QAD maintains general translation capabilities and achieves better or comparable quality to TOWER-V2-70B, outlining the importance of in-domain adaptation of translation LLMs on Chat MT data.

## 2 Chat Translation Shared Task: Dataset and Challenges

This year's chat MT dataset includes bilingual online customer service chats between an English-

speaking agent and clients who speak Portuguese, French, Italian, Dutch, or Korean. These conversations are often unplanned, informal, and nonstandard, contrasting with the well-formed text of most other translation domains. An example conversation is shown in Figure 1.

We present the general statistics from this year's shared task datasets in Table 1, including (i) the number of instances in the dataset for each language pair; (ii) the average character length of the source segments; (iii) the average number of segments in a conversation and (iv) the percentage of segments tagged with MUDA (Fernandes et al., 2023), an automatic tagger for identifying tokens belonging to certain discourse classes (lexical cohesion, verb forms, pronouns, formality) of potentially ambiguous translations. While the development and test sets exhibit a similar distribution in terms of segment length and count, they differ significantly from the training dataset. Furthermore, up to 30% en↔fr instances are tagged as requiring disambiguation according to MUDA, highlighting the complexity and the need for contextual information to generate high-quality translations.[2]

Next, we describe the process of building TOWERCHAT, which was conditioned by the aforementioned inherent complexities of Chat MT.

## 3 Adapting TOWER for Chat Translation

LLMs have shown the potential to use contextual information to perform many NLP tasks (Karpinska and Iyyer, 2023). In this work, we investigate whether providing contextual information can improve translation quality for bilingual chats using strong translation LLMs like TOWERINSTRUCT. Contrary to our expectations, our preliminary results indicate that incorporating context into the prompt instruction diminishes overall translation quality. We believe this is due to TOWERINSTRUCT's training data lacking chat-specific MT examples, which results in the model's unfamiliarity with the context format and the inability to adequately use context (Section 5). To mitigate this and improve the usage of contextual information, we propose two strategies – one for training and one for inference.

---

| Language Pair | # Instances | | | Avg. Source Length | | | Avg. # Segments per Conversation | | | % MuDA tagged | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Train | Dev | Test | Train | Dev | Test | Train | Dev | Test | Dev | Test |
| en↔de | 17805 | 2569 | 2041 | 47.40 | 52.26 | 53.09 | 36.12 | 31.33 | 30.46 | 15.65 | 15.78 |
| en↔fr | 15027 | 3007 | 2091 | 41.84 | 54.90 | 56.23 | 56.92 | 33.41 | 32.17 | 29.43 | 29.65 |
| en↔pt-br | 15092 | 2550 | 2040 | 42.72 | 46.46 | 46.49 | 34.69 | 26.56 | 27.95 | 13.02 | 12.99 |
| en↔ko | 16122 | 1935 | 1982 | 39.86 | 47.67 | 46.90 | 38.11 | 50.92 | 47.19 | 0.41 | 0.50 |
| en↔nl | 15463 | 2549 | 2015 | 45.40 | 52.31 | 54.31 | 25.99 | 35.40 | 34.74 | 22.01 | 23.13 |

Table 1: Statistics for each language pair and split of the data for the WMT24 Chat MT shared task.

---

Context: {context}
Translate the following **{source_lang}** source text to **{target_lang}**, given the context.
**{source_lang}**: {source_seg}
**{target_lang}**: {target_seg}

Figure 2: Instructions with context for Chat MT. Parts in purple are only included when a context is available.

## 3.1 Finetuning on Context-augmented Chats

For a conversation $C$ of length $L$ with segments $\{(x_t, y_t, c_t)\}_{i=1}^{L}$, where $x_t$ is a text generated by the agent or the customer, $y_t$ is its reference translation in the target language, and $c_t$ is the preceding bilingual context, we train the model to minimize the cross-entropy loss using the input prompt shown in Figure 2:

$$\mathcal{L} = -\log P(y_t | x_t, c_t). \quad (1)$$

The context $c_t$ includes all previous turns of the conversation, capturing important discourse-level information such as pronoun references, formality, and other pragmatic elements that influence the translation. For the first turn, no context is available, so the prompt reduces to the standard format used for zero-shot MT, as described in Alves et al. (2024). We train TOWERCHAT by finetuning TOWERBASE 7B on the concatenation of TOWERBLOCKS and the entire training dataset of the shared task, using context-aware prompts. This endows the model with the capacity to better understand and leverage conversational context, enabling it to generate high-quality translations.

## 3.2 QAD with Context-aware Metrics

Decoding strategies informed by translation quality metrics such as Minimum Bayes Risk Decoding (MBR) and Tuned Reranking (TRR) have been shown to consistently improve output quality over greedy decoding (Fernandes et al., 2022; Freitag et al., 2022; Nowakowski et al., 2022; Farinhas et al., 2023). In QAD, the goal is to find

a translation among a set of candidates that maximizes the expected utility function, often measured using an MT metric like reference-based COMET. Recently, Agrawal et al. (2024) showed that context-aware MT metrics correlate better with human judgments compared to their non-contextual counterparts, especially when evaluating out-of-English chat translations. The context-aware versions of COMET (Vernikos et al., 2022; Agrawal et al., 2024) compute quality scores for a source-reference-hypothesis tuple, $(x, y, \hat{y})$, using the representations extracted from a context-augmented input, $([c; x], [c; y], [c; \hat{y}])$.

As such, we use CONTEXTCOMET for MBR decoding in our submission. For a given source text $x$, the previous bilingual context, $c$, and a set of candidate translations sampled from the model $\mathcal{Y}$, the utility of each candidate $\hat{y} \in \mathcal{Y}$, is given by

$$u = \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \text{CONTEXTCOMET}([c; x], [c; y], [c; \hat{y}]).$$
$$(2)$$

The best translation is selected using:

$$y_{\text{mbr}} := \arg\max_{\hat{y} \in \mathcal{Y}} [u(\hat{y}, \mathcal{Y})]. \quad (3)$$

This enables the model to select a translation amongst alternative hypotheses, potentially leading to more accurate and contextually appropriate outputs. QAD with TOWERCHAT serves as our primary submission to the Chat Shared Task.

## 4 Experimental Configurations

**Baselines.** We report the shared task's official baseline: NLLB-3.3B with beam search decoding (beam width: 4). Additionally, we report greedy decoding results with TOWERINSTRUCT-7B, and TOWER-V2-70B, the strongest TOWER model. The former serves as a direct baseline for our method, while the latter is a state-of-the-art baseline for MT.

**TOWERCHAT.** We report greedy and QAD results with the TOWERCHAT-7B model. For

| MODEL | EN-XX | | | | XX-EN | | | |
|---|---|---|---|---|---|---|---|---|
| | CHRF↑ | COMET↑ | MetricX↓ | Context-QE↑ | CHRF↑ | COMET↑ | MetricX↓ | Context-QE↑ |
| **Baselines** | | | | | | | | |
| NLLB | 59.78 9 | 88.61 8 | 1.04 6 | 4.95 6 | 70.76 9 | 88.16 7 | 0.74 5 | 5.06 6 |
| TOWERINSTRUCT 7B (0-shot) | | | | | | | | |
| *w/o context* | 64.95 8 | 91.69 6 | 0.38 3 | 16.29 4 | 76.04 6 | 92.17 5 | 0.56 4 | 15.73 4 |
| *w/ context* | 63.39 9 | 91.09 7 | 0.49 5 | 14.53 5 | 74.32 8 | 91.36 6 | 0.60 4 | 14.67 4 |
| TOWERINSTRUCT 7B (5-shot) | | | | | | | | |
| *w/o context* | 65.20 8 | 91.75 6 | 0.39 3 | 16.62 3 | 75.84 7 | 92.22 5 | 0.54 4 | 15.97 3 |
| *w/ context* | 63.62 9 | 91.03 7 | 0.50 5 | 15.02 4 | 73.52 9 | 91.64 6 | 0.59 4 | 14.67 4 |
| TOWER-V2 70B (5-shot) | | | | | | | | |
| *w/o context* | 68.26 5 | 92.68 4 | 0.30 2 | 18.24 2 | 77.17 4 | 92.69 3 | 0.47 2 | 17.71 1 |
| *w/ context* | 68.26 6 | 92.50 4 | 0.30 2 | 17.53 2 | 76.03 6 | 92.37 4 | 0.46 2 | 17.28 2 |
| **TOWERCHAT** | | | | | | | | |
| *w/o context* | 71.68 5 | 93.01 4 | 0.32 3 | 16.77 3 | 77.97 3 | 92.72 4 | 0.51 3 | 16.40 3 |
| *w/ context* | 75.93 3 | 93.63 3 | 0.32 3 | 16.61 3 | 78.87 2 | 93.01 3 | 0.47 2 | 16.15 3 |
| + QAD (COMET) | 76.36 2 | **94.18** 1 | **0.25** 2 | **18.78** 1 | **78.92** 2 | **93.39** 1 | **0.44** 1 | 18.18 1 |
| + QAD (CONTEXTCOMET) | **76.56** 1 | 94.05 2 | 0.26 2 | 18.68 1 | **78.92** 2 | 93.24 2 | **0.44** 1 | **18.24** 1 |

Table 2: Main Results on Official Test Set: QAD with TOWERCHAT outperforms all baselines across the board. Models with statistically significant performance improvements are grouped in quality clusters

QAD, we perform MBR with COMET or CON-TEXTCOMET on 100 candidates obtained via epsilon sampling with $\epsilon = 0.02$ (Hewitt et al., 2022).

**Instruction settings.** To assess whether systems can properly leverage conversational context, we prompt the LLM-based MT with two instruction formats (see Figure 2): 1) **w/o context**, where the model is prompted without any conversational context (without the purple highlighted text). 2) **w/ context**, where the entire previous bilingual conversation is provided as the context in the prompt.[3]

**Evaluation.** We report the final results on the shared task's test set on all ten language pairs. As exemplified in Figure 1, ambiguous contextual phenomena often arise in Chat MT that require nuanced evaluation. As such, we leverage three types of assessments: 1) automatic metrics for measuring overall translation quality – two neural and one lexical – COMET-22 (Rei et al., 2022), CHRF (Popović, 2015) and METRICX-XL (Juraska et al., 2023); 2) a reference-free neural metric that uses context for quality assessment, CONTEXT-QE (Agrawal et al., 2024); 3) F1-score on MuDA tags for measuring whether models correctly resolve lexical ambiguities (Fernandes et al., 2023). Considering MET-RICX, CHRF, and MuDA is crucial in our case, as COMET may favor the QAD strategies we use.

On Tables 2 and 3 we report performance clusters based on statistically significant performance

gaps at a 95% confidence threshold.[4] We create perlanguage groups for systems with similar performance, following Freitag et al. (2023), and obtain system-level rankings using a normalized Borda count (Colombo et al., 2022), which is defined as an average of the obtained clusters. Note that a first cluster will not exist if no model significantly outperforms all others on a majority of languages.

## 5 Main Results

Table 2 presents the average results for EN→XX and XX→EN translation directions. TOWERCHAT with QAD outperforms all baselines across all settings on automatic metrics and human evaluation.

**TOWERCHAT leverages context more adeptly than TOWERINSTRUCT.** Our primary goal in this task was to create a model that can effectively leverage context to generate high-quality translations with LLMs. As shown in Table 2, TOWER-CHAT consistently outperforms TOWERINSTRUCT across all settings, language pairs and evaluation metrics. Furthermore, TOWERCHAT shows an average improvement of 4 CHRF points for en-xx when using context (*w/ context*), compared to a context-agnostic prompt (*w/o context*).[5] This trend also holds when evaluating translation quality using the primary metric, COMET, for 8 out of 10 language

---

[3]Note that {target_seg} is unavailable during inference and the model is asked to perform prompt completion.

[4]For segment-level metrics, such as COMET, we perform significance testing at the segment level. For CHRF, we substitute segment-level scores with corpus-level scores calculated over 100 random samples, each with a size equal to 50% of the total number of segments.

[5]The improvement is statistically significant with a 92.1% accuracy (Kocmi et al., 2024).

| MODEL | EN-XX | | | | | XX-EN | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | FR | PT | KO | NL | DE | FR | PT | KO | NL |
| **Baselines** | | | | | | | | | | |
| NLLB | 90.56 7 | 91.06 6 | 86.33 9 | 87.26 9 | 87.86 8 | 89.03 6 | 89.18 6 | 86.1 8 | 88.05 9 | 88.45 8 |
| TOWERINSTRUCT 7B (0-shot) | | | | | | | | | | |
| *w/o context* | 91.71 5 | 91.89 5 | 91.9 7 | 91.64 5 | 91.3 7 | 92.08 4 | 92.78 2 | 90.43 7 | 93.13 6 | 92.45 5 |
| *w/ context* | 91.48 6 | 91.08 6 | 90.79 8 | 91.13 7 | 91.0 7 | 91.33 5 | 91.89 5 | 90.63 6 | 91.88 7 | 91.08 7 |
| TOWERINSTRUCT 7B (5-shot) | | | | | | | | | | |
| *w/o context* | 91.75 5 | 91.75 5 | 92.32 6 | 91.41 6 | 91.55 6 | 92.06 4 | 92.28 4 | 90.63 6 | 93.55 5 | 92.6 5 |
| *w/ context* | 91.41 6 | 90.88 6 | 90.85 8 | 90.45 8 | 91.58 6 | 92.06 4 | 92.14 5 | 90.82 5 | 90.89 8 | 92.29 6 |
| TOWER-V2 70B (5-shot) | | | | | | | | | | |
| *w/o context* | 92.81 2 | 92.21 4 | 93.06 5 | 92.55 4 | 92.76 5 | **92.68 1** | **93.23 1** | 91.46 4 | 93.08 6 | 92.98 3 |
| *w/ context* | 92.61 3 | 92.08 4 | 93.03 5 | 91.76 5 | 93.02 4 | 92.07 4 | 92.44 4 | 91.42 4 | 93.05 6 | 92.89 3 |
| **TOWERCHAT** | | | | | | | | | | |
| *w/o context* | 92.36 4 | 92.26 4 | 93.89 4 | 93.73 3 | 92.81 4 | 92.28 3 | 92.79 2 | 91.06 5 | 94.69 4 | 92.78 4 |
| *w/ context* | 92.74 2 | 92.64 3 | 94.53 3 | 94.16 2 | 94.09 3 | 92.24 3 | 92.67 3 | 92.09 3 | 94.98 3 | 93.06 3 |
| said    + QAD (COMET) | **93.28 1** | **93.13 1** | **94.91 1** | **95.01 1** | **94.54 1** | 92.58 1 | 92.95 2 | **92.63 1** | **95.32 1** | **93.49 1** |
|    + QAD (CONTEXTCOMET) | 93.22 1 | 92.96 2 | 94.76 2 | 94.96 1 | 94.36 2 | 92.48 2 | 92.71 3 | 92.46 1 | 95.16 2 | 93.38 2 |
| Official Rank (COMET) | 2nd | 1st | 1st | 1st | 1st | 2nd | 1st | 1st | 1st | 1st |
| Official Rank (Human) | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st |

Table 3: Main Results by COMET on Official Test Set by Language Pair. Models with statistically significant performance improvements are grouped in quality clusters.

| Model | Lexical Cohesion | Verb Form | Pronouns | Formality |
|---|---|---|---|---|
| NLLB | 72.43 | 52.97 | 72.31 | 56.44 |
| TOWERCHAT | 85.13 | 47.80 | 79.71 | 81.93 |
| QAD (COMET) | 85.94 | 61.22 | **80.56** | 82.46 |
| QAD (CONTEXTCOMET) | **86.21** | **64.38** | 79.28 | **83.16** |

Table 4: MuDA F1 results. On average, QAD with CONTEXTCOMET has the best F1 score.

| Models | EN→XX | xx→en |
|---|---|---|
| TOWERINSTRUCT-7B | 84.28 | 82.77 |
| TOWERCHAT-7B | 83.95 | 82.54 |

Table 5: COMET scores for TOWERINSTRUCT and TOWERCHAT on the WMT23 test set.

pairs as shown in Table 3. We attribute this to the inclusion of context-augmented Chat MT instruction dataset in TOWERCHAT's training, highlighting the effectiveness of in-domain fine-tuning.

**QAD results in consistent gains over greedy decoding, surpassing 70B models.** The highest-quality translations according to all metrics considered are obtained after performing QAD with COMET or CONTEXTCOMET on top of TOWER-CHAT-7B, even outperforming the much larger TOWER-V2-70B, which uses few-shot examples. Moreover, QAD closes the gap in quality as measured by METRICX and CONTEXT-QE between TOWERCHAT-7B (greedy) and TOWER-V2-70B models, demonstrating that advanced inference techniques can effectively make smaller models competitive against much larger ones.

**Context-aware QAD improves MUDA F1 over Context-agnostic QAD.** While all neural and lexical metrics indicate that QAD with CON-TEXTCOMET and COMET perform comparably, these metrics may not fully capture nuanced dif-

ferences in translation quality. To address this, we evaluate MUDA F1 accuracy scores for a subset of models in Table 4. The results show that QAD with CONTEXTCOMET consistently outperforms QAD with COMET across all dimensions, except pronouns. Our qualitative analysis suggests that the pronoun accuracy might have been lower due to potential paraphrasing. Coupled with the previous results, these findings strongly motivate further exploration of QAD with context-aware metrics.

**Finetuning on Chat data does not degrade general translation capabilities.** To ensure that adding chat MT dataset in the mix does not impact the generic translation capabilities of LLMs, we report COMET on the standard WMT23 benchmark (Kocmi et al., 2023) averaged across EN→XX and XX→EN directions for TOWERINSTRUCT and TOWERCHAT in Table 5. TOWERCHAT suffers only minor degradation ($-0.3$) relative to TOW-ERINSTRUCT, validating the viability and effectiveness of our finetuning approach.

| SYSTEM | EN-DE | | | EN-FR | | | EN-NL | | | EN-PT | | | EN-KO | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | T (XX) | T (EN) | C | T (XX) | T (EN) | C | T (XX) | T (EN) | C | T (XX) | T (EN) | C | T (XX) | T (EN) | C |
| Baseline | 78.05 | 87.57 | 74.50 | 80.59 | 77.82 | 67.81 | 82.66 | 90.98 | 53.07 | 61.27 | 73.98 | 56.37 | 79.13 | 90.47 | 85.63 |
| Unbabel-IT | 89.42 | 92.74 | 84.22 | 90.24 | 90.00 | 79.62 | 98.16 | 97.40 | 92.22 | 82.04 | 82.37 | 78.00 | 93.39 | 96.31 | 93.21 |

Table 6: Human Evaluation results on the official test set. T and C represent aggregated turn-level and conversation-level direct-assessment scores respectively.

## 6 Human Evaluation

TOWERCHAT is the winner of the WMT24 Chat MT Shared Task across all language pairs according to human evaluation. Table 6 shows that our model significantly surpasses the baseline on both turn-level (T) and conversation-level (C) evaluations in all language directions. Notably, it reaches an average direct assessment score of $> 90$ at both turn-level and conversation-level for EN-FR, EN-NL, and EN-KO translation pairs. The victory on conversation-level evaluation outlines the superior capacity of TOWERCHAT to incorporate bilingual conversational context when translating.

That said, there is a visible drop between turn-level and conversation-level scores, leaving room for improvement on how well TOWERCHAT leverages context for translation. In future work, we wish to explore thoroughly under what circumstances context is useful to produce a better translation, and to what extent TOWERCHAT can leverage it appropriately.

## 7 Conclusion

In this work, we present two strategies for improving context usage for bilingual chat translation using LLMs. Our training strategy involves finetuning LLMs on context-augmented instructions resulting in higher-quality translations during inference when using bilingual context. Second, we propose a novel quality-aware decoding strategy with a context-aware metric (CONTEXTCOMET) that significantly improves translation quality across the board, surpassing a state-of-the-art 70B translation model and all other baselines. Our findings show successful usage of contextual information as measured by MUDA in resolving ambiguities for the highly contextual domain of chat translation. Crucially, our system finished first in human evaluation across all the shared task's language pairs.

## Acknowledgments

## References

Sweta Agrawal, Amin Farajian, Patrick Fernandes, Ricardo Rei, and André FT Martins. 2024. Is context helpful for chat translation evaluation? *arXiv preprint arXiv:2403.08314*.

Duarte Alves, Nuno Guerreiro, João Alves, José Pombal, Ricardo Rei, José de Souza, Pierre Colombo, and Andre Martins. 2023. Steering large language models for machine translation with finetuning and in-context learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11127–11148, Singapore. Association for Computational Linguistics.

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *Preprint*, arXiv:2402.17733.

Pierre Colombo, Nathan Noiry, Ekhine Irurozki, and Stéphan Clémençon. 2022. What are the best systems? new perspectives on nlp benchmarking. In *Advances in Neural Information Processing Systems*.

M Amin Farajian, António V Lopes, André FT Martins, Sameen Maruf, and Gholamreza Haffari. 2020. Findings of the wmt 2020 shared task on chat translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 65–75.

Ana C Farinha, M. Amin Farajian, Marianna Buchicchio, Patrick Fernandes, José G. C. de Souza, Helena Moniz, and André F. T. Martins. 2022. Findings of the WMT 2022 shared task on chat translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 724–743, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

António Farinhas, José de Souza, and Andre Martins. 2023. An empirical study of translation hypothesis ensembling with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11956–11970, Singapore. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Patrick Fernandes, Kayo Yin, Emmy Liu, André Martins, and Graham Neubig. 2023. When does translation require context? a data-driven, multilingual exploration. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 606–626, Toronto, Canada. Association for Computational Linguistics.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of wmt23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, Singapore. Association for Computational Linguistics.

Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2024. Exploring human-like translation strategy with large language models. *Transactions of the Association for Computational Linguistics*, 12:229–246.

John Hewitt, Christopher Manning, and Percy Liang. 2022. Truncation sampling as language model desmoothing. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3414–3427, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Marzena Karpinska and Mohit Iyyer. 2023. Large language models effectively leverage document-level context for literary translation, but critical errors persist. In *Proceedings of the Eighth Conference on Machine Translation*, pages 419–451, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024. Navigating the metrics maze: Reconciling score magnitudes and accuracies. *arXiv preprint arXiv:2401.06760*.

Samuel Läubli, Rico Sennrich, and Martin Volk. 2018. Has machine translation achieved human parity? a case for document-level evaluation. *arXiv preprint arXiv:1808.07048*.

Artur Nowakowski, Gabriela Pałka, Kamil Guttmann, and Mikołaj Pokrywka. 2022. Adam Mickiewicz University at WMT 2022: NER-assisted and quality-aware neural machine translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 326–334, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Raphael Reinauer, Patrick Simianer, Kaden Uhlig, Johannes E. M. Mosig, and Joern Wuebker. 2023. Neural machine translation models can learn to be few-shot learners. *Preprint*, arXiv:2309.08590.

Antonio Toral, Sheila Castilho, Ke Hu, and Andy Way. 2018. Attaining the unattainable? reassessing claims of human parity in neural machine translation. *arXiv preprint arXiv:1808.10432*.

Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any

pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 118–128, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.

Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li, Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhiwei Cao, Binbin Xie, Tianxiang Hu, Shangjie Li, Binyuan Hui, Bowen Yu, Dayiheng Liu, Baosong Yang, Fei Huang, and Jun Xie. 2023. Polylm: An open source polyglot large language model. *Preprint*, arXiv:2307.06018.

Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhengrui Ma, Yan Zhou, Langlin Huang, Mengyu Bu, Shangtong Gui, Yunji Chen, Xilin Chen, and Yang Feng. 2023. Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models. *Preprint*, arXiv:2306.10968.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2765–2781, Mexico City, Mexico. Association for Computational Linguistics.

# A  Validation Results

| MODEL | EN-XX | | XX-EN | |
|---|---|---|---|---|
| | CHRF | COMET | CHRF | COMET |
| **Baselines** | | | | |
| NLLB 3.3B | 58.41 | 86.97 | 65.39 | 85.51 |
| TOWERINSTRUCT 7B (0-shot) | | | | |
| *w/o context* | 63.69 | 90.69 | 71.57 | 90.62 |
| *w/ context* | 63.51 | 90.53 | 70.16 | 89.84 |
| TOWER-V2 70B (5-shot) | | | | |
| *w/o context* | 67.08 | 91.95 | 73.41 | 91.41 |
| *w/ context* | 66.85 | 91.69 | 71.87 | 90.94 |
| **TOWERCHAT** | | | | |
| *w/o context* | 70.63 | 92.21 | 73.42 | 91.13 |
| *w/ context* | 74.17 | 92.76 | 73.81 | 91.35 |
| + QAD (COMET) | 74.49 | **93.49** | 73.93 | **91.85** |
| + QAD (CONTEXTCOMET) | **74.54** | 93.31 | **74.15** | 91.70 |

Table 7: Results on the Validation Set: TOWERCHAT with QAD outperforms all baselines.

# B  Test Results by Language Pair

| MODEL | EN-XX | | | | | XX-EN | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | FR | PT | KO | NL | DE | FR | PT | KO | NL |
| **Baselines** | | | | | | | | | | |
| NLLB | 70.22 | 76.03 | 58.60 | 34.50 | 59.55 | 71.79 | 76.37 | 67.13 | 69.87 | 68.62 |
| TOWERINSTRUCT 7B (0-shot) | | | | | | | | | | |
| *w/o context* | 71.81 | 74.59 | 72.26 | 43.18 | 62.90 | 77.57 | 79.02 | 72.06 | 75.73 | 75.80 |
| *w/ context* | 71.16 | 74.38 | 68.50 | 41.70 | 61.23 | 75.68 | 78.31 | 71.83 | 72.63 | 73.15 |
| TOWERINSTRUCT 7B (5-shot) | | | | | | | | | | |
| *w/o context* | 71.38 | 74.72 | 72.59 | 42.76 | 64.55 | 76.64 | 78.67 | 71.68 | 76.23 | 75.96 |
| *w/ context* | 71.48 | 73.66 | 66.15 | 40.94 | 65.86 | 75.05 | 77.56 | 70.39 | 70.87 | 73.74 |
| TOWER-V2 70B (5-shot) | | | | | | | | | | |
| *w/o context* | 75.58 | 75.53 | 75.02 | 47.16 | 68.00 | 78.07 | 80.49 | 73.58 | 76.57 | 77.12 |
| *w/ context* | 74.60 | 75.28 | 74.05 | 46.99 | 70.38 | 77.54 | 77.63 | 73.16 | 75.69 | 76.10 |
| **TOWERCHAT** | | | | | | | | | | |
| *w/o context* | 74.04 | 77.12 | 79.71 | 57.63 | 69.91 | 79.31 | 79.36 | 74.00 | 80.17 | 77.01 |
| *w/ context* | 76.41 | 79.97 | 82.24 | 61.27 | 79.78 | 79.91 | 79.26 | 75.72 | 81.30 | 78.15 |
| + QAD (COMET) | 77.09 | 80.34 | 82.25 | 61.79 | 80.33 | 79.70 | 78.78 | 75.88 | 81.56 | 78.67 |
| + QAD (CONTEXTCOMET) | 77.23 | 80.51 | 82.55 | 62.29 | 80.25 | 79.87 | 78.57 | 76.01 | 81.57 | 78.60 |

Table 8: Results by CHRF (higher is better) on Official Test Set by Language Pair.

| MODEL | EN-XX | | | | | XX-EN | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | FR | PT | KO | NL | DE | FR | PT | KO | NL |
| **Baselines** | | | | | | | | | | |
| NLLB | 0.62 | 0.38 | 1.57 | 1.51 | 1.13 | 0.65 | 0.70 | 1.07 | 0.59 | 0.68 |
| TOWERINSTRUCT 7B (0-shot) | | | | | | | | | | |
| *w/o context* | 0.28 | 0.23 | 0.43 | 0.57 | 0.37 | 0.50 | 0.53 | 0.86 | 0.37 | 0.53 |
| *w/ context* | 0.38 | 0.29 | 0.69 | 0.60 | 0.49 | 0.56 | 0.55 | 0.74 | 0.46 | 0.69 |
| TOWERINSTRUCT 7B (5-shot) | | | | | | | | | | |
| *w/o context* | 0.29 | 0.25 | 0.39 | 0.62 | 0.39 | 0.50 | 0.55 | 0.79 | 0.37 | 0.52 |
| *w/ context* | 0.35 | 0.32 | 0.69 | 0.75 | 0.39 | 0.54 | 0.59 | 0.72 | 0.60 | 0.51 |
| TOWER-v2 70B (5-shot) | | | | | | | | | | |
| *w/o context* | 0.24 | 0.22 | 0.27 | 0.45 | 0.30 | 0.48 | 0.46 | 0.63 | 0.33 | 0.45 |
| *w/ context* | 0.25 | 0.21 | 0.28 | 0.45 | 0.29 | 0.50 | 0.48 | 0.58 | 0.31 | 0.42 |
| **TOWERCHAT** | | | | | | | | | | |
| *w/o context* | 0.27 | 0.24 | 0.29 | 0.42 | 0.37 | 0.50 | 0.51 | 0.71 | 0.33 | 0.52 |
| *w/ context* | 0.34 | 0.26 | 0.27 | 0.45 | 0.27 | 0.47 | 0.48 | 0.60 | 0.30 | 0.48 |
| + QAD (COMET) | 0.30 | 0.22 | 0.24 | 0.31 | 0.21 | 0.46 | 0.46 | 0.55 | 0.27 | 0.45 |
| + QAD (CONTEXTCOMET) | 0.31 | 0.22 | 0.24 | 0.29 | 0.23 | 0.47 | 0.47 | 0.56 | 0.27 | 0.45 |

Table 9: Results by METRICX (lower is better) on Official Test Set by Language Pair.

| MODEL | EN-XX | | | | | XX-EN | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | DE | FR | PT | KO | NL | DE | FR | PT | KO | NL |
| **Baselines** | | | | | | | | | | |
| NLLB | 15.56 | 1.24 | -5.51 | 4.11 | 9.35 | 19.09 | 0.77 | -6.75 | 4.13 | 8.04 |
| TOWERINSTRUCT 7B (0-shot) | | | | | | | | | | |
| *w/o context* | 21.84 | 8.96 | 9.11 | 19.73 | 21.83 | 23.41 | 7.46 | 7.49 | 18.66 | 21.64 |
| *w/ context* | 21.26 | 7.22 | 6.89 | 17.50 | 19.79 | 22.52 | 7.45 | 7.12 | 17.72 | 18.53 |
| TOWERINSTRUCT 7B (5-shot) | | | | | | | | | | |
| *w/o context* | 21.75 | 9.41 | 10.47 | 19.50 | 21.95 | 23.37 | 8.29 | 8.51 | 18.35 | 21.33 |
| *w/ context* | 21.83 | 8.20 | 8.17 | 15.22 | 21.68 | 22.93 | 6.75 | 7.16 | 15.49 | 21.02 |
| TOWER-v2 70B (5-shot) | | | | | | | | | | |
| *w/o context* | 23.42 | 10.47 | 12.38 | 20.84 | 24.07 | 25.21 | 9.77 | 10.26 | 20.08 | 23.21 |
| *w/ context* | 23.13 | 9.74 | 12.30 | 18.91 | 23.56 | 25.11 | 9.91 | 9.72 | 18.88 | 22.80 |
| **TOWERCHAT** | | | | | | | | | | |
| *w/o context* | 22.31 | 9.15 | 10.55 | 20.08 | 21.75 | 24.12 | 7.72 | 8.81 | 19.48 | 21.85 |
| *w/ context* | 22.39 | 8.69 | 11.36 | 18.58 | 22.05 | 24.28 | 7.45 | 9.06 | 17.96 | 21.97 |
| + QAD (COMET) | 24.27 | 10.92 | 13.01 | 21.65 | 24.04 | 26.12 | 9.67 | 10.77 | 21.02 | 23.31 |
| + QAD (CONTEXTCOMET) | 24.41 | 10.67 | 12.74 | 21.64 | 23.93 | 26.15 | 10.00 | 10.59 | 21.08 | 23.39 |

Table 10: Results by CONTEXT-QE (higher is better) on Official Test Set by Language Pair.

# C  MUDA F1 Scores by Language Pair

Figure 3: MUDA accuracy scores by LPs. Plots are left empty for the cases MUDA does not return tags (e.g., verb form for Korean).

# Optimising LLM-Driven Machine Translation with Context-Aware Sliding Windows

**Xinye Yang, Yida Mu, Kalina Bontcheva, Xingyi Song**
School of Computer Science, The University of Sheffield
{xyang138, y.mu, k.bontcheva, x.song}@sheffield.ac.uk

## Abstract

This paper describes SheffieldGATE's submission to WMT 2024 Chat Shared Translation Task. We participate in three language pairs: English-German, English-Dutch, and English-Portuguese (Brazil). In this work, we introduce a context-aware sliding window decoding method to track dependencies between chat messages. We fine-tune a large pre-trained language model based on the training data provided by the shared task Our experiments (i) compare the model performance between multilingual and bilingual fine-tuning and (ii) assess the impact of different window sizes. Our experimental results demonstrate that utilising contextual information yields superior performance in document-level translation compared to translating documents as isolated text segments, and that models fine-tuned with multilingual data perform better than those fine-tuned with bilingual data.

## 1 Introduction

Translating chat text is an important and challenging application of machine translation technology (Farajian et al., 2020; Farinha et al., 2022). The purpose of this task is to build a translation model that addresses the challenges of multilingual customer support for multinational companies. In informal conversations, people often use abbreviations and incomplete sentences and may include spelling errors, leading to significant noise in the dialogue text (Varnhagen et al., 2010). These factors complicate the translation of such texts, a challenge that traditional machine translation methods struggle to address (Almansor et al., 2020).

Recently, large language models (LLMs) have gradually taken over the mainstream in the field of natural language processing (Ouyang et al., 2022). LLMs have demonstrated impressive capabilities in a wide range of domains such as computational social science (Mu et al., 2024), question answering (Tan et al., 2023), and machine translation(Wang et al., 2023). Their ability to be well robust to noise in the input data provides new ideas to address the challenges of chat translation.

At the sentence level, Neural Machine Translation (NMT), represented by pre-trained large language models, is approaching the quality of professional human translations or even exceeding that of crowd-sourced non-professional translations in a few resource-rich languages (Hassan et al., 2018). For document-level translation, NMT systems still have certain errors that are difficult to detect in sentence-level translation (Läubli et al., 2018). Such as language ambiguity, which frequently results in numerous translation errors. Depending on the context, a single word or phrase can have multiple meanings (Abeysiriwardana and Sumanathilaka, 2024). Without the use of contextual information, problems including co-reference (Guillou and Hardmeier, 2016), lexical cohesion (Carpuat, 2009), or lexical disambiguation (Rios Gonzales et al., 2017) will be difficult to address (Jin et al., 2023).

In this work, we focus on modelling strategies based on contextual information. Our submission is based on an existing pre-trained model and fine-tuned using multilingual chat data, behaviour without incorporating additional contextual information during the fine-tuning process. We implemented context-aware sliding windows for the inference stage to perform translation tasks. We also conducted the following experiments (i) to compare the performance difference between using multilingual data and bilingual data in the fine-tuning process and (ii) the impact of window size, or the extent of contextual information, on the quality of translation.

With this study, we aim to shed light on the great potential of large language models for machine translation tasks and their ability to utilise contextual information for document-level translation and learn from migrating across linguistic data.

| Language Pair | Train | Val. | Test |
|---|---|---|---|
| EN <-> DE | 17,805 | 2,569 | 2,041 |
| EN <-> FR | 15,027 | 3,007 | 2,091 |
| EN <-> PT-BR | 15,092 | 2,550 | 2,040 |
| EN <-> KO | 16,122 | 1,935 | 1,982 |
| EN <-> NL | 15,463 | 2,549 | 2,015 |

Table 1: Number of source segments in the released dataset.

## 2 Data

The dataset for this task comprises authentic bilingual customer support conversations across five language pairs: English-German, English-French, English-Korean, English-Dutch, and English-Portuguese (Brazil). Table 1 displays the number of training, validation and test samples for each language pair in the dataset.

### 2.1 Dataset Characteristics

The chat content flows freely without strict format constraints, authentically reflecting the characteristics of real conversations. This natural language use includes incomplete sentences, interjections, and context-dependent responses, which, while representative of genuine dialogue, increases the complexity of processing and translation.

## 3 System Description

### 3.1 Context-Aware Sliding Window

To effectively utilise contextual information, we use a context-aware sliding window mechanism. This approach allows model to consider context sentences when translating each individual message, thereby enhancing the overall coherence and accuracy of the translation. In addition, we improve translation efficiency by reusing the Key-Value (KV) cache. KV caching is a crucial technique in transformer models, involves storing and reusing previously computed Key and Value matrices in the self-attention mechanism. This method significantly enhances inference speed by eliminating redundant calculations, particularly beneficial for long sequences or auto-regressive generation tasks such as machine translation. It enables the model to efficiently leverage information from the source language when generating the target sequence, substantially reducing computational overhead, especially for longer texts.

**Structure of the Sliding Window** Our context-aware sliding window comprises four key components:

- Task Description: Provides the model with clear instructions about the translation task.

- Source language tag: Identifies the beginning of the original text.

- Original Text: Contains the message to be translated along with its context.

- Target Language Label: Indicates the end of the original text and directs the model to give the translation.

Figure 1 illustrates the structure of the Context-Aware Sliding Window. This system comprises a task description and a window containing a sequence of source sentences, which together function as input to the model. The model generates new translations based on the contextual information available within the window. After each translation is produced, it is inserted into the list of translated sentences, and the window shifts to incorporate new source sentences. If the number of sentences in the source text window exceeds a predefined limit, the earliest sentence in the window is removed to maintain the set window size. This sliding mechanism ensures that the model consistently has track dependencies throughout the translation process.



Figure 1: Context-Aware Sliding Window

**Prompt** We used the following prompt for translation:
You are a translation specialist serving multinational companies. Your task is to translate the given text from [source language] to [target language]. Provide the translation result in [target language] directly without including any additional content.

**Workflow** The operation of our context-aware sliding window can be described as follows:

- Initialisation: The sliding window starts empty and gradually fills with sentences from the chat log up to the predefined window size.

- Generation: The language model generates the translation for the most recent sentence, considering both the original sentences in the window and their existing translations.

- Window Shift: After generating a translation, the window shifts by one position. It incorporates the next sentence from the chat log and removes the earliest one and its corresponding translation if the window is full. If the translation direction of the next sentence changes , the windows storing the original text and the translated text are swapped. This approach allows for seamless handling of bidirectional translations within the same conversation, maintaining context in both languages.

- Iteration: Steps 2 and 3 are repeated until all sentences in the chat log have been processed.

The workflow of the context-aware sliding window is illustrated in pseudocode in Algorithm 1.

**Advantages**   This approach offers several benefits:

- Improved Coherence: By considering the surrounding context, the model can maintain better consistency in tone, style, and terminology across the translation.

- Enhanced Accuracy: Contextual information helps resolve ambiguities and choose more appropriate translations for words or phrases with multiple meanings.

## 4   Experiments

In this section, we describe the experiments conducted to select the fine-tuning strategy and determine the optimal window size for our system. The hyperparameters used in this experiment are listed in Table 2. All experiments were executed on a single Nvidia A100 GPU equiped with 40GB of memory.

Three evaluation metrics are used in this experiment, aligned with the automatic evaluation metrics of the shared task, they are:

- BLEU (Bilingual Evaluation Understudy) (Papineni et al., 2002): Measures translation qual-

---

**Algorithm 1** Context-Aware Sliding Window Translation Algorithm with Bidirectional Support

1: **Initialise:**
2:     *source-window* ← [ ]
3:     *target-window* ← [ ]
4:     *window-size* ← predefined window size
5:     *translation-result* ← [ ]
6:     *current-direction* ← initial translation direction
7: **for** each *sentence* in *input text* **do**
8:       **if** *sentence-direction* ≠ *current-direction* **then**
9:             *source-window*, *target-window* ← *target-window*, *source-window*
10:            *current-direction* ← *sentence-direction*
11:      **end if**
12:      **if** *len(source-window)* < *window-size* **then**
13:            *source-window.append(sentence)*
14:            *translation* ← *Generate(source-window, target-window)*
15:            *target-window.append(translation)*
16:            *translation-result.append(translation)*
17:      **else**
18:            *source-window.pop(0)*
19:            *target-window.pop(0)*
20:            *source-window.append(sentence)*
21:            *translation* ← *Generate(source-window, target-window)*
22:            *target-window.append(translation)*
23:            *translation-result.append(translation)*
24:      **end if**
25: **end for**
26: **Output** *translation-result*

---

ity based on n-gram overlap between the candidate and reference translations. BLEU primarily assesses fluency and adequacy at the phrase level. It is widely used but may not always capture deeper semantic nuances.

- chrF (Character n-gram F-score) (Popović, 2015): Evaluates translation quality at the character level. It is particularly effective for capturing morphological accuracy and subtle differences in word forms. chrF is sensitive to grammatical correctness and precise word choice.

- COMET (Cross-lingual Optimised Metric for Evaluation of Translation) (Rei et al., 2020): A more recent metric that focuses on seman-

| Hyperparameter | Value |
|---|---|
| LoRA rank (r) | 8 |
| LoRA alpha | 16 |
| LoRA dropout | 0.05 |
| Learning rate | 2.5e-5 |
| Weight decay | 0.001 |
| Batch size | 8 |
| Training epochs | 10 |
| Warmup ratio | 0.3 |
| Max gradient norm | 0.3 |
| LR scheduler | Linear |

Table 2: Fine-tuning Hyperparameters

tic similarity between the source, translation, and reference. COMET uses contextual embedding to evaluate meaning preservation and overall translation quality, aiming to correlate better with human judgements.

### 4.1 Multilingual and bilingual Fine-tuning

Given the computational resources and time constraints, we choose the LLaMA3-8B instruct model (LLaMA) (Dubey et al., 2024) as our base model. We fine-tune LLaMA using Low-Rank Adaptation (LoRA) (Hu et al., 2022) with training and validation data provided by the shared task. We employed two distinct fine-tuning strategies, i.e., (i) multilingual fine-tuning and (ii) bilingual fine-tuning.

For the multilingual fine-tuning, we feed five language pairs simultaneously: English <-> German, English <-> French, English <-> Brazilian Portuguese, English <-> Korean, and English <-> Dutch. This strategy allows the model to learn from multiple languages concurrently and potentially leverage cross-lingual information.

In contrast, our bilingual strategy involved fine-tuning separate models for each language pair, using solely the training and validation data specific to that pair. This approach enables more focused adaptation to each language pair.

The motivation for employing these two strategies was to explore the cross-linguistic learning and transfer capabilities of large language models (Lample and Conneau, 2019). By comparing these approaches, we aim to investigate whether the model can extract universally applicable translation patterns and linguistic features from multiple language pairs, thereby potentially improving its performance on new language pairs.

The experiment results are shown in Table 3.

The multilingual fine-tuned models outperform bilingual fine-tuned models. This may be because multilingual dataset provide more samples than each bilingual datasets, offering a broader and more diverse set of data, which helps prevent the model from overfitting. Also, the model can learn translation patterns through transfer learning from other languages. Hence, in our final submission, the model was fine-tuned using the multilingual dataset.

### 4.2 Impact of Window Size

We also investigated the effect of different window sizes on the translation quality. In this work, the window size determines the amount of context available to the model during the translation process.

To that end, we conducted experiments with window sizes ranging from 1 to 3 sentences. For each window size, we translated five language pairs from the validation set provided by shared task and evaluated the results using automated metrics. Table 4 presents the detailed results for chrF, BLEU, and COMET scores across different window sizes and language pairs.

The window size used in our submission is 3. Our findings indicate that the translation quality generally improves as the window size increases, but the extent and nature of improvement varies across translation directions and metrics. We observe that the COMET metric tends to favour larger window sizes more consistently than chrF or BLEU.

COMET scores show improvement or maintain high performance with larger windows in 5 out of 6 translation directions (de-en, en-de, pt-br-en, nl-en, en-nl).

For en-pt-br, small windows have the best performance across all metrics. This unique behavior might be attributed to several factors. Firstly, the structural similarities between English and Brazilian Portuguese allow for effective translation with minimal context.(Angeli and Mota, 2023) The relatively simple morphology of English compared to Portuguese's more complex system might also contribute to this phenomenon. Additionally, the direct lexical correspondence between many English and Portuguese words could lead to high accuracy in word-to-word translations, which is particularly well-captured by chrF and BLEU metrics.

In contrast, chrF and BLEU metrics often peak

| Language Pair | multilingual | | | bilingual | | |
|---|---|---|---|---|---|---|
| | chrF | Bleu | COMET | chrF | Bleu | COMET |
| de->en | **67.45** | **44.46** | **88.13** | 65.63 | 41.11 | 86.55 |
| en->de | **60.95** | **35.74** | **86.41** | 60.03 | 34.82 | 85.59 |
| pt-br->en | **65.50** | **43.74** | **87.10** | 63.17 | 36.52 | 84.68 |
| en->pt-br | **66.94** | **42.02** | **89.43** | 65.21 | 39.43 | 87.67 |
| nl->en | **68.05** | **45.94** | **88.66** | 65.77 | 42.58 | 86.38 |
| en->nl | **62.26** | **35.94** | **89.29** | 59.65 | 32.41 | 87.09 |

Table 3: Translation Quality Metrics for Multilingual and bilingual Models. The highest scores for each metric are marked in bold.

| Language Pair | Window Size = 1 | | | Window Size = 2 | | | Window Size = 3 | | |
|---|---|---|---|---|---|---|---|---|---|
| | chrF | BLEU | COMET | chrF | BLEU | COMET | chrF | BLEU | COMET |
| de-en | 64.37 | 39.75 | 84.78 | **68.16** | **45.53** | **88.35** | 67.45 | 44.46 | 88.12 |
| en-de | 60.86 | 35.47 | 86.31 | **61.15** | **35.92** | 86.11 | 60.95 | 35.74 | **86.40** |
| pt-br-en | 62.96 | 39.24 | 83.44 | **65.62** | **44.24** | **87.35** | 65.50 | 43.74 | 87.10 |
| en-pt-br | **67.82** | **45.49** | **89.94** | 67.34 | 43.04 | 89.48 | 66.94 | 42.02 | 89.43 |
| nl-en | 64.02 | 40.94 | 83.23 | **68.15** | **48.01** | 88.26 | 68.05 | 45.94 | **88.66** |
| en-nl | 60.16 | 33.06 | 87.67 | 60.32 | 33.34 | 88.15 | **62.26** | **35.94** | **89.29** |

Table 4: Translation Quality Metrics for Different Window Sizes. The highest scores for each metric are marked in bold.

at window size 2 or even size 1 for some translation directions. For example, en-pt-br achieves its highest chrF and BLEU scores with window size 1. The en-nl pair is a notable exception, showing consistent improvement across all metrics as the window size increases.

This pattern suggests that the COMET metric may be more sensitive to the broader context provided by larger window sizes, while chrF and BLEU might prioritise local fluency or accuracy that can sometimes be captured effectively with smaller windows.

## 5 Conclusion

In this paper, we compared the performance of fine-tuning using multilingual data and bilingual data. Additionally, we conducted an ablation study by evaluating the translation quality with different window sizes. Our research indicates that fine-tuning models on multilingual data results in superior translation capabilities compared to fine-tuning on a single language. This approach could improve translation quality for low-resource languages. Furthermore, we also found that increasing the contextual information provided to the model can enhance its semantic performance in translation. Our future work will focus on:

- **Named Entity Handling** We plan to integrate a named entity recognition system and leverage external knowledge resources, such as Wikipedia, to ensure accurate translations of named entities.

- **Model Fine-tuning Comparison** We also aim to conduct a comparative analysis between fine-tuning the foundation model and the instruction-tuned model, exploring the trade-offs between general and task-specific performance.

## Acknowledgements

## References

Miuru Abeysiriwardana and Deshan Sumanathilaka. 2024. A survey on lexical ambiguity detection and word sense disambiguation. *arXiv preprint arXiv:2403.16129*.

Ebtesam Almansor, Ahmed Al-Ani, and Farookh Hussain. 2020. *Transferring Informal Text in Arabic as*

[1] https://www.veraai.eu/home

*Low Resource Languages: State-of-the-Art and Future Research Directions*, pages 176–187. Springer International Publishing.

Natália Pinheiro De Angeli and Mailce Borges Mota. 2023. Cross-linguistic priming effects during the comprehension of the passive voice: Two primes are enough. *Ilha do Desterro*, 76(3):17–39.

Marine Carpuat. 2009. One translation per discourse. In *Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions (SEW-2009)*, pages 19–27, Boulder, Colorado. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and et al. 2024. The llama 3 herd of models.

M Amin Farajian, António V Lopes, André FT Martins, Sameen Maruf, and Gholamreza Haffari. 2020. Findings of the wmt 2020 shared task on chat translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 65–75.

Ana C Farinha, M. Amin Farajian, Marianna Buchicchio, Patrick Fernandes, José G. C. de Souza, Helena Moniz, and André F. T. Martins. 2022. Findings of the WMT 2022 shared task on chat translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 724–743, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Liane Guillou and Christian Hardmeier. 2016. PROTEST: A test suite for evaluating pronouns in machine translation. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 636–643, Portorož, Slovenia. European Language Resources Association (ELRA).

Hany Hassan, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, Shujie Liu, Tie-Yan Liu, Renqian Luo, Arul Menezes, Tao Qin, Frank Seide, Xu Tan, Fei Tian, Lijun Wu, Shuangzhi Wu, Yingce Xia, Dongdong Zhang, Zhirui Zhang, and Ming Zhou. 2018. Achieving human parity on automatic chinese to english news translation. *CoRR*, abs/1803.05567.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Linghao Jin, Jacqueline He, Jonathan May, and Xuezhe Ma. 2023. Challenges in context-aware neural machine translation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15246–15263, Singapore. Association for Computational Linguistics.

Guillaume Lample and Alexis Conneau. 2019. Cross-lingual language model pretraining. *ArXiv*, abs/1901.07291.

Samuel Läubli, Rico Sennrich, and Martin Volk. 2018. Has machine translation achieved human parity? a case for document-level evaluation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4791–4796, Brussels, Belgium. Association for Computational Linguistics.

Yida Mu, Ben P. Wu, William Thorne, Ambrose Robinson, Nikolaos Aletras, Carolina Scarton, Kalina Bontcheva, and Xingyi Song. 2024. Navigating prompt complexity for zero-shot classification: A study of large language models in computational social science. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 12074–12086, Torino, Italia. ELRA and ICCL.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Annual Meeting of the Association for Computational Linguistics*.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Annette Rios Gonzales, Laura Mascarell, and Rico Sennrich. 2017. Improving word sense disambiguation in neural machine translation with sense embeddings. In *Proceedings of the Second Conference on Machine Translation*, pages 11–19, Copenhagen, Denmark. Association for Computational Linguistics.

Yiming Tan, Dehai Min, Y. Li, Wenbo Li, Nan Hu, Yongrui Chen, and Guilin Qi. 2023. Can chatgpt replace traditional kbqa models? an in-depth analysis of the question answering performance of the gpt llm family. In *International Workshop on the Semantic Web*.

Connie K Varnhagen, G Peggy McFall, Nicole Pugh, Lisa Routledge, Heather Sumida-MacDonald, and Trudy E Kwong. 2010. Lol: New language and

spelling in instant messaging. *Reading and writing*, 23:719–733.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.

# Context-Aware LLM Translation System Using Conversation Summarization and Dialogue History

**Mingi Sung[1],[*], Seungmin Lee[2],[*], Jiwon Kim[2], Sejoon Kim[1]**
[1]PwC Korea, Seoul, South Korea
[2]Yonsei University, Seoul, South Korea
*{mingi.sung@pwc.com, elplaguister@yonsei.ac.kr, jwkim808@yonsei.ac.kr, sejoon.s.kim@pwc.com}*
[*]Equal contribution.

## Abstract

Translating conversational text, particularly in customer support contexts, presents unique challenges due to its informal and unstructured nature. We propose a context-aware LLM translation system that leverages conversation summarization and dialogue history to enhance translation quality for the English-Korean language pair. Our approach incorporates the two most recent dialogues as raw data and a summary of earlier conversations to manage context length effectively. We demonstrate that this method significantly improves translation accuracy, maintaining coherence and consistency across conversations. This system offers a practical solution for customer support translation tasks, addressing the complexities of conversational text.

## 1 Introduction

The WMT 2024 Chat Shared Task addresses the unique challenges of translating conversational text, with a particular focus on customer support chats. Unlike formal or structured texts, conversations are typically spontaneous and casual, presenting several key challenges. First, the system must comprehend the dialogue's flow while accurately translating content from one language to another. Second, it is crucial to maintain logical continuity throughout entire conversations, preserving the context and intent of each exchange. Third, the task requires effectively handling the inherent noise and colloquial nature of chat data.

To tackle these challenges, we developed a context-aware LLM translation system that leverages both dialogue history and conversation summarization. Our approach is designed to maintain coherence and accuracy in translation by referencing two key elements: (1) 'History' field: The two most recent dialogues of the target conversation, provided as raw data. (2) 'History Summary' field: A concise summary (maximum 200 characters) of

earlier conversations, excluding the two most recent dialogues.

We utilize both history and history summary in our approach for the following reasons. Dialogues often require multi-turn information for accurate understanding, as context within a single turn can be insufficient or misleading. Furthermore, while referencing all previous conversations would be ideal, it is often prohibited by the context length limitations of LLMs. Our method addresses these challenges by using recent parts of the conversation verbatim and summarizing earlier parts of the dialogue, effectively reducing context length while maintaining overall contextual information.

Our approach is informed by previous research demonstrating the effectiveness of context-aware models. Current study has shown that stimulating LLMs to memorize small dialogue contexts first and then recursively produce new memory using previous memory helps the chatbot generate highly consistent response (Wang et al., 2023). The History-Aware Hierarchical Transformer (Zhang et al., 2022) also used historical information to improve the understanding of the current conversation context. The TiM (Think-in-Memory) framework (Liu et al., 2023), a LLM agent also recalls relevant thoughts from memory before generating response, and then integrates both historical and new thoughts to update the memory. By incorporating these insights, our system aims to produce translations that are not only accurate in language conversion but also maintain the coherent tone and appropriate word selection crucial in conversational contexts.

The Gemma-2-27B-it model (Team, 2024) is used as the foundation for our translation system, specifically focusing on the English-Korean language pair. Our experiments demonstrated that incorporating recent dialogues and previous dialogue summaries significantly improved translation performance compared to methods that did not utilize

this contextual information. We further refined our system by implementing more detailed instructions, which yielded additional improvements in translation accuracy. Furthermore, we used the GPT-4o mini (OpenAI, 2024) model for efficient conversation summarization. These combined methods resulted in substantial enhancements to overall translation quality, clearly demonstrating the effectiveness of our approach in boosting translation accuracy.

## 2 Methodology

### 2.1 Data Preparation

Among the five language pairs provided by the WMT 2024 Chat Shared Task, we selected the en-ko dataset for our experiments. The dataset provided by WMT consists of 16,122 training instances, 1,935 validation instances, and 1,982 test instances.

---

**Example 1**

```
"source_language": "ko",
"target_language": "en",
"source": "비밀번호 재설정 메일이 도착하지
않습니다.",
"reference": "I don't receive a password
reset email.",
"doc_id": "64619c16ab8523e90010b544",
"client_id": "0015800001EMz0vAAD",
"sender": "customer",
"history": "As I understand you are unable
to login to your account as it asks you to
reset the password and you are not getting
reset password email.",
"제가 알기로는 비밀번호를 재설정하라는 메시
지가 표시된 후 비밀번호 이메일을 재설정하지
않기 때문에 계정에 로그인할 수 없으십니다."
"Am I correct?",
"맞습니까?"

"Instruction": "You are tasked with
translating the following sentences from
Korean to English. These sentences are part
of conversations between a customer and a
customer service agent.\nWhen translating,
keep the following instructions in
mind:\n- Provide only the translation of
the 'source' text.\n- Keep the translated
text in a single line.\n- The context
involves a game user contacting a game
company's customer service center online.
Since the inquiries are typed, there
may be many typos. Please translate with
this in mind.\n- Consider the summary
of the previous conversation, referred
to as 'Dialogue Context', if it is
given.\n- Refer to the context from the
previous conversation if it is provided.\n-
Ensure your translations maintain the
intended meaning and tone of the original
```

---

```
dialogue.\nDialogue Context: The customer,
NAME-N, contacted PRS-ORG for help signing
in and reported not receiving a password
reset email.",
"History_summary": "Dialogue Context: The
customer, NAME-N, contacted PRS-ORG for
help signing in and reported not receiving
a password reset email.",
"System": "You are a professional
translator fluent in both Korean and
English."
```

---

During the data preparation process, we utilized several fields from the provided dataset and introduced two new ones to enhance context awareness. The original fields are `source_language`, `target_language`, `source`, `reference`, and `doc_id`. The newly inserted fields are `history` and `history_summary`. Example 1 shows the final preprocessed dataset used for model training. The `source_language` and `target_language` fields specify the language pair to be translated, `source` contains the text to be translated, `reference` provides the correct translation, and `doc_id` is used to uniquely identify each conversation session. The `history` and `history_summary` refer to the conversation context as described in the following sections.

### 2.1.1 History

The `history` field includes the raw data from the two previous dialogues of the targeted dialogue that needs to be translated. This information enables the model to capture the conversation's flow and maintain coherence in the generated translation.

### 2.1.2 History Summary

The `history_summary` field contains a concise summary of earlier conversations, excluding the two most recent dialogues. It helps to understand the overall context and background of the current conversation. For summary generation, we used the GPT-4o mini model with a prompt that limits the summary to a maximum of 200 characters. This approach allows the model to focus on the essential part of the previous content without being overwhelmed by excessive details.

### 2.1.3 System Prompt and Translation Instructions

To ensure consistency in model training and translation tasks, we developed a prompt strategy. We defined a base prompt that positioned the model as a professional translator fluent in both Korean and

English. We provided detailed guidelines for instruction, making the model focus on maintaining the context and tone of customer service dialogues, accounting for possible typos, and incorporating provided conversation summaries. To effectively guide the model, we supplied detailed prompts and instructions separately. The content of the system prompt and translation instructions can be found in the 'system' and 'instruction' fields in Example 1. This approach enhances translation quality by providing clear, context-specific guidance to the model.

## 2.2 Context-Aware LLM Translation System Using Conversation Summarization and Dialogue History

After data preparation, we used the provided Chat Template for Gemma's Instruct-Tuning to structure our data and fine-tuned the Gemma-2-27B-it model for translating customer support dialogues in Korean and English. We used DeepSpeed library to quantize the model and applied LoRA (Low-Rank Adaptation) for model compression. The key parameters of our training setup were as follows: per device train batch size of 4, gradient accumulation steps of 8, learning rate of 1.0e-4, 5 training epochs, cosine learning rate scheduler, warmup ratio of 0.1 and bfloat16 precision enabled. The GPU we used was NVIDIA H100 and the training process took about one hour to complete.

While optimizing our model's performance, we also addressed the unique challenges of dialogue translation. Context plays a pivotal role, significantly influencing the accuracy of interpreting and translating each turn. However, this importance presents a dual challenge. On one hand, preserving the conversation's history is crucial for coherent translations. On the other hand, as dialogues extend, managing this context becomes increasingly complex. Including all previous turns becomes impractical and can degrade the quality of subsequent translations.

To address these challenges, we implemented the following strategy:

- **Recent Dialogues (history)**: We utilized the history field to include the two most recent dialogue in their raw form. This preserves the immediate context necessary for accurate and coherent translations.

- **Dialogue Context (history_summary)**: For earlier parts, we provided a condensed summary of essential points, generated prior to inference. This helps the model grasp the broader context without being overwhelmed by excessive information (Bae et al., 2022).

This approach balances detailed immediate context with summarized background, allowing the model to capture both current dynamics and overall dialogue context.

Our prompt structure consists of three key components: a system role-play definition, a task instruction, and the sentence to be translated. This setup was critical for guiding the model's performance in customer support dialogue contexts. By structuring the context using natural language and leveraging the model's instruction-tuned capabilities, we aimed to enhance its ability to generate translations that are not only accurate but also contextually appropriate. This method allowed us to capture the natural flow and nuances of conversations more effectively.

## 3 Experimental Results and Application

The performance of our translation model was evaluated through both human assessment and various automated metrics. Table 1 shows translation performance scores by human evaluation. The 'sentence' columns indicate the evaluation scores for translation quality at the individual sentence level, while the 'Document' column reflects how well the translation maintains consistency and context across a full conversation. Our DeepText-Lab team received notably high evaluations, with scores of 91.35 for translating sentences from English to Korean and 95.71 for translating sentences from Korean to English. Our team also received a score of 90.04 at the document level. These results demonstrate our system's strong performance at both sentence and document levels.

Besides human evaluation, the performance on the test dataset was also evaluated using several automated metrics, including COMET, chrF, BLEU, and Contextual-Comet-QE. The results are summarized in Table 2. We achieved strong results across all metrics. The COMET score of 93.5 indicates high translation quality, while the chrF score of 66.0, BLEU score of 47.6, and Contextual-Comet-QE score of 0.161 demonstrate solid performance.

Beyond assessing overall performance, we explored how the inclusion of conversation history, history summaries, and detailed prompts influenced our model's translation quality. Table 3 illustrates

| Team | Sentence (en→ko) | Sentence (ko→en) | Document |
|---|---|---|---|
| unbabel+it | 93.39 | 96.31 | 93.21 |
| DeepText_Lab | 91.35 | 95.71 | 90.04 |
| DCUGenNLP | 89.71 | 96.15 | 89.83 |
| baseline | 79.13 | 90.47 | 85.63 |

Table 1: Human Evaluation of Sentence and Document-Level Translation (Test Dataset Results)

| Team | COMET | chrF | BLEU | C-COMET-QE |
|---|---|---|---|---|
| unbabel+it | 95.0 | 70.2 | 51.5 | 0.214 |
| DeepText_Lab | 93.5 | 66.0 | 47.6 | 0.161 |
| DCUGenNLP | 92.3 | 59.8 | 39.4 | 0.158 |
| baseline | 87.6 | 48.9 | 26.0 | 0.041 |

Table 2: Automatic Evaluation Using Multiple Metrics (Test Dataset Results)

| Configuration | Direction | COMET | chrF | BLEU | C-COMET-QE |
|---|---|---|---|---|---|
| **w/ recent dialogues and dialogue context** | en→ko | 0.916 | 51.52 | 32.97 | 0.138 |
| | ko→en | 0.893 | 61.57 | 40.73 | - |
| **w/o recent dialogues and dialogue context** | en→ko | 0.910 | 48.96 | 29.82 | 0.126 |
| | ko→en | 0.889 | 59.65 | 38.86 | - |
| **w/o prompt modification** | en→ko | 0.911 | 49.96 | 31.30 | 0.135 |
| | ko→en | 0.894 | 61.15 | 40.64 | - |

Table 3: Impact of Contextual Elements on Translation Performance (Validation Dataset Results)

the significant impact of these elements on various performance metrics, especially chrF and BLEU scores. The absence of these contextual elements led to a significant decrease in translation quality, emphasizing the importance of context preservation and precise guidance in producing high-quality translations.

In addition to the above evaluations, we also assessed our model's performance in terms of formality and lexical cohesion using the MuDA (Fernandes et al., 2023) framework. The results of these assessments are presented in Table 4 and Table 5. For formality, the model achieved a precision score of 73.0, and a recall of 35.1, resulting in an overall F1-Score of 47.4. For lexical cohesion, the model demonstrated strong performance with a precision of 70.5, and a recall of 73.8, leading to an F1-Score of 72.1.

| Team | Precision | Recall | F1 |
|---|---|---|---|
| unbabel+it | 69.4 | 44.2 | 54.0 |
| DeepText_Lab | 73.0 | 35.1 | 47.4 |
| DCUGenNLP | 25.5 | 18.2 | 21.2 |
| baseline | 50.0 | 10.4 | 17.2 |

Table 4: Formality Results

| Team | Precision | Recall | F1 |
|---|---|---|---|
| unbabel+it | 73.3 | 76.2 | 74.7 |
| DeepText_Lab | 70.5 | 73.8 | 72.1 |
| DCUGenNLP | 73.5 | 68.3 | 70.8 |
| baseline | 66.1 | 65.1 | 65.6 |

Table 5: Lexical Cohesion Results

## 4 Conclusion

We participated in the WMT English-Korean Chat Translation Task using the Gemma-2-27B-it model enhanced with dialogue history for context-aware translations. We effectively reduced the context length by summarizing earlier conversations and enhanced the model's translation performance by including the history of the two most recent dialogues and the summary of the previous dialogues, excluding the most recent two.

As a result, the translation performance has significantly improved, though there is still room for enhancement. Despite our team's high score, certain issues were identified in the generated translations. For instance, the Gemma 2 model occasionally produces translations in unexpected languages like Turkish, French, and Polish. This stems from the model's multilingual pretraining and presents an area for further exploration in future work.

# References

Sanghwan Bae, Donghyun Kwak, Soyoung Kang, Min Young Lee, Sungdong Kim, Yuin Jeong, Hyeri Kim, Sang-Woo Lee, Woomyoung Park, and Nako Sung. 2022. Keep me updated! memory management in long-term conversations. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, Abu Dhabi, UAE. Association for Computational Linguistics.

Patrick Fernandes, Kayo Yin, Emmy Liu, André Martins, and Graham Neubig. 2023. When does translation require context? a data-driven, multilingual exploration. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 606–626, Toronto, Canada. Association for Computational Linguistics.

Lei Liu, Xiaoyan Yang, Yue Shen, Binbin Hu, Zhiqiang Zhang, Jinjie Gu, and Guannan Zhang. 2023. Think-in-memory: Recalling and post-thinking enable llms with long-term memory. *arXiv preprint arXiv:2311.08719*.

OpenAI. 2024. Gpt-4o mini: advancing costefficient intelligence. *OpenAI Website*.

Team. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.

Qingyue Wang, Liang Ding, Yanan Cao, Zhiliang Tian, Shi Wang, Dacheng Tao, and Li Guo. 2023. Recursively summarizing enables long-term dialogue memory in large language models. *arXiv preprint arXiv:2308.15022*.

Tong Zhang, Yong Liu, Boyang Li, Zhiwei Zeng, Pengwei Wang, Yuan You, Chunyan Miao, and Lizhen Cui. 2022. History-aware hierarchical transformer for multi-session open-domain dialogue system. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, Abu Dhabi, UAE. Association for Computational Linguistics.

# Enhancing Translation Quality: A Comparative Study of Fine-Tuning and Prompt Engineering in Dialog-Oriented Machine Translation Systems. Insights from the MULTITAN-GML Team

Lichao Zhu[1], Maria Zimina-Poirot[1], Behnoosh Namdarzadeh[1], Nicolas Ballier[1,3] and Jean-Baptiste Yunès[2]

[1]CLILLAC-ARP, [2]IRIF, [3]LLF

Université Paris Cité, F-75013 Paris, France

Contact: `lichao.zhu@u-paris.fr`

## Abstract

For this shared task, we have used several machine translation engines to produce translations (en ⇔ fr) by fine-tuning a dialog-oriented NMT engine and having NMT baseline translations post-edited through prompt engineering. Our objectives are to assess the effectiveness of a fine-tuning strategy with a robust NMT model, to advance towards a comprehensive pipeline that covers the entire translation process (from fine-tuning and machine translation to automatic post-editing (APE)), and to evaluate the strengths and weaknesses of NMT systems.

## 1 Introduction

We had three research objectives in carrying out our experiments. The first objective was to assess the feasibility of fine-tuning an in-domain neural machine translation (NMT) baseline model using minimal unlabelled data. The second objective involved utilising large language models (LLMs) and prompt engineering techniques to post-edit translations within the same domain. The third objective was to examine the linguistic features of various models' erroneous translations, particularly in bilingual customer service conversations. For example, in their description of the data of the first edition of the Chat Task, (Farajian et al., 2020) noted the excessive use of pronouns in the dataset.

The remaining sections of the paper are organised as follows: section 2 mentions previous research, section 3 outlines our methods and describes our NMT systems, section 4 delves into our results[1], section 5 provides a discussion of these results, and section 6 outlines future work.

## 2 Previous Research

### 2.1 Fine-tuning Strategies for NMT and Domain Adaptation

Fine-tuning a pre-trained LLM baseline model with low-resource NMT has been the subject of previous MT empirical studies (Galiano-Jiménez et al., 2023) and the back-translation approach is often used to improve the accuracy of models (Hoang et al., 2018). Open source toolkits are available for building pipelines, such as `fairseq`[2]. However, some models require a higher level of expertise in pipeline construction and rely on cutting-edge hardware for optimal performance[3]. In terms of domain adaptation, filtering back-translations is considered one of the most frugal and efficient techniques (Kumari et al., 2021). In addition, more and more domain adaptations rely on prompt engineering.

Based on what was reported in the findings of the Chat Task 2022 (Farinha et al., 2022), MT systems handle source-related issues more or less similarly. Analysing the distribution of error types presented in the task indicates that "mistranslation" is the most frequent error across all systems. Furthermore, prompt-based machine translation has shown a significant impact in medical domains. For example, Ramachandran et al. (2023) demonstrated that using GPT-4 for extracting Social Determinants of Health (SDOH) from electronic health records achieved a 0.652 F1 score, which is comparable to the 7th best system among traditional supervised approaches.

---

[1] `https://github.com/lichaozhu/team_MULTITAN-GML_WMT24_Chat_Shared_Task`

[2] `https://github.com/facebookresearch/fairseq`

[3] For example, NLLB-200-3.3.B requires Hydra (Yadan, 2019) and very high GPU resources. We were unable to load and train the model using a dual A100 40GB setup due to persistent memory overflow problems.

## 2.2 Automatic Post-editing of MT and Prompt Engineering

Automatic post-editing (APE) systems are designed to enhance the quality of machine translation (MT) by *leveraging* data (Raunak et al., 2023; Gao et al., 2023). These systems work by taking both the source text and the initial MT output as inputs, then applying learned post-editing patterns to refine the translation, and the final output is an improved translation (Chollampatt et al., 2020; Sharma et al., 2021; Bhattacharyya et al., 2023). To further improve performance, APE systems often employ domain adaptation and fine-tuning on in-domain data (Moslem et al., 2023). Based on previous studies, prompting for machine translation still suffers from issues such as copying, mistranslation of entities, and hallucinations (Zhang et al., 2023). Furthermore, previous comprehensive evaluations of GPT models for machine translation across various language pairs indicate that GPT models perform competitively for high-resource languages, but face limitations with low-resource languages (Hendy et al., 2023; Jiao et al., 2023; Peng et al., 2023).

## 3 Methods and Tools

### 3.1 Fine-tuning via NMT Engine

For our primary submission, we have used a neural machine translation (NMT) engine, its in-domain baseline model, and in-domain training data to fine-tune the model. To create our fine-tuning dataset, we used the Chat Task 2022's valid and test sets (en ⇔ fr) as well as the Chat Task 2024's `train` and `valid` sets and compiled 13,622 aligned segments (122,905 words in English and 127,335 words in French). We used this dataset to fine-tune the *Dialog* in-domain model on the training server Model Studio Lite of Systran®[4] since we did not manage to fine-tune Facebook's NLLB-200-3.3B model, which was our first choice.

### 3.2 Translation and Post-editing with LLMs

For our two contrastive submissions, we have used NLLB-200-3.3B (NLLB Team et al.) baseline model and deep-translator[5] which was used by ChatGPT (GPT-4-turbo) to generate translations. All translations are then post-edited using prompt engineering via ChatGPT-4o.

---

## 4 Results

### 4.1 Qualitative Assessment

We have then compared three models in Systran Model Studio Lite to verify whether the in-domain Dialog model is adapted or not to the custom service conversation domain, by using the test set and reference translations published by the organisers of the Chat Task 2024. Table 1 compares the performance of three different models for language translation tasks: a fine-tuned model, an in-domain baseline model, and a generic baseline model. The performance is measured for two translation directions: English to French (en → fr) and French to English (fr → en).

|  | Fine-tuned model | In-domain baseline model | Generic baseline model |
|---|---|---|---|
| en → fr | **57.19** | 48.05 | 50.47 |
| fr → en | **55.02** | 48.28 | 48.19 |

Table 1: Comparison of generic baseline, in-domain baseline and fine-tuned models of Systran®

The fine-tuned model shows a significant improvement over both baseline models in both translation directions. This highlights the effectiveness of fine-tuning in enhancing model performance for specific tasks. The in-domain baseline model performs slightly worse than the generic baseline model for en → fr but slightly better for fr → en. This suggests that the in-domain data may not always provide a consistent advantage over generic data without further fine-tuning. The results indicate the importance of model fine-tuning in achieving superior translation quality and accuracy, especially in specialised domains. They seem to support our approach and the effectiveness of our fine-tuning dataset.

To compare translations, we used quantitative methods such as *vocabulary growth*, *characteristic elements computation*, and *correspondence analysis* (Lebart et al., 1997; Fleury and Zimina, 2014; Zimina-Poirot et al., 2020) implemented in *iTrameur*[6] and *Voyant Tools*[7]. In Figure 1, generated with *iTrameur*, the vocabulary growth curves of three predictions, fine-tuned Systran (*systran_ft*), NLLB-200-3.3B (*nllb*), and Deep translator (*deep-translator*) can be compared

---

Figure 1: Vocabulary growth curve of reference translation and predictions of fine-tuned Systran, NLLB-200-3.3B and Deep translator.



Figure 2: *Characteristic elements* computation for comparison of specific lexical features of reference translation and predictions of fine-tuned Systran (`systran_ft`), NLLB-200-3.3B (`nllb`) and Deep translator (`deep_translator`).

with the (*reference*) translation. While the reference translation is the longest (Nb occurrences: 22,834), it is followed by fine-tuned Systran (Nb occurrences: 22,291), which is the closest to the reference in terms of vocabulary growth.

In Figure 2 generated with *iTrameur*, we used *characteristic elements* computation to compare three predictions with the reference translation. The results show that many translation errors (including the occurrences of *E, S, t, Thank*, etc.) are over-represented in NLLB-200-3.3B prediction, while the reference translation and fine-tuned Systran prediction share common lexical features, such as identical translations *Are you still there?* ⇒ *Êtes-vous toujours là ?* attested by the over-representation of *Êtes*.

In Figure 3, we used correspondence analysis in *Voyant Tools* to compare our three predictions with the reference translation. The results suggest that the reference translation was carried out with human intervention, as it is clearly opposed



Figure 3: Correspondence analysis of the *reference* translation and tree predictions: fine-tuned Systran (*systran_ft*), NLLB-200-3.3B (*nllb*), and Deep translator *deep-translator*.

to three predictions (Zimina-Poirot et al. (2020) provides a discussion on this phenomenon). Although fined-tuned Systran is closer to `reference`, it is also very close to Deep translator, with NLLB-200-3.3B having a distinct profile.

Table 2 presents examples of segments that were incorrectly translated in our primary submission. It includes a comparison between the original source text, the reference translation, and our system's primary output, along with corresponding sentence-level BLEU and TER scores.

## 4.2 Comparisons of Primary and Contrastive Translations

In Table 3, we compared sentenceBLEU and TER scores of our *Primary* predicted by fine-tuned Systran model and two *Contrastives* predicted respectively by NLLB-200-3.3 baseline and Deep Translator. Except NLLB-200-3.3's predictions which have noticeably lower score, Deep Translator and fine-tuned Systran model have higher similar scores, which confirms our analysis of Figure 3. Deep Translator gets a slightly higher mean sentenceBLEU score, but its TER score is also higher. We noticed however that Deep Translator provided more literal or inaccurate translations of pragmatic expressions. It has translated *Bonjour* (greetings in French used in the daytime) by *Good morning*, and wrongly translated *You're welcome* by *Vous êtes les bienvenus*, which means "You are most welcome" in French.

Following the release of human evaluations, we have focused on mistranslations which scored 0 points, e.g. *I hope you have an excellent day* (source) is translated to *Merci pour l'information*

| | Source | Reference | Primary | sentenceBLEU | TER |
|---|---|---|---|---|---|
| 1 | Is there anything else I can assist you with to-day? | Avez-vous besoin d'aide pour autre chose aujourd'hui ? | Y a-t-il autre chose que je puisse faire pour vous aider aujourd'hui ? | 0.25 | 1.125 |
| 2 | I am so sorry to hear that. | Je regrette sincèrement d'apprendre cela. | Je suis vraiment désolé de l'apprendre. | 0.00 | 1.0 |
| 3 | You are welcome! | Avec plaisir ! | Je vous en prie. | 0.00 | 1.33 |
| 4 | You are welcome! | Ce fut un plaisir de vous parler. | C'était agréable de parler avec vous. | 0.00 | 1.0 |
| 5 | ok merci | Ok, thanks | Ok, thank you | 0.00 | 1.0 |

Table 2: Mistranslated segments in our primary submission

("Thank you for the information"). The presence of these translation segments probably reflects misalignments in the fine-tuning data, as Systran Model Studio Lite does not necessarily filter out mismatching segments during the training process. These segments of the translation memory can be deemed correct as part of the normalisation process.

## 5 Discussion

### 5.1 Automatic Post-editing vs. Prompt Engineering

Pipelines for translation and post-editing using LLM engines were proposed with LLM engines (Vidal et al., 2022). The primary submission and the two contrasting submissions were subsequently post-edited by ChatGPT-4o using instructions such as:

```
"Post-edit the translations in
file XX according to the source
texts in file YY where English
sentences are translated into
French, and French sentences
translated into English. Send me
back in one single file",
```

where two raw text files are given: XX is line-separated source file and YY translation file. We noticed that when we asked ChatGPT-4o to post-edit by performing domain adaptation considering our dataset as a reference or knowledge base, it did not work.

The default instructions are ineffective when used with Anthropic Claude. To detect the language accurately, it is necessary to use language columns. In this context, using tags enhances the precision of the translation (without them, the translation will default to a single language). Adhering to the token limit is crucial, as failure to do so may lead to overlooking the total number of tokens in the input. Although the tag has been modified to "tear", it still functions as the translated message.

Another hallucination occurred when the instructions themselves were translated. Figure 4 illustrates the interface and the applied prompt. The French text contained several misspellings, homophonic confusions, such as *est* versus *ait*, participle versus infinitive confusions, and various conjugation errors. We also attempted to prompt LLMs to translate from the initial CSV file, but this strategy has limitations. The LLMs may suggest Python code to extract sentences in both languages, translate only one language, or perform the task for a limited number of sentences.

Using Anthropic Claude for translation also highlights the variability in LLM translations. For example, for the sentence *Pardonnez-moi je n'ai pas du bien formuler ma question.* Three translations were obtained: 1) "I'm sorry I must not have formulated my question well." 2) "I apologize I must not have phrased my question well." 3) "I apologize, I may not have phrased my question well."

## 6 Further Research

### 6.1 Retrieval-Augmented Generation (RAG)

The database serves as a vital resource for addressing the challenges posed by rare or complex structures that may not be well-represented in translation models (Gao et al., 2024). Retrieval-augmented generation (RAG) is a technique for enhancing the accuracy and reliability of generative AI models with facts fetched from external sources. Future improvements could involve aug-

| Primary (Systran® fine-tuned) | | Contrastive 1 (NLLB-200-3.3) | | Contrastive 2 (Deep translator) | |
|---|---|---|---|---|---|
| sentenceBLEU | TER | sentenceBLEU | TER | sentenceBLEU | TER |
| 0.70 | 0.25 | 0.57 | 0.50 | 0.71 | 0.28 |

Table 3: Primary and Contrastives metrics comparison (arithmetic mean)



Figure 4: Anthropic Claude's interface with a prompt based on the URL of the WMT shared task test set

menting the training set with more examples, either through synthetic data or diverse real-world instances, to enhance the model's performance to translate challenging constructions, such as dislocations.

## 6.2 Explainability: Probing MT Systems for Trustworthy Outputs

Controlling LLM outputs and their repeatability is crucial for trustworthy AI. We tried to probe LLMs with (a) the detection of explicit representations and (b) their potential use in the LLM outputs. Similarly, in NMT, information might be available but not used by the system, as seen in the case of gender information discrepancies (Wisniewski et al. (2022a,b).

## 7 Conclusion

In this paper, we outline our methods for participating in the Chat Task 2024, focusing on enhancing translation quality in dialog-oriented machine translation systems through fine-tuning and prompt engineering. Our translation data files are available on GitHub[8]. Key findings indi-

cate that fine-tuning an in-domain NMT model is feasible with minimal unlabelled data, resulting in significant improvements in translation quality. The research also emphasises the importance of analysing linguistic features in translations to identify strengths and weaknesses of different machine translation models. The study also highlights the necessity of ensuring explainability in LLM outputs to foster trust in AI systems.

## Acknowledgements

---

[8] https://github.com/lichaozhu/team_MULTITAN-GML_WMT24_Chat_Shared_Task

[9] https://plateformes.u-paris.fr/category/plateformes/traitement-automatique

[10] https://u-paris.fr/eila/actualites-projet-multitan-gml

[11] https://u-paris.fr/plateforme-paptan

1020

# References

Pushpak Bhattacharyya, Rajen Chatterjee, Markus Freitag, Diptesh Kanojia, Matteo Negri, and Marco Turchi. 2023. Findings of the WMT 2023 shared task on automatic post-editing. In *Proceedings of the Eighth Conference on Machine Translation*, pages 672–681, Singapore. Association for Computational Linguistics.

Shamil Chollampatt, Raymond Hendy Susanto, Liling Tan, and Ewa Szymanska. 2020. Can automatic post-editing improve NMT? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2736–2746, Online. Association for Computational Linguistics.

M. Amin Farajian, António V. Lopes, André F. T. Martins, Sameen Maruf, and Gholamreza Haffari. 2020. Findings of the WMT 2020 shared task on chat translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 65–75, Online. Association for Computational Linguistics.

Ana C Farinha, M. Amin Farajian, Marianna Buchicchio, Patrick Fernandes, José G. C. de Souza, Helena Moniz, and André F. T. Martins. 2022. Findings of the WMT 2022 shared task on chat translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 724–743, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Serge Fleury and Maria Zimina. 2014. Trameur: A framework for annotated text corpora exploration. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: System Demonstrations*, pages 57–61, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.

Aarón Galiano-Jiménez, Felipe Sánchez-Martínez, Víctor M. Sánchez-Cartagena, and Juan Antonio Pérez-Ortiz. 2023. Exploiting large pre-trained models for low-resource neural machine translation. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 59–68, Tampere, Finland. European Association for Machine Translation.

Yuan Gao, Ruili Wang, and Feng Hou. 2023. How to design translation prompts for chatgpt: An empirical study. *arXiv preprint arXiv:2304.02182*.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Awadalla. 2023. How good are GPT models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Vu Cong Duy Hoang, Philipp Koehn, Gholamreza Haffari, and Trevor Cohn. 2018. Iterative back-translation for neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 18–24, Melbourne, Australia. Association for Computational Linguistics.

Wenxiang Jiao, Wenxuan Wang, Jen-Tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is ChatGPT a good translator? a preliminary study. *arXiv preprint arXiv:2301.08745*.

Surabhi Kumari, Nikhil Jaiswal, Mayur Patidar, Manasi Patwardhan, Shirish Karande, Puneet Agarwal, and Lovekesh Vig. 2021. Domain adaptation for NMT via filtered iterative back-translation. In *Proceedings of the Second Workshop on Domain Adaptation for NLP*, pages 263–271, Kyiv, Ukraine. Association for Computational Linguistics.

Ludovic. Lebart, André Salem, and Lisette Berry. 1997. *Exploring Textual Data*. Text, Speech and Language Technology. Springer Netherlands.

Yasmin Moslem, Gianfranco Romani, Mahdi Molaei, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. Domain terminology integration into machine translation: Leveraging large language models. *ArXiv*, abs/2310.14451.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. Towards making the most of ChatGPT for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5622–5633, Singapore. Association for Computational Linguistics.

Giridhar Kaushik Ramachandran, Yujuan Fu, Bin Han, Kevin Lybarger, Nic Dobbins, Ozlem Uzuner, and Meliha Yetisgen. 2023. Prompt-based extraction of social determinants of health using few-shot learning. In *Proceedings of the 5th Clinical Natural Language Processing Workshop*, pages 385–393, Toronto, Canada. Association for Computational Linguistics.

Vikas Raunak, Amr Sharaf, Yiren Wang, Hany Awadalla, and Arul Menezes. 2023. Leveraging GPT-4 for automatic translation post-editing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12009–12024, Singapore. Association for Computational Linguistics.

Abhishek Sharma, Prabhakar Gupta, and Anil Nelakanti. 2021. Adapting neural machine translation for automatic post-editing. In *EMNLP 2021 Sixth Conference on Machine Translation (WMT21)*, pages 315–319.

Blanca Vidal, Albert Llorens, and Juan Alonso. 2022. Automatic post-editing of MT output using large language models. In *Proceedings of the 15th Biennial Conference of the Association for Machine Translation in the Americas (Volume 2: Users and Providers Track and Government Track)*, pages 84–106, Orlando, USA. Association for Machine Translation in the Americas.

Guillaume Wisniewski, Lichao Zhu, Nicolas Ballier, and François Yvon. 2022a. Analyzing gender translation errors to identify information flows between the encoder and decoder of a NMT system. In *Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 153–163, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Guillaume Wisniewski, Lichao Zhu, Nicolas Ballier, and François Yvon. 2022b. Biais de genre dans un système de traduction automatique neuronale : une étude des mécanismes de transfert cross-langue [gender bias in a neural machine translation system: a study of crosslingual transfer mechanisms]. In *Traitement Automatique des Langues, Volume 63, Numéro 1 : Varia [Varia]*, pages 37–61, France. ATALA (Association pour le Traitement Automatique des Langues).

Omry Yadan. 2019. Hydra - a framework for elegantly configuring complex applications. *Github https://github.com/facebookresearch/hydra*.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: a case study. pages 41092–41110.

Maria Zimina-Poirot, Nicolas Ballier, and Jean-Baptiste Yunès. 2020. Approches quantitatives de l'analyse des prédictions en traduction automatique neuronale (TAN). In *JADT 2020 : 15èmes Journées Internationales d'Analyse statistique des Données Textuelles*, Toulouse, France. Université de Toulouse.

# The SETU-ADAPT Submissions to WMT 2024 Chat Translation Tasks

**Maria Zafar, Antonio Castaldo[a], Prashanth Nayak[b], Rejwanul Haque, Andy Way[c]**

South East Technological University, Carlow, Ireland
[a]University of Pisa, Tuscany, Italy
[b]KantanAI, Dublin, Ireland
[c]ADAPT Centre, Dublin City University, Dublin, Ireland
c00304029@setu.ie,antonio.castaldo@phd.unipi.it,pnayak@kantanai.io
rejwanul.haque@setu.ie,andy.way@adaptcentre.ie

## Abstract

This paper presents the SETU-ADAPT submissions to the WMT24 Chat Translation Task. Large language models (LLM) currently provides the state-of-the-art solutions in many natural language processing (NLP) problems including machine translation (MT). For the WMT24 Chat Translation Task we leveraged LLMs for their MT capabilities. In order to adapt the LLMs for a specific domain of interest, we explored different fine-tuning and prompting strategies. We also employed efficient data retrieval methods to curate the data used for fine-tuning. We carried out experiments for two language pairs: German-to-English and French-to-English. Our MT models were evaluated using three metrics: BLEU, chrF and COMET. In this paper we describes our experiments including training setups, results and findings.

## 1 Introduction

There have been drastic transformation in many areas of natural language processing (NLP) in recent times mainly due to the emergence of powerful LLMs. The LLM-based solutions are becoming more powerful and accurate than ever before. Notably, we have seen the unprecedented successes in many MT tasks in recent years, thanks to multilingual LLMs. In sum, the LLMs are the current state-of-the-art in MT research and development.

In our submission for the French-to-English and German-to-English Chat Translation Tasks, we built our MT systems using multilingual LLMs such as NLLB-200-600M (Team et al., 2022),[1]

Llama-3-8B (Dubey et al., 2024), [2] and mBART-50 (Tang et al., 2021).[3] We fine-tuned these models using both domain-specific and synthetic back-translated data.

Due to the lack of high-quality domain parallel data, we used a data generation approach. For this, we utilised a freely available monolingual data. We retrieved domain-specific monolingual sentences of target language and translated them back to source language for creating new synthetic data (Sennrich et al., 2016). This synthetic data was then combined with the original data for fine-tuning the LLMs. This approach ensured that our MT models are better adapted to the domain, thereby improving the quality of translations. We retrieved domain-specific monolingual German sentences from OPUS ELRC-4992 Customer Support MT [4] for creating our synthetic data. We also explored the idea of in-context learning by fine-tuning LLMs with a few-shot approach. These techniques helped our MT systems better adapt in translating agent input from the source language to the target language and customer response from the target language to the source.

The rest of the paper is organised as follows. Section 2 describes our related work. Our datasets are explained in Section 3 and Section 4 tells about the models and their fine-tuning. Section 5 discusses the experimental setup describing the parameters tuned in our systems. In Section 6, we discuss our results. Finally, Section 7 presents the conclusion of our work.

---

[1]NLLB-200: `https://ai.meta.com/research/no-language-left-behind/`

[2]Llama-3: `https://github.com/unslothai/unsloth/`
[3]mBART-50: `https://huggingface.co/facebook/mbart-large-50`
[4]OPUS: `https://opus.nlpl.eu/ELRC-4992-Customer_Support_MT/de&en/v1/ELRC-4992-Customer_Support_MT`

## 2   Related Works

In this section we discuss the papers that are related to our work. Alves et al. (2022) conducted experiments on fine-tuning `mBART-50` using domain-specific data retrieved through semantic search. For this, they used `LaBSE` (Language-Agnostic BERT Sentence Embedding) (Feng et al., 2020). They demonstrated how this approach leads them to large gains across all language pairs under evaluation. They also performed experiments using this data to further adapt the model using KNN-MT (Khandelwal et al., 2021). Note that this approach involves a nearest neighbor retrieval strategy, through which a set of relevant examples are provided at decoding time. They demonstrate how combining these approach leads to improved translation quality, over regular fine-tuning.

Liang et al. (2022a) used pre-trained LLMs and fine-tuned them to the domain of interest. For this, they first trained their models on general domain data and then fine-tuned them with chat translation training data. They used strategies such as including the multi-encoder framework, speaker tag prompt-based fine-tuning and boosted Self-COMET-based (Rei et al., 2020a) ensemble models to incorporate the potential context. They found their strategies helpful in improving the quality of translations produced by their MT models.

Yang et al. (2022) participated in the English-to-German task of the WMT22 Chat Translation Task. For this, they utilised the models previously submitted to the WMT21[5] news task (Wei et al., 2021) as their MT baseline systems. These baseline models are built upon a deep Transformer architecture (Vaswani et al., 2017). They used widely adopted optimisation strategies to improve model performance, including domain transfer, data selection, back-translation, self-training, noisy self-training, fine-tuning, and model averaging. Their results showed the effectiveness of their approached in improving the quality of translations.

Zhou et al. (2022) presented a multi-task multi-stage transitional (MMT) training framework, where they trained their model using the bilingual chat translation dataset and additional monolingual dialogues. To incorporate dialogue coherence and speaker characteristics in their model, they designed two auxiliary tasks: utterance discrimination and speaker discrimination. Their training had three stages: sentence-level pre-training on the large-scale parallel corpus, intermediate training with auxiliary tasks using additional monolingual dialogues and context-aware fine-tuning with a gradual transition. They found that the second stage served as a medium to reduce the training discrepancy between the pre-training and fine-tuning stages. They also trained their model using a gradual transition strategy, i.e. gradually transitioning from monolingual to bilingual dialogues, to make their stage transition smoother. Their results demonstrated the effectiveness of their framework, giving them better translations.

Liang et al. (2022b) contributed to the two large-scale in-domain paired bilingual dialogue corpora (28M for English-to-Chinese and English-to-German) through their framework. Their framework consisted of scheduled multi-task learning with three training stages, in which a gradient-based scheduling strategy was designed to take advantage of the auxiliary tasks for their model for higher translation quality. They conducted extensive experiments on four chat translation tasks, and their model achieved new state-of-the-art performance and outperformed the existing chat MT models by a significant margin.

## 3   Data Statistics

For our experiments we used the data provided by the WMT-24 Chat Translation Task[6] organisers. The dataset consists of authentic bilingual customer support conversations. This includes parallel data of interactions between an agent and a customer within the customer support domain. We detail the data description in Table 1. Note that we removed duplicates from the training data.

## 4   The LLMs

This section details the configurations of the LLMs that were used for our experiments.

### 4.1   mBART

mBART (Liu et al., 2020) is a pre-trained encoder-decoder Transformer model that was first trained on an auto-denoising task with monolingual data of twenty five languages. For adapting the mBART to the MT task, Tang et al. (2021) performed multilingual fine-tuning using data from fifty supported languages. For our experiments we used `facebook/mbart-large-50-many-to-many-mmt`

---

[5] https://www.statmt.org/wmt21/index.html

[6] https://www2.statmt.org/wmt24/chat-task.html

| Dataset | EN–to–DE | EN–to–FR |
|---|---|---|
| **WMT-24** | | |
| Train | 10,556 | 7,856 |
| Validation | 2,569 | 3,007 |
| Blind Test | 2,041 | 2,091 |
| + Back-translation | 1,317 | - |
| **WMT-22** | | |
| Train | 2,110 | 2,754 |
| Validation | - | - |
| **WMT-20** | | |
| Train | 10,248 | - |
| Validation | 1,619 | - |

Table 1: Overview of datasets.

checkpoint. We used the following hyperparameters setup for our experiments: `batch size: 4, number of training epochs: 5, predict_with_generate: True, evaluation strategy: epoch, logging steps: 2,000, and checkpoint save steps: 500`. The remaining parameters were set to default values.

### 4.2 NLLB

NLLB is a cutting-edge multilingual translation model developed to support many languages, mainly low-resource languages. Initially, the model was trained using diverse, multilingual data that includes various underrepresented languages. This comprehensive pre-training allows NLLB to effectively handle translation tasks across many languages that typically lack sufficient data. For our experiments we used `facebook/nllb-200-distilled-600M` checkpoint for building our MT systems. Our training configuration is as follows: `batch size: 4, 8; max sequence length: 128 tokens; training steps: 10,000, 20,000, 40,000; learning rate: 0.0001; optimiser: Adafactor; weight decay and gradient clipping applied; and model saved every 1000 steps`.

### 4.3 Llama

Llama is an auto-regressive language model that pretrained and fine-tuned in different sizes of data. We used `unsloth/llama-3-8b-bnb-4bit` checkpoint for building our MT systems. Our training parameters we set are as follows: `max seq length:`

`2048 tokens, batch size: 2 per device, gradient accumulation steps: 4, learning rate: 2e-4, mixed precision training enabled: (fp16 or bf16), learning rate scheduler: linear with 5-step warmup, maximum training steps: 500, optimiser: adamw-8bit, logging steps: 1, seed: 3407`.

## 5 Methodology

In this section, we discuss our methodologies.

### 5.1 mBART50

We used mBART for building three different MT systems for German-to-English. More specifically, mBART was fine-tuned on three distinct datasets: WMT-20,[7] WMT-22,[8] and WMT-24[9]. For French-to-English we used two datasets from WMT: WMT-22 and WMT-24. we detailed the datasets and hyperparameters setups in Section 3 and 4, respectively.

We performed data preprocessing the original data such as normalisation by removing special characters, removing duplicates and performing lowercase conversions. The source and target sentences were then tokenized using a predefined tokenizer.

In order to handle data during training and evaluation, a collator named as `DataCollatorForSeq2Seq`[10] is instantiated with the tokenizer and pretrained model checkpoint from Transformers library. This collator is designed to dynamically pad inputs to the maximum length within a batch, ensuring efficient processing. The `Seq2SeqTrainer` is then instantiated with the pretrained model checkpoint, training arguments, tokenized datasets, evaluation function, data collator, and tokenizer. This setup ensures a structured and efficient fine-tuning process, evaluating the model's performance at each epoch.

Fine-tuning is performed using the `Seq2SeqTrainer` class from the Transformers library. The training arguments are specified through `Seq2SeqTrainingArguments`, where parameters such as the output directory, batch sizes for training and evaluation and the number of training epochs were defined in Section 4.

---

[7] https://www.statmt.org/wmt20/chat-task.html
[8] https://wmt-chat-task.github.io
[9] https://www2.statmt.org/wmt24/chat-task.html
[10] https://huggingface.co/docs/transformers/en/main_classes/data_collator

## 5.2 NLLB

We built four MT systems for English-to-German and two MT systems for English-to-French considering NLLB as the baselines. We evaluated our MT systems on the development and blind test sets. The MT training setups are detailed below. Our first MT systems involves fine-tuning the baseline NLLB model on the original data. For our second MT system we used normalised data (i.e. removing special characters and duplicates, and lowercasing) for fine tuning to observe any impact of data cleaning on performance.

We build two additional MT systems for for English–German. We back-translated monolingual data in order to create a synthetic bilingual data. For this, we mined monolingual data from OPUS [11]. The domain of the monolingual data is customer support. We combined the generated synthetic data with the original data for fine-tuning. The first and second MT systems were fine-tuned on the combined data, and this gave us the third and fourth MT systems, respectively.

For training we handled out-of-memory errors by dynamically creating training batches, i.e. the Adafactor (Shazeer and Stern, 2018) optimizer is employed instead of AdamW (Loshchilov and Hutter, 2017) to save GPU memory. Weight decay and gradient clipping (Loshchilov and Hutter, 2017) were applied to stabilize the training. Training batches were created by randomly choosing the translation direction (source to target or reverse) and sampling sentence pairs. To enhance robustness against memory issues, a function was implemented to release memory, with parameters set to different batch-sizes, maximum sentence length, and different training-steps. For the German-to-English and French-to-English tasks the best performing models were found to be the ones with 40,000 and 10,000 training steps, respectively. The model is saved every 1,000 steps, allowing for interruptions to adjust parameters or evaluate translations. Training typically runs for a short period of time, which is sufficient for a language similar to those already known by NLLB.

Post-training evaluation involves testing translation quality using parameters like $num-beams = 4$, which affects accuracy, speed, and memory consumption, and parameters $a$ and $b$, which control

the maximum length of a generated text. The number of beams (or beam size) controls how many alternative sequences are kept during the search. This means that the model keeps the top 4 translations at each step during decoding.

## 5.3 Llama

We also used LLaMA for English-to-German and English-to-French. We built two MT systems for each of the translation tasks. This time, we focused on a specific learning technique, i.e. few-shot in-context learning. For this, we constructed a sentence retrieval system based on dense vector embeddings. Initially, sentence embeddings were generated using SentenceTransformer. More specifically, we used `all-MiniLM-L6-v2` [12] for our task. This model was applied to the source sentences of the dataset, transforming them into high-dimensional vector representations. These embeddings were then indexed using FAISS (Facebook AI Similarity Search) (Douze et al., 2024)[13], creating a searchable database of vectors. In other words, in order to create in-context learning examples, we encode the source test sentence using the pre-trained SentenceTransformer model. The resulting embedding is then used to query the FAISS index, which retrieves the most semantically similar sentences from the training dataset (note that we set $k = 3$). For constructing prompts, we retrieve the source sentences, their corresponding target sentences, and the associated language labels from the training dataset using the indices returned by FAISS. These sentences are then iteratively combined to construct three-shot prompts. Figure 1 shows the structure of a prompt. As can be seen from Figure 1, the prompt consists of initial instruction followed by three components: an instruction, an input, and a response. The instruction guides the task of translating from English-to-German/French or German/French-to-English. The language in an instruction is set dynamically based on labels provided with the sentence in our dataset. The input provides context, and the response is the desired output.

For the fine-tuning process tokenizer was instantiated using the `FastLanguageModel` class with parameters tailored to support efficient training on large sequences. The model is loaded from the `unsloth/llama-3-8b-bnb-4bit` pre-trained

---

[11] https://opus.nlpl.eu/ELRC-4992-Customer_Support_MT/de&en/v1/ELRC-4992-Customer_Support_MT

[12] https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

[13] https://github.com/facebookresearch/faiss

checkpoint. Subsequently, the model was further configured with $get-peft-model$, which applies Parameter-Efficient Fine-Tuning (PEFT) (Xu et al., 2023) techniques.

The fine tuning process was managed by `SFTTrainer`, which was integrated into `TrainingArguments` from the Transformers library. The training configurations were discussed in Section 4. Throughout the training, logging was performed at every step and the training process was seeded with a fixed value to ensure reproducibility. This training approach leveraged the state-of-the-art techniques to enable fine-tuning LLMs on extensive datasets while minimizing the computational overhead. The same parameters were tuned for both language-pairs.

The construction of a prompt for inference is identical to those constructed for training. The prompt instructs the model to translate a chat abstract from one language (source) to another (target). The instruction is specifically tailored to the languages involved in the translation, which are specified in the input-row. The source text is included in the prompt, while the output field is left blank, allowing the model to generate the translation. The prompt is tokenized using the pre-trained tokenizer, and the inputs are formatted as tensors compatible with PyTorch, with the computation offloaded to a GPU. The tokenized inputs are then passed to the model's $generate$ method, which performs the translation. The generated output is then decoded from the tokenized format back into a human-readable string using the batch-decode method.

# 6  Results

This section describes the results obtained. Table 2 shows the performance of our MT systems on the validation sets. Tables 3 and 4 show the results obtained on the blind test set provided by the task organisers.

As mentioned in Section 5, we normalised the original data by removing special characters and duplicate sentences and lowercasing to see the impact of data cleaning on performance. We see from Table 2 that the MT models fine-tuned on normalised data are better than those fine-tuned on original data for the German-to-English translation task. This clearly shows us the effectiveness of data normalisation. Our primary submission for German-to-English was based on fine-tuned Llama with

```
Below is an instruction that
    describes a task, paired with
    an input that provides further
     context. Write a response
    that appropriately completes
    the request.
Instruction:
Translate this chat from German
    to English:
Input:
German 1: <German sentence 1>
English 1: <English sentence 1>
German 2: <German sentence 2>
English 2: <English sentence 2>
German 3: <German sentence 3>
English 3: <English sentence 3>
German 4: <German sentence 4>
English 4:
Response:
<English sentence 4>
```

Figure 1: The structure of a generated prompt.

few-shot prompting. For this setup, we obtained 42.08 BLEU, 66.84 chrF and 85.25 COMET points on the validation set (cf. row 12 of Table 2). We also submitted two contrastive systems, i.e. NLLB and mBART50 fine-tuned on augmented data. The performance fine-tuned NLLB and mBART50 are shown in rows 9 and 3 of Table 2, respectively. Our constrastive submission 1 was based on NLLB. For this setup, we obtained 48.21 BLEU, 70.31 chrF and 84.60 COMET points on the validation set (cf. row 9 of Table 2). Our constrastive submission 2 was based on mBART50. For this setup, we obtained 47.73 BLEU, 69.17 chrF and 84.09 COMET points on the validation set (cf. row 3 of Table 2).

As for primary submission of the French-to-English task, we considered NLLB as our baseline and its performance is reported in row 19 of Table 2). We see from the table that this setup provided us 54.79 BLEU, 73.88 chrF, and 85.5 COMET points on the validation set. As in the German-to-English-task, we also submitted two contrastive systems for the French-to-English task. We fine-tuned Llama-3-8B and mBART50 following the few-shots prompt generation strategies described in Section 5. Our constrastive submission 1 was based on fine-tuned Llama-3-8B. For this setup, we obtained 38.23 BLEU, 66.54 chrF and 89.08 COMET points on the validation set (cf. row

| Model | BLEU | chrF | COMET |
|---|---|---|---|
| **German-English** | | | |
| mBART50 WMT20 | 53.27 | 72.43 | 86.042 |
| mBART50 WMT22 | 32.50 | 55.66 | 76.94 |
| mBART50 WMT24 | **47.73** | **69.17** | **84.096** |
| NLLB WMT24 SrcNorm $\rightarrow$ TgtNorm | 35.54 | 61.37 | 82.93 |
| NLLB WMT24 TgtNorm $\rightarrow$ SrcNorm | 46.08 | 68.75 | 83.73 |
| NLLB WMT24 Src $\rightarrow$ Tgt | 20.43 | 50.81 | 77.69 |
| NLLB WMT24 Tgt $\rightarrow$ Src | 16.10 | 51.44 | 76.16 |
| NLLB WMT24 + BT SrcNorm $\rightarrow$ TgtNorm | 39.62 | 64.02 | 84.30 |
| NLLB WMT24 + BT TgtNorm $\rightarrow$ SrcNorm | **48.21** | **70.31** | **84.60** |
| NLLB WMT24 + BT Src $\rightarrow$ Tgt | 23.32 | 53.10 | 79.05 |
| NLLB WMT24 + BT Tgt $\rightarrow$ Src | 17.33 | 52.50 | 77.10 |
| LLaMA WMT24 FS SrcNorm $\rightarrow$ TgtNorm | **42.08** | **66.84** | **85.25** |
| LLaMA WMT24 FS TgtNorm $\rightarrow$ SrcNorm | 20.05 | 52.90 | 83.03 |
| LLaMA WMT24 FS Src $\rightarrow$ Tgt | 20.07 | 57.60 | 85.71 |
| LLaMA WMT24 FS Tgt $\rightarrow$ Src | 35.54 | 59.79 | 87.66 |
| **French-English** | | | |
| mBART50 WMT22 | 43.51 | 64.64 | 80.27 |
| mBART50 WMT24 | **53.15** | **72.68** | **84.55** |
| NLLB WMT24 SrcNorm $\rightarrow$ TgtNorm | 46.24 | 68.78 | 85.42 |
| NLLB WMT24 TgtNorm $\rightarrow$ SrcNorm | **54.79** | **73.88** | **85.50** |
| NLLB WMT24 Src $\rightarrow$ Tgt | 31.45 | 59.69 | 80.65 |
| NLLB WMT24 Tgt $\rightarrow$ Src | 34.07 | 63.20 | 79.80 |
| LLaMA WMT24 FS SrcNorm $\rightarrow$ TgtNorm | 31.11 | 60.02 | 85.90 |
| LLaMA WMT24 FS TgtNorm $\rightarrow$ SrcNorm | 5.64 | 30.87 | 80.42 |
| LLaMA WMT24 FS Src $\rightarrow$ Tgt | **38.23** | **66.54** | **89.08** |
| LLaMA WMT24 FS Tgt $\rightarrow$ Src | 24.71 | 58.38 | 82.91 |

Table 2: Performance of our MT systems on the validation set. SrcNorm and TgtNorm stand for Source and Target normalised, respectively. BT stands for back-translation and FS stands for Few-Shot.

| Tag | Precision | Recall | F1 |
|---|---|---|---|
| **French-English** | | | |
| formality | 90.2 | 78.8 | 84.1 |
| lexical cohesion | 46.4 | 42.5 | 44.3 |
| pronouns | 90.8 | 72 | 80.3 |
| verb form | 62.9 | 56.8 | 59.7 |

Table 3: Official results for the French-English translation task (blind set).

24 of Table 2). Our constrastive submission 2 was based on mBART50. For this setup, we obtained 53.15 BLEU, 72.67 chrF and 84.55 COMET points on the validation set (cf. row 17 of Table 2).

Our primary submission of the German-to-English task was based on Llama. As can be seen from Table 4, we obtained 55.0 BLEU, 72.1 chrF, 90.8 COMET and 0.167827 CONTEXT-COMET-QE (Rei et al., 2020b) points on the WMT 2024 blind test sets. Our primary submission of French-to-English was based on NLLB. For this setup we obtained 31.3 BLEU, 60.9 chrF, 82.4 COMET and -0.23095 CONTEXT-COMET-QE points. Our best-performing system of the German-to-English task is Llama with few-shot learning. We secured the top place for German-to-English in this competition. Table 3 shows our precision for pronouns in the French-to-English system. Our submission for the French-to-English translation task is in fact the best-performing system in terms of pronoun translatiosn.

| Model | COMET | ChrF | BLEU | COMET-QE |
|---|---|---|---|---|
| **German-English** | | | | |
| LlaMa WMT24 FS | 90.8 | 72.1 | 55.0 | 0.16 |
| **French-English** | | | | |
| NLLB WMT24 | 82.4 | 60.9 | 31.3 | -0.23 |

Table 4: Official results. Performance of our MT systems on the blind set (primary submissions).

## 7 Conclusion

This paper described our submissions to the WMT 2024 Chat Translation Task for German-to-English and French-to-English language pairs. We applied several training and fine-tuning strategies such as standard fine-tuning and fine-tuning with few-shot prompting. We investigated our approaches using three different LLMs: NLLB, Llama and mBART. This allowed us to make a comparative analysis between different architectures and strategies. One of the key findings of our investigation is that the performance of the MT systems on translating conversational messages can be improved with knowledge transfer. We also found that our MT systems exhibit robustness on this *difficult-to-translate* domain.

For future investigations, given the shortage of conversational data, we plan to focus on exploring the use of advanced data augmentation techniques. We also intend to further investigate to what extent synthetic data can be beneficial in chat translation scenarios.

## References

João Alves, Pedro Henrique Martins, José G. C. de Souza, M. Amin Farajian, and André F. T. Martins. 2022. Unbabel-IST at the WMT chat translation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 943–948, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The Faiss library. ArXiv:2401.08281 [cs].

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2020. Language-agnostic bert sentence embedding. *arXiv preprint arXiv:2007.01852*.

Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2021. Nearest neighbor machine translation.

Yunlong Liang, Fandong Meng, Jinan Xu, Yufeng Chen, and Jie Zhou. 2022a. BJTU-WeChat's systems for the WMT22 chat translation task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 955–961, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Yunlong Liang, Fandong Meng, Jinan Xu, Yufeng Chen, and Jie Zhou. 2022b. Scheduled multi-task learning for neural chat translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4375–4388, Dublin, Ireland. Association for Computational Linguistics.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual Denoising Pre-training for Neural Machine Translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Ilya Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020a. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020b. Unbabel's participation in the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 911–920, Online. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data.

Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. *CoRR*, abs/1804.04235.

Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2021. Multilingual Translation from Denoising Pre-Training. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3450–3466, Online. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Daimeng Wei, Zongyao Li, Zhanglin Wu, Zhengzhe Yu, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Minghan Wang, Lizhi Lei, Min Zhang, Hao Yang, and Ying Qin. 2021. HW-TSC's participation in the WMT 2021 news translation shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 225–231, Online. Association for Computational Linguistics.

Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. 2023. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment.

Jinlong Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, Xiaoyu Chen, Zhengzhe Yu, Zhiqiang Rao, Shaojun Li, Zhanglin Wu, Yuhao Xie, Yuanchang Luo, Ting Zhu, Yanqing Zhao, Lizhi Lei, Hao Yang, and Ying Qin. 2022. HW-TSC translation systems for the WMT22 chat translation task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 962–968, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Chulun Zhou, Yunlong Liang, Fandong Meng, Jie Zhou, Jinan Xu, Hongji Wang, Min Zhang, and Jinsong Su. 2022. A multi-task multi-stage transitional training framework for neural chat translation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(7):7970–7985.

# Exploring the traditional NMT model and Large Language Model for chat translation

**Jinlong Yang, Hengchao Shang, Daimeng Wei, Jiaxin Guo,**
**Zongyao Li, Zhanglin Wu, Zhiqiang Rao, Shaojun Li,**
**Yuhao Xie, Yuanchang Luo, Jiawei Zheng, Bin Wei, Hao Yang**
Huawei Translation Service Center, Beijing, China
{yangjinlong7,shanghengchao,weidaimeng,guojiaxin1,lizongyao,
wuzhanglin2,raozhiqiang,lishaojun18,xieyuhao2,luoyuanchang,
zhengjiawei15,weibin29,yanghao30}@huawei.com

## Abstract

This paper describes the submissions of Huawei Translation Services Center(HW-TSC) to WMT24 chat translation shared task on English↔Germany (en-de) bidirection. The experiments involved fine-tuning models using chat data and exploring various strategies, including Minimum Bayesian Risk (MBR) decoding and self-training. The results show significant performance improvements in certain directions, with the MBR self-training method achieving the best results. The Large Language Model also discusses the challenges and potential avenues for further research in the field of chat translation.

## 1 Introduction

Neural machine translation (NMT) (Sutskever et al., 2014; Bahdanau et al., 2015; Gehring et al., 2017; Wu et al., 2023) has made substantial progress in recent years, largely due to the adoption of the transformer (Vaswani et al., 2017) architecture. NMT has demonstrated promising translation results across various scenarios. However, research in the field of chat translation remains limited, primarily due to the scarcity of chat data. In prior chat-related tasks, we utilized data from related domains, such as spoken dialogue and subtitles, to augment our translation models, but the outcomes were only mediocre.

Like the preceding two chat shared tasks, the WMT24 chat shard task concentrates on translating conversations between consumers and servers in different languages. We participated in the en-de bidirectional translation task. The en-de bidirectional models we submitted to the WMT22 chat task (Yang et al., 2022) function as our baseline models, leveraging the deep transformer (Dou et al., 2018) architecture. Building on this foundation, we employed the Minimum Bayesian Risk (MBR) strategy to select the optimal translation outcomes,

and iterative self-training yielded the best results on the development set.

Beyond traditional NMT models, the emergence of large language model(LLM) has introduced a new paradigm to translation tasks(Wang et al.; Moslem et al., 2023; Guo et al., 2024). Due to its extensive context length and powerful language modeling capabilities, large language models significantly outperform NMT in the translation of lengthy texts and the fluency of translation results. We input the translation output from the NMT model into the LLM as a prompt, allowing the LLM to combine the reference translatio from traditional NMT model to produce an improved translation. However, the comet metric of the LLM's output did not surpass the optimal results of the NMT model.

Recognizing that chat translation is a context-aware task, we conducted a series of context-aware experiments(Wu et al., 2024) using LLMs with WMT and IWSLT document data . We fine-tuned the LLM by constructing streamed translations and contextualized translation data, and translated the development set in the same format. Unfortunately, the results were unsatisfactory.

The structure of this paper is as follows: Section 2 describes our data volume and format for fine-tuning the LLM. The model structure and key methods utilized are presented in Section 3. Section 4 outlines our experiment setting. Results and analysis are presented in Section 5, and we conclude our work in Section 6.

## 2 Data

### 2.1 Data Size

All experiments conducted for this task are based on the model developed by our team, as participated in the WMT22 chat shared task. For details on the training data and strategies used for this model, please refer to the system report Yang et al. (2022); Wei et al. (2021). Table 1 and Table 2 list all

| 24 train | 24 valid | 22 valid | 22 test |
|---|---|---|---|
| 17805 | 2569 | 2109 | 2488 |

Table 1: Chat shared task en-de bilingual data lines used for training

| Dataset | lines | documents |
|---|---|---|
| iwslt_2017_ted | 209522 | 1705 |
| news-commentary-v18 | 449333 | 11396 |

Table 2: Document-level data used for LLM related experiments

the data used in this experiment. Based on the prior tasks experience, the contribution of out-domain data to the improvement of translation quality is limited. Therefore, we only further optimize our translation model using the data shown in Table 1, which consists of historical chat tasks. The data in Table 2 is used for fine-tuning the LLM, enabling it to translate context-aware texts and validate the impact of paragraph information on dialogue translation quality.

## 2.2 Data pre-processing

Since the domain-specific data listed in Table 1 is limited, no special treatment was applied to this portion of the data; it was simply tokenized and input into the NMT model. For the document data in Table 2, we constructed the two formats shown in Table 3 by considering the characteristics of chat tasks, and used them to fine-tune the LLM, separately validating the impact of only preceding information and both preceding and context information on chat translation quality.

In the format of streamlined translation, during each translation session, only preceding information is visible. The LLM generates results based on this preceding information and the previews translation output, resulting in a translation that leans more towards the style of the reference.

In the context-aware translation format, during each translation session, preceding and following N sentences are provided along with the output of the NMT model, guiding the LLM to combine context information to produce a more natural translation.

## 3 System Overview

### 3.1 Model

The baseline models for WMT24 chat task use the Transformer-Big architecture. Deep transformer is an improvement of Transformer, which increases the number of encoder layers and uses pre-layer-normalization to further improve model performance. Therefore, in this task, we adopt the following model architecture:

- Deep 25-6 large Model: This model features 25-layer encoder, 6-layer decoder, 1024 dimensions of word vector, 4096 domensions of FFN, 16-head self-attention, and pre-layer-normalization.

For experiments related to large language model, we choose llama2-8b as the base.

### 3.2 MBR Decoding

Minimum Bayesian Risk (MBR) decoding was initially introduced during the era of statistical machine translation(Kumar and Byrne, 2004; Jinnai et al., 2024). This strategy calculates the output with the minimum expected error among multiple candidates, rather than simply selecting the result with the highest probability during the decoding process. In our experimental approach, we utilize the outputs of 10 distinct models as candidates. These candidates are then used to score each other's comet, and the candidate with the highest average comet is chosen as the final output. Algorithm 1 show the detail.

### 3.3 Regularized Dropout

Regularized Dropout (R-Drop) [1](Liang et al., 2021) presents a simple yet more effective approach to regulate the training inconsistency caused by dropout (Srivastava et al., 2014). Specifically, during each mini-batch training, each data sample is processed twice through the forward pass, with each pass utilizing a distinct sub-model and randomly dropping out some hidden units. R-Drop minimizes the bidirectional Kullback-Leibler (KL) divergence (van Erven and Harremos, 2014) between the two distributions outputted by the two sub-models for the same data sample, thereby regulating the outputs of two sub-models randomly sampled from dropout for each data sample in training. This method effectively alleviates the inconsistency between the training and inference stages.

### 3.4 Self-Training

Self-Training(ST) (Imamura and Sumita, 2018), also known as forward translation (FT) (Wu et al., 2019), typically involves utilizing a forward NMT

---

[1]https://github.com/dropreg/R-Drop

| Streaming Translation Data Format |
| --- |
| Natural English: \<src1\>, Translated German: \<mt1\>, Natural German:\<ref1\> |
| Natural German: \<src2\>, Translated English: \<mt2\>, Natural English:\<ref2\> |
| Natural English: \<src3\>, Translated German: \<mt3\>, Natural German:\<ref3\> |
| Translate the following sentence into German with a style bias towards Natural: |
| Natural English: \<src4\>, Translated German: \<mt4\>, Natural German: \<ref4\> |
| **Context-aware Translation Data Format** |
| Natural English: \<src1\>, Translated German: \<mt1\> |
| Natural German: \<src2\>, Translated English: \<mt2\> |
| Natural English: \<src3\>, Translated German: \<mt3\> |
| Natural German: \<src4\>, Translated English: \<mt4\> |
| Natural English: \<src5\>, Translated German: \<mt5\> |
| Translate the following sentence into German with a style bias towards Natural: |
| Natural English: \<src3\>, Natural German: \<ref3\> |

Table 3: LLM Supervised fine-tuning(SFT) data format

---

**Algorithm 1** MBR decoding algorithm

**Input:**

The set of translation candidates file, $MT_n$;

The source text file, $SRC$;

Comet metric model, $M_{comet}$;

**Output:** final translation output

1: initialize output list $out[]$
2: **for** each $line \in [MT_1, ..., MT_n, SRC]$ **do**
3:     initialize $tmp\_max\_comet = 0$
4:     initialize $candidate\_mt = ''$
5:     **for** each $candidate \in [mt_1, mt_2, ..., mt_n]$ **do**
6:         let each $mt_x$ as ref, $candidate$ as mt and calculate the comet score with source text using $M_{comet}$
7:         $mean\_comet = \frac{\sum_{x=1}^{n} comet_x}{n}$
8:         **if** $mean\_comet > tmp\_max\_comet$ **then**
9:             $tmp\_max\_comet = mean\_comet$
10:             $candidate\_mt = candidate$
11:         **end if**
12:     **end for**
13:     out.append($candidate\_mt$)
14: **end for**
15: **return** out

---

model to translate source-side monolingual data into target-side text, thereby generating synthetic bilingual data. The generated data is then employed to train the forward translation model. Typically, beam search (Freitag and Al-Onaizan, 2017) is applied for forward translation. In our experimental approach, we set the beam size to 4. Furthermore, we utilized the MBR selection results as self-training data, which led to the best results on the validation set.

### 3.5 Back Translation

Back-translation (Edunov et al., 2018; Wei et al., 2023) is acknowledged as a highly effective data augmentation strategy to boost NMT model performance. Unlike forward translation, back-translation converts target-side monolinguals into source-side text, thereby producing synthetic parallel corpora. Numerous back-translation techniques have been explored, with sampling (Graça et al., 2019), noise (Edunov et al., 2018), and tagged back-translation (Caswell et al.) demonstrating superior results. In our experimental setup, we opted for sampling back-translation.

### 3.6 Model Averaging

Model averaging (Dormann et al., 2018) is a widely utilized technique to enhance translation quality. Typically, models (in our experiment, 5 models) that exhibit the highest performance on the development set are chosen for parameter averaging, which leads to substantial improvements.

### 3.7 LLM Few-shot Prompting

Although large language models exhibit impressive zero-shot capabilities, they still struggle with more complex tasks in the zero-shot setting. To address this, few-shot prompting can be employed as a technique for in-context learning, where demonstrations are provided in the prompt to guide the model towards enhanced performance. In our approach, we provide 5 reference translations to assist the large language model in producing superior results.

### 3.8 LLM SFT with LoRA

LLM SFT (Supervised Fine-Tuning) is a technique for fine-tuning large language models using specific datasets, which effectively enhances the performance of large language models on tasks such as text generation, machine translation, or sentiment analysis. LoRA (Low-Rank Adaptation)(Hu et al., 2022) is a technique that reduces the computational burden during large language model training by decreasing the number of model parameters through matrix decomposition. This technique maintains performance while lowering computational and memory requirements. By applying LoRA, large language models can perform better under limited computational resources, reducing training costs and resource consumption.

## 4 Experiment Setting

During the NMT model training phase, we use Pytorch-based Fairseq[2] (Ott et al., 2019) open-source framework as our benchmark system. Each model is trained using 8 GPUs with a batch size of 2048. The update frequency is 4 and the learning rate is 5e-4. The label smoothing rate is set to 0.1, the warm-up steps to 4000, and the dropout to 0.3. Adam optimizer (Kingma and Ba, 2015) with $\beta1$=0.9 and $\beta2$=0.98 is also used. Beyond that, we have configured the hyper parameter reg-alpha of the R-Drop technique to a value of 5. In the evaluation phase, We employ the official automatic evaluation scripts and primarily base our model and result selection on the comet metric(Rei et al., 2022)[3].

In the experiments related to large models, we utilize the open-source model llama2_8b_instruct from Meta and the training scripts from HF to train our models, setting the max_seq_length to 1024.

For inference on large models, we employ the vllm tool.

## 5 Result and Analysis

Table 4 displays the results of the official test set, ranked according to the comet-22 score, where our system achieved the top position in comet-22, chrF, and BLEU metrics.

The primary results we submit are obtained by translating the source text of the test set with multiple NMT models, selecting the optimal output using MBR strategy, then training on the best models from the validation set using self-training method. The models are averaged over 5 epochs before being used to translate the test set to yield the final results.

### 5.1 Sentence-level NMT

In the previous chat tasks, we have tried various strategies to optimize the model, and the results from the validation set indicate that the baseline model from 2022 was already sufficiently powerful. On this basis, we combined this year's training set, the 2022 validation and test sets, and conducted BT and ST reinforcement strategies, only in the direction of translation from English to German has there been a noticeable improvement. The results shown in Table 5.

To further improve the results, we attempted the MBR decoding strategy, generating 10 alternative outputs for the validation set using different NMT models in previous steps. These outputs were scored using comet, and the output with the lowest Bayesian risk was selected as the final result. The results in Table 5 indicate that improvement was only seen in the en→de direction. Further, we utilized the MBR results to perform another ST on each direction, ultimately achieving the best results in both directions in the validation set. The reason for the improvement we observed is that the MBR algorithm can integrate the capabilities of multiple models. When performing self training, it essentially utilizes the optimal results of multiple models for a round of knowledge distillation.

### 5.2 Document-level MT with LLM

According to the test results shown in Table 6, on the chat task valid set, the results of LLM (Large Language Model) are significantly worse than sentence-level under both comet or doc-comet metrics. The few-shot capabilities of LLM is in-

| team | comet↑ | chrf↑ | bleu↑ | context-comet-qe↑ |
|------|--------|-------|-------|-------------------|
| HW-TSC | **93.4** | **83.2** | **69.8** | 0.221 |
| unbabel+it | 92.9 | 78.2 | 62 | **0.253** |
| clteam | 91.3 | 71.9 | 53 | 0.204 |
| ADAPT | 90.8 | 72.1 | 55 | 0.168 |
| DCUGenNLP | 90.8 | 71.2 | 53 | 0.188 |
| baseline | 89.8 | 70.8 | 51.1 | 0.173 |
| SheffieldGate | 89.4 | 67.5 | 45.2 | 0.177 |

Table 4: The official automatic evaluation results of the test set, ranked based on the COMET-22 score

| System | en→de | de→en |
|--------|-------|-------|
| baseline | 86.76 | 85.88 |
| 22_denoise | 90.06 | 91.42 |
| + ST | 91.23 | 91.40 |
| + ST&BT | 91.23 | 91.53 |
| + MBR ST | **91.91** | **91.86** |
| MBR | 91.75 | 90.87 |

Table 5: Sentence-level NMT results.

deed far better than zero-shot, but it still falls short of sentence-level results. After using the document-level data for LLM SFT, the results became even worse. We analyzed that the reason is the large domain shift, as the IWSLT and WMT datasets we used are far from the domain of the chat task.

To validate the capability of LLM in translating document-level content, we tested the results on the iwslt2017 en-de document-level test set. The results in the right half of Table 6 demonstrate that LLM's few-shot capability surpassed that of the chat task's sentence-level model on this test set. Further, by fine-tuning the large model with document-level data, we obtained better results.

Comparing the results of stream translation and context-aware translation, we originally expected context-aware format data to yield better results because the model could refer to contextual information during translation. However, we analyzed that stream translation sees the previous step's translation result each time, which is more consistent with the translation style of large model. On the contrary, context-aware requires input of the reference MT result from sentence-level model in one go, which is less consistent with the style of large model, causing the model to fail to effectively utilize these information.

## 6 Conclusion

This paper presents the submissions of HW-TSC to the WMT 2024 Chat Translation Shared Task. For both direction in en↔de translation task, we perform experiments with a series of training strategies. The results show that MBR self-training achieves the best results. In the future, we will continue to explore the applicability of MBR strategy mentioned in this paper.

Beyond that, due to time constraints, further fine-tuning of large language models using chat task data was not conducted to assess its performance. Additionally, there is room for continued exploration of the translation capabilities of large language models.

## References

Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.

Isaac Caswell, Ciprian Chelba, and David Grangier. Tagged back-translation. *WMT 2019*, page 53.

Carsten F. Dormann, Justin M. Calabrese, Gurutzeta Guillera-Arroita, Eleni Matechou, Volker Bahn, Kamil Bartoń, Colin M. Beale, Simone Ciuti, Jane Elith, Katharina Gerstner, Jérôme Guelat, Petr Keil, José J. Lahoz-Monfort, Laura J. Pollock, Björn Reineking, David R. Roberts, Boris Schröder, Wilfried Thuiller, David I. Warton, Brendan A. Wintle, Simon N. Wood, Rafael O. Wüest, and Florian Hartig. 2018. Model averaging in ecology: a review of bayesian, information-theoretic, and tactical approaches for predictive inference. *Ecological Monographs*, 88(4):485–504.

Zi-Yi Dou, Zhaopeng Tu, Xing Wang, Shuming Shi, and Tong Zhang. 2018. Exploiting deep representations for neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4253–4262.

| System | chat en→de | | chat de→en | | iwslt en→de | | iwslt de→en | |
|---|---|---|---|---|---|---|---|---|
| | comet | d-comet | comet | d-comet | comet | d-comet | comet | d-comet |
| Baseline | 86.76 | 79.40 | 85.88 | 79.77 | - | - | - | - |
| MBR ST | **91.91** | **85.41** | **91.86** | **86.21** | 84.70 | 77.55 | 87.05 | 80.81 |
| llama2_8b_instruce | 87.56 | 79.99 | 86.96 | 80.89 | 82.53 | 75.07 | 86.21 | 79.74 |
| + 5 best | 90.05 | 83.34 | 88.72 | 83.11 | 85.10 | 77.98 | 87.20 | 81.04 |
| stream | 85.47 | 78.50 | 83.98 | 78.80 | **85.69** | **78.91** | **87.45** | **81.73** |
| context-aware | 81.82 | 73.81 | 83.89 | 77.37 | 84.80 | 77.51 | 86.65 | 80.43 |

Table 6: The results of LLM MT

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 489–500.

Markus Freitag and Yaser Al-Onaizan. 2017. Beam search strategies for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 56–60.

Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. 2017. Convolutional sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1243–1252.

Miguel Graça, Yunsu Kim, Julian Schamper, Shahram Khadivi, and Hermann Ney. 2019. Generalizing back-translation in neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 45–52.

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Kenji Imamura and Eiichiro Sumita. 2018. Nict self-training approach to neural machine translation at nmt-2018. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 110–115.

Yuu Jinnai, Tetsuro Morimura, Ukyo Honda, Kaito Ariu, and Kenshi Abe. 2024. Model-based minimum Bayes risk decoding for text generation. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 22326–22347. PMLR.

Diederik P Kingma and Jimmy Lei Ba. 2015. Adam: A method for stochastic optimization.

Shankar Kumar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 169–176, Boston, Massachusetts, USA. Association for Computational Linguistics.

Xiaobo Liang, Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, and Tie-Yan Liu. 2021. R-Drop: Regularized Dropout for Neural Networks. *arXiv e-prints*, page arXiv:2106.14448.

Yasmin Moslem, Rejwanul Haque, and Andy Way. 2023. Adaptive machine translation with large language models.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems-Volume 2*, pages 3104–3112.

Tim van Erven and Peter Harremos. 2014. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. Document-level machine translation with large language models.

Daimeng Wei, Zongyao Li, Zhanglin Wu, Zhengzhe Yu, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Minghan Wang, Lizhi Lei, Min Zhang, Hao Yang, and Ying Qin. 2021. HW-TSC's participation in the WMT 2021 news translation shared task. In *Proceedings of the Sixth Conference on Machine Translation*, pages 225–231, Online. Association for Computational Linguistics.

Daimeng Wei, Zhanglin Wu, Hengchao Shang, Zongyao Li, Minghan Wang, Jiaxin Guo, Xiaoyu Chen, Zhengzhe Yu, and Hao Yang. 2023. Text style transfer back-translation.

Lijun Wu, Yiren Wang, Yingce Xia, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2019. Exploiting monolingual data at scale for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4207–4216.

Minghao Wu, Thuy-Trang Vu, Lizhen Qu, George Foster, and Gholamreza Haffari. 2024. Adapting large language models for document-level machine translation. *arXiv preprint arXiv:2401.06468*.

Zhanglin Wu, Zongyao Li, Daimeng Wei, Hengchao Shang, Jiaxin Guo, Xiaoyu Chen, Zhiqiang Rao, Zhengzhe Yu, Jinlong Yang, Shaojun Li, Yuhao Xie, Bin Wei, Jiawei Zheng, Ming Zhu, Lizhi Lei, Hao Yang, and Yanfei Jiang. 2023. Improving neural machine translation formality control with domain adaptation and reranking-based transductive learning. In *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*, pages 180–186, Toronto, Canada (in-person and online). Association for Computational Linguistics.

Jinlong Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, Xiaoyu Chen, Zhengzhe Yu, Zhiqiang Rao, Shaojun Li, Zhanglin Wu, Yuhao Xie, Yuanchang Luo, Ting Zhu, Yanqing Zhao, Lizhi Lei, Hao Yang, and Ying Qin. 2022. HW-TSC translation systems for the WMT22 chat translation task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 962–968, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

# Graph Representations for Machine Translation in Dialogue Settings

**Lea Krause[1]\*, Selene Baez Santamaria[1], Jan-Christoph Kalo[2]**

[1]Vrije Universiteit Amsterdam, [2]Universiteit van Amsterdam
{l.krause, s.baezsantamaria}@vu.nl, j.c.kalo@uva.nl

## Abstract

In this paper, we present our approach to the WMT24 - Chat Task, addressing the challenge of translating chat conversations. Chat conversations are characterised by their informal, ungrammatical nature and strong reliance on context posing significant challenges for machine translation systems. To address these challenges, we augment large language models with explicit memory mechanisms designed to enhance coherence and consistency across dialogues. Specifically, we employ graph representations to capture and utilise dialogue context, leveraging concept connectivity as a compressed memory. Our approach ranked second place for Dutch and French, and third place for Portuguese and German, based on COMET-22 scores and human evaluation.

## 1 Introduction

Machine translation (MT) has been a prominent area of research, leading to the development of various approaches over the years (Maruf et al., 2021). While significant progress has been made, the majority of research has concentrated on refining methodologies rather than exploring the different types of text that require translation. A notable gap exists in the automatic translation of chat conversations—a gap that the WMT24 - Chat task specifically aims to address.

Chat conversations present unique challenges due to their informal, spontaneous nature, and frequent grammatical inconsistencies (Gonçalves et al., 2022). These characteristics starkly contrast with the more structured and formal text types, such as news articles, technical manuals, and political or medical documents, which have been the traditional focus of MT systems. In the context of chat translation, it is crucial to incorporate dialogue context effectively and to model the speakers and their language direction.



Figure 1: Approach 1: Triple-TowerInstruct

Recent advances in machine translation have increasingly leveraged large language models (LLMs). However, as noted by Maharana et al. (2024), LLMs often struggle with tasks requiring long-term memory, reasoning over historical context, and establishing long-range temporal or causal connections. These limitations are particularly problematic to dialogue tasks, where maintaining coherence and consistency across a conversation is vital.

To address these challenges, our system proposes enhancing LLMs with explicit memory mechanisms designed to support the generation of more consistent and coherent translations in dialogue settings[1]. We hypothesise that utilising graph representations will further improve the translation of chat conversations by capturing the connectivity between concepts, thus serving as a compressed memory of the dialogue context.

## 2 Related Work

In this section, we provide a brief overview of related work in the areas of conversational NLP, machine translation of conversational text, and text generation methods that incorporate knowledge graphs as an additional source of information.

---

\* Corresponding author.

[1]All code and data related available at https://github.com/selBaez/chat-task-2024-data.

**Conversational NLP** Dialogue systems have a long-standing history in NLP. The advent of LLMs has led to significant improvements in the quality of these systems. However, a persistent challenge has been the limited context window of LMs, which restricts their ability to manage long chat histories effectively (Xu et al., 2021). To address this, retrieval-augmented models have been developed, which retrieve relevant passages from prior interactions to maintain coherence in dialogue over extended conversations (Xu et al., 2021). Recently, advancements in model architecture have resulted in substantially larger context windows, enabling state-of-the-art dialogue systems, such as ChatGPT, to operate effectively with this extensive LMs (Achiam et al., 2023).

**Machine Translation** Machine translation has seen remarkable advancements with the rise of large language models (Wang et al., 2023; Robinson et al., 2023). However, translating dialogues remains a particularly challenging task due to the informal and often context-dependent nature of conversational text (Gonçalves et al., 2022). The findings of recently shared tasks highlight ongoing difficulties and emerging solutions in this area (Farinha et al., 2022).

Our work is particularly related to the use of knowledge graphs in translation tasks (Moussallem et al., 2018; Zhao et al., 2021). In most existing approaches, multilingual knowledge graphs are leveraged to disambiguate and translate key entities within the text. This approach differs significantly from our method, as we employ a monolingual graph to store key information from the dialogue in a compressed format, facilitating more accurate and context-aware translations.

**Graph-based Dialogue Systems** Knowledge graphs have proven to be a valuable resource for grounding dialogue systems. The most common approach involves integrating large, external knowledge graphs to provide additional context and information that can enhance the dialogue's quality and relevance (Liu et al., 2019; Tuan et al., 2019; Zhang et al., 2020). While these approaches share a similar objective with our work, they fundamentally differ in that the knowledge graphs used are independent of the dialogue content itself.

In contrast, other approaches leverage graphs to represent the dialogue history, offering a structured way to maintain and utilise past interactions (Xu et al., 2020; Chen et al., 2023). This method enhances transparency, reduces the likelihood of hallucinations, and improves the system's ability to manage long-term conversations (Baez Santamaria et al., 2023). Our work aligns with this approach by utilising a graph to capture and organise key dialogue information, enabling more effective and contextually grounded dialogue systems.

## 3 Shared Task description

A dataset of original bilingual customer support conversations is provided. The language pairs available are English ⇌ German (en-de), English ⇌ Dutch (en-nl), English ⇌ French (en-fr), English ⇌ Brazilian Portuguese (en-pt_br), and English ⇌ Korean (en-ko). Due to our team's language expertise, we decided to focus on the first four pairs.

## 4 System Overview

All our systems work with graphs extracted from dialogues. We employ a multi-step process to extract entities and relationships from the dialogue data and utilise these in various model settings. Our primary submission, **Triple-TowerInstruct**, integrates dialogue history into the translation process at inference, leveraging contextual cues to enhance performance across four language pairs. In addition to this, we explored an ablation study (TowerInstruct without dialogue history) and a novel model, **GraphFlanT5**, which combines graph and text embeddings within a unified framework.

### 4.1 Pre-processing

For generating the graphs, we perform entity and relation extraction by prompting GPT-4o. The prompt used for this process (see Prompt 1) is designed to extract relevant triples from the dialogue data, capturing the essence of interactions in a structured format. The system is instructed to analyse the dialogue and break it down into triples, each consisting of a subject, predicate, and object. These triples serve as the fundamental building blocks of the graph, representing the interactions between speakers.

In addition to extracting these triples, the prompt also instructs the system to annotate each triple with several attributes that provide deeper insights into the nature of the interactions. These annotations include:

- **Sentiment**: This attribute captures the emotional tone of the interaction, with values rang-

ing from -1 for negative sentiment, 0 for neutral, and 1 for positive sentiment. This allows us to understand the emotional context in which the interaction takes place.

- **Polarity**: Polarity indicates whether the interaction involves a negation, affirmation, or is neutral or questioning. It is coded as -1 for negation, 0 for neutral or questioning, and 1 for affirmation. This helps in identifying the stance or intent behind the speaker's words and keeps the predicates uniform across negation, statements and questions (e.g. "don't travel" and "travel" receive the same predicate *travel* with different polarity scores)

- **Certainty**: This attribute is on a scale from 0 (uncertain) to 1 (certain), reflecting the speaker's confidence or the definitiveness of the statement. This helps in distinguishing between statements of fact and those that are speculative or uncertain and can subsequently be used by the model to communicate certainty about its knowledge more effectively.

- **Dialogue Act**: Dialogue acts categorise the type of speech act being performed, with predefined categories such as greeting, farewell, negative reaction, positive reaction, concern, query, and others.

## 4.2   Approach 1: Triple-TowerInstruct

In our first approach, we use the TowerInstruct-7B-v0.2[2] model, a variant of the Tower (Alves et al., 2024) family specifically designed for translation-related tasks.

**TowerInstruct-7B-v0.2**   is based on the LLaMA-2 architecture, which has been extended through additional pretraining and fine-tuning to enhance its multilingual capabilities, outperforming other open models of similar scale. The model's foundation, TowerBase, was developed by continuing the pre-training on a diverse multilingual dataset across 10 languages (including Dutch, German, French, and Portuguese) incorporating both monolingual and parallel data to improve translation quality. Subsequently, TowerInstruct was fine-tuned using the TowerBlocks dataset, which includes a broad range of translation-related tasks and, relevant for the task

of chat translation, multi-turn dialogue data from UltraChat (Ding et al., 2023). This fine-tuning process tailored the model specifically for translation workflows, making it adept at handling complex, multilingual interactions.

**Prompt 1: Triple extraction with GPT-4o**

```
system_prompt =
You will analyze a dialogue and break it down
into triples consisting of a subject, predicate,
and object. Each triple should capture the
essence of interactions between speakers.
Additionally, annotate each triple with:
- Sentiment (-1 for negative, 0 for neutral,
1 for positive)
- Polarity (-1 for negation, 0 for neutral/
questioning, 1 for affirmation)
- Certainty (a scale between 0 for uncertain
and 1 for certain)
- Dialogue act (
  0 : "greeting",
  1 : "farewell",
  2 : "negative_reaction",
  3 : "positive_reaction",
  4 : "concern",
  5 : "query",
  6 : "other")

Ensure that predicates are semantically
meaningful. Separate multi-word items with
an underscore.

Save it as a JSON with this format:
{
"Conversation ID": "60250de4b",
"dialogue": [
    {
      "sender": "customer",
      "text": "I can't find my order. It was
      supposed to arrive yesterday.",
      "triples": [
        {
          "subject": "I",
          "predicate": "cannot_find",
          "object": "my_order",
          "sentiment": -1,
          "polarity": -1,
          "certainty": 1,
          "dialogue_act": 4
        },
        {
          "subject": "It",
          "predicate": "was_supposed_to_arrive",
          "object": "yesterday",
          "sentiment": -1,
          "polarity": 1,
          "certainty": 0.7,
          "dialogue_act": 4
        }]},
    {
      "sender": "agent",
      "text": "I will help you with that.",
      "triples": [
        {
          "subject": "I",
          "predicate": "will_help",
          "object": "you_with_that",
          "sentiment": 1,
          "polarity": 1,
          "certainty": 1,
          "dialogue_act": 3
        }]}]}

user_prompt = f"Analyze the following con-
versation with ID {conversation_id}:
{conversation_text}"
```

**Triple-TowerInstruct**   During inference, we merge the triple-based dialogue history, generated in the pre-processing stage (see Section 4.1), with

Figure 2: Approach 2: GraphFlanT5

the current source sentence. This combined input, which encapsulates both the conversational context and the immediate translation task, is then fed into the model to ensure that the output accurately reflects the dialogue's ongoing flow and context.

As an ablation, we also evaluated the model's performance without providing dialogue history graphs.

## 4.3 Approach 2: GraphFlanT5

We have developed a unified framework named GraphFlanT5 (see Figure 2), which integrates both graph and text input into a single architecture, similar to (Yao et al., 2023). This model is designed to generate target sequences in text based on the dialogue history and the source sequence represented in text and graph forms.

As further preprocessing for this approach, we use spaCy's NeuralCoref[3] to resolve co-references, limiting the number of nodes to a maximum of 100. These are then represented as an adjacency matrix and fed into the main model.

To encode the graph, we employ a Graph Attention Network (GAT) (Veličković et al., 2018) with a single attention layer, followed by a dense layer and normalization. On the text side, we use a Transformer encoder for encoding. We specifically used FlanT5-base[4] for its multilingual capabilities. After obtaining the encoded features from both the graph and text, we apply cross-attention to align the text representation with the graph representation. A gated fusion mechanism (Wu et al., 2021) is then used to combine the outputs of the cross-attention. Finally, the fused features are passed into the Transformer decoder to generate the final textual answer.

We fine-tuned our model for 25 epochs with a learning rate of 5e-5 and a weight decay of 0.05.

Training was conducted using mixed precision on two A10 GPUs.

## 5 Results & Discussion

### 5.1 Automated Metrics

Our primary submission, Triple-TowerInstruct, and its ablation variant without dialogue history graphs (NH) are compared against our second approach GraphFlanT5, the baseline (NLLB-200's (Team et al., 2022) 3.3B variant[5]), and the top-performing Unbabel system, using COMET-22 (Rei et al., 2022), Contextual-COMET-QE (Vernikos et al., 2022), BLEU (Papineni et al., 2002), and ChrF (Popović, 2015) scores [6].

Tables 1 and 2 show the results from our experiments across four language pairs: en-de, en-nl, en-nl, and en-pt_br. While we only submitted Approach 1 (Triple-TowerInstruct), we include the evaluation of the other approaches which were conducted after the shared task submission deadline. From the submitted approach, our team ranked second place for en-nl and en-fr, and third place for en-pt_br and en-de on the COMET-22 (Rei et al., 2020a) score.

**Triple-TowerInstruct** performed well across all language pairs, consistently outperforming the baseline based on COMET and in the majority of instances for the other metrics. For instance, in the en-de task, Triple-TowerInstruct achieved a COMET score of 91.3, outperforming the baseline's 89.8. The BLEU and ChrF scores further support this, with Triple-TowerInstruct scoring 53.0 in BLEU and 71.9 in ChrF for en-de, both above the baseline scores of 51.1 and 70.8, respectively. The

---

[3] https://github.com/huggingface/neuralcoref
[4] https://huggingface.co/google/flan-t5-base

[5] https://huggingface.co/facebook/nllb-200-3.3B
[6] Sacrebleu is used for the implementation of BLEU and ChrF (Post, 2018).

| Model | en-de | | | | en-nl | | | |
|---|---|---|---|---|---|---|---|---|
| | COMET | ChrF | BLEU | Context-COMET-QE | COMET | ChrF | BLEU | Context-COMET-QE |
| **Triple-TowerInstruct** | 91.3 | 71.9 | 53.0 | 0.2039 | 90.9 | 70.6 | 48.0 | 0.0816 |
| **TowerInstruct NH** | 91.2 | 72.2 | 53.9 | 0.2128 | 91.3 | 66.2 | 44.7 | 0.1982 |
| **GraphFlanT5** | 85.3 | 65.1 | 44.5 | 0.0120 | 88.4 | 68.5 | 48.7 | 0.0697 |
| **Baseline** | 89.8 | 70.8 | 51.1 | 0.1730 | 88.1 | 62.6 | 38.7 | 0.0873 |
| **Unbabel+it** | 92.9 | 78.2 | 62.0 | 0.2526 | 93.6 | 79.8 | 63.9 | 0.1167 |

Table 1: Translation Results for German (en-de) and Dutch (en-nl). NH models refer to ablations without dialogue history. Results for the baseline and best performing system in the task (Unbabel+it) are included for comparison.

| Model | en-fr | | | | en-pt | | | |
|---|---|---|---|---|---|---|---|---|
| | COMET | ChrF | BLEU | Context-COMET-QE | COMET | ChrF | BLEU | Context-COMET-QE |
| **Triple-TowerInstruct** | 91.6 | 75.7 | 58.8 | 0.0775 | 91.3 | 66.8 | 45.3 | 0.1909 |
| **TowerInstruct NH** | 91.7 | 75.2 | 57.9 | 0.0756 | 90.6 | 71.0 | 50.9 | 0.0686 |
| **GraphFlanT5** | 85.8 | 67.4 | 47.0 | -0.1007 | 90.4 | 75.0 | 56.7 | -0.0095 |
| **Baseline** | 90.1 | 76.2 | 58.7 | 0.0101 | 86.2 | 62.2 | 35.3 | -0.0613 |
| **Unbabel+it** | 92.8 | 79.8 | 65.7 | 0.1034 | 93.9 | 79.7 | 65.0 | 0.2367 |

Table 2: Translation Results for French (en-fr) and Portuguese (en-pt). NH models refer to ablations without dialogue history. Results for the baseline and the best performing system in the task (Unbabel+it) are included for comparison.

NH variant, which omits dialogue history, saw a minor drop in performance for en-de and en-pt_br, with a drop in COMET score of 0.1 and 0.7 respectively, and slightly lower BLEU and ChrF scores. Interestingly, the opposite is true for the en-nl and en-nl language pairs. The Context-COMET-QE scores (Rei et al., 2020b), which are intended for reference-free machine translation evaluation and trained to reflect human judgements of the quality of translations, also demonstrated variability. For en-de, Triple-TowerInstruct scored 0.2039 in Context-COMET-QE (Rei et al., 2020b), while the NH variant scored 0.2128, showing a slight improvement when dialogue history was removed. While for en-pt_br including the history increased the score by 0.0383[7]. We also observed that COMET-based metrics and n-gram matching metrics (ChrF and BLEU) disagreed in ranking our

TowerInstruct variants. When COMET favoured one variant, the n-gram metrics ranked it lower, and vice-versa. Underscoring the importance of using a combination of metrics, as relying on a single metric could give an incomplete picture of model performance.

**GraphFlanT5** which integrates graph and text input within a unified framework, showed moderate results and did not outperform our TowerInstruct variants or the baseline in most cases. In the en-de task, GraphFlanT5 recorded a COMET score of 85.3, lower than both TowerInstruct and the baseline. Its BLEU and ChrF scores were also lower, at 44.5 and 65.1, respectively. However, in some tasks like en-nl, GraphFlanT5 performed competitively with a BLEU score of 48.7, suggesting that the integration of graph representations may offer benefits in certain contexts, but requires further optimisation to be competitive to more traditional approaches.

---

[7]See Kocmi et al. (2024) for an explanation of the different dynamic ranges of the mentioned metrics.

| Model | en-de | | | | en-nl | | | |
|---|---|---|---|---|---|---|---|---|
| | **Formality** | **Lexical Cohesion** | **Pronouns** | **Verb Form** | **Formality** | **Lexical Cohesion** | **Pronouns** | **Verb Form** |
| **Triple-TowerInstruct** | 86.3 | 74.1 | 78.5 | – | 35.5 | 66.4 | – | 40.0 |
| **Baseline** | 79.4 | 76.0 | 79.1 | – | 53.0 | 57.4 | – | 35.7 |
| **Unbabel+it** | 88.6 | 82.9 | 70.5 | – | 93.9 | 87.7 | – | 54.5 |

Table 3: F1 Scores for German (en-de) and Dutch (en-nl) across different evaluation dimensions of MUDA. Where entries are left blank, the metric does not evaluate the language for that dimension.

| Model | en-fr | | | | en-pt | | | |
|---|---|---|---|---|---|---|---|---|
| | **Formality** | **Lexical Cohesion** | **Pronouns** | **Verb Form** | **Formality** | **Lexical Cohesion** | **Pronouns** | **Verb Form** |
| **Triple-TowerInstruct** | 89.6 | 78.6 | 88.6 | 68.1 | 78.7 | 88.5 | 55.0 | – |
| **Baseline** | 86.9 | 82.1 | 82.0 | 70.2 | 45.7 | 81.0 | 55.8 | – |
| **Unbabel+it** | 91.3 | 90.2 | 92.9 | 74.2 | 88.0 | 95.5 | 74.4 | – |

Table 4: F1 Scores for French (en-fr) and Portuguese (en-pt) across different evaluation dimensions of MUDA.

### 5.1.1 MUDA

Tables 3 and 4 present the F1 scores for different evaluation dimensions—Formality, Lexical Cohesion, Pronouns, and Verb Form—of the Multilingual Discourse-Aware (MuDA) benchmark (Fernandes et al., 2023). We compared our primary model, Triple-TowerInstruct, against the baseline and the top-performing system, Unbabel+it. MuDA is designed to systematically evaluate machine translation models on their handling of discourse phenomena that require context. Unlike traditional metrics that focus broadly on translation accuracy, it specifically targets the model's ability to correctly translate discourse elements, such as pronouns and verb forms, that depend heavily on the surrounding context.

The performance of our model varied across different dimensions and language pairs, outperforming the baseline in 7 out of 13 cases. Overall, it demonstrated relatively strong performance on the **Formality** dimension, achieving competitive F1 scores in language pairs such as en-de, en-nl, and en-pt_br, with a notable increase of 33 points over the baseline for the latter. The exception was the en-nl pair, where the model's formality score was notably lower compared to both the baseline and top-performing systems, indicating a need for targeted improvements in handling formality specific to Dutch translations. However, performance on **Lexical Cohesion**, **Pronouns**, and **Verb Form** was less consistent across language pairs, with the model outperforming the baseline in only half of the cases.

### 5.2 Human Evaluation

Human evaluation confirms that our approach outperforms the baseline, and ranked second place for en-nl and en-fr, and third place for en-pt_br and en-de across all submitted approaches.

| | **en-de** | **en-nl** | **en-fr** | **en-pt** |
|---|---|---|---|---|
| **Triple-TowerInstruct** | 78.6 | 84.37 | 73.32 | 69.85 |
| **Baseline** | 74.5 | 53.07 | 67.81 | 56.37 |
| **Unbabel+it** | 84.22 | 92.22 | 79.62 | 78.0 |

Table 5: Human Evaluation Scores on document level for German (en-de), Dutch (en-nl), French (en-fr), and Portuguese (en-pt) across models.

The human evaluation was facilitated by the task organisers. It was conducted by professional linguists and translators using a combination of Direct Assessment and scalar quality metric (DA+SQM) implemented via the Appraise framework (Federmann, 2018).

## 6 Conclusion & Future Work

Our results underscore the importance of incorporating dialogue history in improving translation quality, highlighting its role in maintaining coherence and context throughout chat-based translations. The integration of graph-based representations also shows promise, particularly in capturing and leveraging the structural relationships within dialogue contexts. However, our findings indicate that further optimisation is required to fully realise the benefits of this approach, especially in terms of consistently outperforming more traditional text-based models.

In future work, one of our key objectives is to combine the strengths of TowerInstruct's translation capabilities with the advanced context modelling offered by our graph-based approach. By integrating these two methodologies, we aim to create a more robust system that can better handle the complexities of chat dialogue translation.

Furthermore, we plan to investigate the incorporation of additional contextual information, such as certainty or sentiment scores derived during preprocessing. These scores could potentially enhance the model's ability to weigh different parts of the dialogue based on their reliability and emotional tone, thereby improving overall translation accuracy. By factoring in sentiment, the model can better preserve the nuances of emotional expression within the conversation, leading to more contextually appropriate translations, which is particularly important in the task's customer service domain where frustration is common. By pursuing these directions, we aim to refine our models further, making them more adaptable and effective in real-world chat translation and dialogue tasks.

## Acknowledgements

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An Open Multilingual Large Language Model for Translation-Related Tasks. *Preprint*, arXiv:2402.17733.

Selene Baez Santamaria, Lea Krause, Lucia Donatelli, and Piek Vossen. 2023. The Role of Personal Perspectives in Open-Domain Dialogue: Towards Enhanced Data Modelling and Long-term Memory. *Proceedings of BNAIC/BeNeLearn the Joint International Scientific Conferences on AI and Machine Learning*, pages 1–19.

Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and Jia Li. 2023. Large language models meet harry potter: A dataset for aligning dialogue agents with characters. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8506–8520, Singapore. Association for Computational Linguistics.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3029–3051, Singapore. Association for Computational Linguistics.

Ana C Farinha, M. Amin Farajian, Marianna Buchicchio, Patrick Fernandes, José G. C. de Souza, Helena Moniz, and André F. T. Martins. 2022. Findings of the WMT 2022 shared task on chat translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 724–743, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Christian Federmann. 2018. Appraise evaluation framework for machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 86–88, Santa Fe, New Mexico. Association for Computational Linguistics.

Patrick Fernandes, Kayo Yin, Emmy Liu, André Martins, and Graham Neubig. 2023. When does translation require context? a data-driven, multilingual exploration. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 606–626, Toronto, Canada. Association for Computational Linguistics.

Madalena Gonçalves, Marianna Buchicchio, Craig Stewart, Helena Moniz, and Alon Lavie. 2022. Agent and user-generated content and its impact on customer support MT. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 201–210, Ghent, Belgium. European Association for Machine Translation.

Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024. Navigating the metrics maze: Reconciling score magnitudes and accuracies. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1999–2014, Bangkok, Thailand. Association for Computational Linguistics.

Zhibin Liu, Zheng-Yu Niu, Hua Wu, and Haifeng Wang. 2019. Knowledge aware conversation generation with explainable reasoning over augmented graphs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1782–1792, Hong Kong, China. Association for Computational Linguistics.

Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. Evaluating very long-term conversational memory of llm agents. *Preprint*, arXiv:2402.17753.

Sameen Maruf, Fahimeh Saleh, and Gholamreza Haffari. 2021. A survey on document-level neural machine translation: Methods and evaluation. *ACM Comput. Surv.*, 54(2).

Diego Moussallem, Matthias Wauer, and Axel-Cyrille Ngonga Ngomo. 2018. Machine translation using semantic web technologies: A survey. *Journal of Web Semantics*, 51:1–19.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020a. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020b. Unbabel's participation in the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 911–920, Online. Association for Computational Linguistics.

Nathaniel Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. ChatGPT MT: Competitive for high- (but not low-) resource languages. In *Proceedings of the Eighth Conference on Machine Translation*, pages 392–418, Singapore. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.

Yi-Lin Tuan, Yun-Nung Chen, and Hung-yi Lee. 2019. DyKgChat: Benchmarking dialogue generation grounding on dynamic knowledge graphs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1855–1865, Hong Kong, China. Association for Computational Linguistics.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. *Preprint*, arXiv:1710.10903.

Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 118–128, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.

Zhiyong Wu, Lingpeng Kong, Wei Bi, Xiang Li, and Ben Kao. 2021. Good for misconceived reasons: An empirical revisiting on the need for visual context in multimodal machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6153–6166, Online. Association for Computational Linguistics.

Jing Xu, Arthur Szlam, and Jason Weston. 2021. Beyond goldfish memory: Long-term open-domain conversation. *arXiv preprint arXiv:2107.07567*.

Jun Xu, Haifeng Wang, Zheng-Yu Niu, Hua Wu, Wanxiang Che, and Ting Liu. 2020. Conversational graph grounded policy learning for open-domain conversation generation. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 1835–1845.

Yao Yao, Zuchao Li, and Hai Zhao. 2023. Beyond chain-of-thought, effective graph-of-thought reasoning in large language models. *Computing Research Repository*, arXiv:2305.16582.

Houyu Zhang, Zhenghao Liu, Chenyan Xiong, and Zhiyuan Liu. 2020. Grounded conversation generation as guided traverses in commonsense knowledge graphs. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2031–2043, Online. Association for Computational Linguistics.

Yang Zhao, Jiajun Zhang, Yu Zhou, and Chengqing Zong. 2021. Knowledge graphs enhanced neural machine translation. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 4039–4045.

# Reducing Redundancy in Japanese-to-English Translation: A Multi-Pipeline Approach for Translating Repeated Elements

**Qiao Wang**
Waseda University
Tokyo, Japan
judy.wang@aoni.waseda.jp

**Yixuan Huang**
Waseda University
Tokyo, Japan
yixuan.huang@moegi.waseda.jp

**Zheng Yuan**
King's College London
London, UK
zheng.yuan@kcl.ac.uk

## Abstract

This paper presents a multi-pipeline Japanese-to-English machine translation (MT) system designed to address the challenge of translating repeated elements from Japanese into fluent and lexically diverse English. The system was developed as part of the Non-Repetitive Translation Task at WMT24, which focuses on minimizing redundancy while maintaining high translation quality. Our approach utilizes MeCab, the de facto Natural Language Processing (NLP) tool for Japanese, to identify repeated elements, and Claude Sonnet 3.5, a Large Language Model (LLM), for translation and proofreading. The system effectively accomplishes the shared task by identifying and translating in a diversified manner 89.79% of the 470 repeated instances in the test dataset and achieving an average translation quality score of 4.60 out of 5, significantly surpassing the baseline score of 3.88. The analysis also revealed challenges, particularly in identifying standalone noun-suffix elements and occasional cases of consistent translations or mistranslations.

## 1 Introduction

In the Japanese language, repetition at the word and phrasal levels is frequently employed (Fujimura-Wilson, 2007). One reason for this is that Japanese is a topic-prominent language, where the topic of the sentence is often explicitly stated and reiterated to ensure clarity and prominence (Tsujimura, 2013). Additionally, Japanese is highly context-dependent and typically omits subject pronouns, relying on the repetition of key nouns and verbs to maintain coherence (Maynard, 1997). Specifically for personal names, repetition is commonly used instead of pronouns to convey politeness and respect (Mogi, 2000).

In contrast, English typically favors variety and succinctness to maintain reader engagement (Hinkel, 2002; Halliday, 1994). Research in translation studies emphasizes the importance of lexical variety to ensure fluency and readability in translated texts (Baker, 1992; Newmark, 1988). Therefore, effectively translating repeated elements from Japanese to English may require the use of more diverse expressions while ensuring consistency and clarity.

The Non-Repetitive Translation Task at WMT24 addresses the challenge of translating repeated elements from Japanese into English (Kinugawa et al., 2024). This task aims to develop machine translation (MT) systems capable of identifying repeated expressions in Japanese text and translating them into lexically diverse and fluent English sentences. Participants are provided with training and test datasets comprising Japanese-English parallel corpora, in both raw and annotated formats with repeated targets tagged. Systems are evaluated on their ability to minimize redundancy while maintaining high translation quality.

Our contribution includes the development of a multi-pipeline MT system that effectively avoids redundancy in translating repeated words and phrases from the source Japanese text. Specifically, we utilized MeCab (Kudo, 2005) for tokenization and lemmatization of Japanese sentences to identify repeated elements and adopted the Large Language Model (LLM) Claude Sonnet 3.5 (Anthropic, 2024) for translation and proofreading. When compared with the baseline system provided by the task organizers, our system achieved an average translation score of 4.60, significantly higher than that of the baseline system at 3.88; and a BLEU metric of 24.4

compared with human benchmark translation.

## 2 Related work

The identification of repeated elements in Japanese poses unique challenges due to the language's agglutinative nature (Tsujimura, 2013). Unlike Indo-European languages where word boundaries are clearly delineated (Baker, 1992), Japanese requires sophisticated tokenization methods to accurately segment text into meaningful units. As far as we are aware, no previous NLP studies have specifically focused on identifying repeated tokens in Japanese, although some have explored the identification of repetition at the semantic level, not necessarily of the same words (Kawamoto, 2022). However, the repeated elements in this shared task need to be identical, the most straightforward method is to adopt an Natural Language Processing (NLP) tool designed specifically for the Japanese language. Previous studies have recommended MeCab as the de facto text segmentation library capable of part-of-speech (POS) tagging, tokenization and lemmatization for Japanese (Kudo, 2005). Even with MeCab, challenges in identifying repeated elements persist. For example, for "国文学" (Japanese literature) and "漢文学" (Chinese literature), humans may easily detect the repeated element "文学" (literature). However, for MeCab, "国文学" is regarded as one token, while "漢文学" is regarded as two tokens, "漢" and "文学". As such, no repetition can be detected as the tokens do not match.

Recent advancements in LLMs have prompted us to explore them for identifying repeated elements in Japanese text. Prior to building our system, we experimented with GPT-3.5[1] (fine-tuned on the WMT24 training dataset), GPT-4[2] (used direct prompting), and Claude Sonnet 3.5[3] (also used direct prompting) to assess their ability to detect repeated tokens. However, none of these LLMs consistently identified repetitions. This indicates that rule-based approaches using MeCab remain necessary.

Machine translation (MT) between Japanese and English has long been challenging due to the significant linguistic differences between the two languages (Wang et al., 2022). While Neural Machine Translation (NMT) systems, particularly those based on Transformer models (Vaswani et al., 2023), have demonstrated success in handling structured text, they face limitations when dealing with informal language, idiomatic expressions, and culturally specific references. Meanwhile, LLMs such as GPT-3.5 and GPT-4 have shown considerable promise in improving translation quality in Japanese-English (JA-EN) tasks (Vincent et al., 2024). These models benefit from extensive training on diverse datasets, which allows them to generate more contextually appropriate translations, particularly in cases where traditional supervised NMT systems may struggle (Siu, 2023).

Previous studies have compared NMT systems and LLMs in translating high-resource and low-resource languages and found that LLMs such as GPT series perform better in high but not low-resource languages (Robinson et al., 2023). Since Japanese is not considered a low-resource language, we expect LLMs to perform better in this shared task. As thus, we piloted with Google Translate[4], GPT4o and Claude Sonnet 3.5 with the training dataset and found the Claude Sonnet 3.5 performed the best in translating the Japanese texts.

In addition to translation, we aimed to incorporate a proofreading pipeline to enhance overall translation quality. This idea was inspired by the translation-review-revise methodology, which emphasizes iterative improvement of translated content through multiple stages of refinement (Ng, 2024). In line with this methodology, we aimed to incorporate a proofreading process with LLMs as well. In the following, we'll describe in detail the system design and implementation results of our system for the shared task.

## 3 System description

Our system comprises of four pipelines: 1) the preprocessing pipeline for identifying and subsequently assigning IDs and reference numbers

---

Figure 1: Workflow of the system

(ref numbers) to repeated elements, i.e., targets, in Japanese sentences; 2) the translation pipeline for translating the Japanese sentences into English while trying to reduce redundancy; 3) the proofreading pipeline for revising the translated sentences; and 4) the post-processing pipeline for appending the types of strategies, i.e., substitution or reduction, to the IDs and ref numbers of the repeated targets in the Japanese sentences. The input of the system is the raw Japanese sentences while the output includes the Japanese sentences with repeated targets assigned IDs, ref numbers and types, and raw translated English texts. Figure 1 shows the workflow of the four pipelines. The source code of the developed system can be found at our GitHub repository[5].

## 3.1 Preprocessing pipeline

For identifying and subsequently assigning IDs and ref numbers to repeated targets in Japanese sentences, we adopted MeCab. The POS tags in MeCab are structured into a maximum of three hierarchical layers. For example, a three-layered POS tag can be "名詞-普通名詞-副詞可能" (noun-common-adverbial), with the top layer indicating the token is a noun, the second layer further explaining it is a common noun, and the third layer showing that the noun can also be used as an adverb.

---

[5] https://github.com/judywq/non-repetitive-translation

**Step 1: POS tagging** We performed POS tagging on all the Japanese sentences in the training dataset, creating a pool of tokens with their POS tags.

**Step 2: POS screening** From the pool, we first did a simple token match to find repeated tokens. Then, we compared the repeated tokens with the tagged repeated targets in the training dataset to decide what POS tags should be included and what should not in identifying repeated tokens. The result was a 'whitelist' of POS tags at different layers. In the first layer, we focused only on content words, including nouns, verbs, adjectives, adverbs and prefixes, while excluding function words such as auxiliaries, conjunctions and particles. For the POS tags from the second layer on, we maximized the coverage of POS tags found in the training dataset while reducing noises. Table 1 shows the POS tags in the 'whitelist'. Blanks in the third layer indicate that all the third-layer tags have been selected.

**Step 3: Identifying targets** When identifying the repeated targets, the easiest way is to do exact match. However, for some POS tags, special treatment was necessary. These include Verb, Noun-Suffix, and Prefix-Noun connection.

For Verb, we adopted their lemmatized forms using MeCab. This is due to the fact that Japanese verbs are rich in inflections. For example, the verb "食べる" (*taberu*, meaning "to eat") can appear in various forms depending on the tense, politeness level, and grammatical context, such as "食べた" (*tabeta*, past tense "ate"), "食べます"(*tabemasu*, polite form "eat"), or "食べられる" (*taberareru*, potential form "can eat").

For Noun-Suffix and Prefix-Noun, as they are suffixes or prefixes, they should be bound to another token. We thus added a rule where when tokens with the two POS tags are repeated, their neighboring tokens, i.e., the token before the suffix and the token after the prefix, should also be bound together with them. If the bound elements still match, then they are valid targets. Otherwise, they will be dismissed. For example, with the Noun-Suffix "者" (person), if it is preceded twice by the verb "容疑"(suspect), the compound noun

| Layer 1 | Layer 2 | Layer 3 |
|---|---|---|
| Adverb (副詞) | General (一般) | |
| Verb (動詞) | Independent (自立) | |
| Noun (名詞) | Suffix (接尾) | Adjectival noun stem (形容動詞語幹) |
| | | Personal name (人名) |
| | | Area (地域) |
| | | Special (特殊) |
| | Dependent (非自立) | Adverbial (副詞可能) |
| | Suru verb (サ変接続) | |
| | Nai adjective stem (ナイ形容詞語幹) | |
| | General (一般) | |
| | Proper noun (固有名詞) | |
| | Adverbial (副詞可能) | |
| | Adjectival noun stem (形容動詞語幹) | |
| Adjective (形容詞) | Independent (自立) | |
| Prefix (接頭詞) | Noun connection (名詞接続) | |

Table 1: POS tags in the 'whitelist'

"容疑者"(a suspected person) will be the valid target. However, if it is preceded by "容疑" once and by "被爆"(be bombed) the second time, they will be dismissed as "容疑者" does not match "被爆者"(an atomic bomb victim).

**Step 4: Assigning IDs and ref numbers** We assigned IDs and ref numbers to targets based on their order of occurrence in the Japanese sentence. The output of step 4 are Japanese sentences with repeated targets tagged.

## 3.2 Translation pipeline

The translation pipeline is responsible for translating the Japanese sentences into English while trying to adopt diversified expressions for the repeated targets tagged in the Japanese sentences; singling out the translation for each occurrence of the targets; and deciding which type and subtype of strategy was used. It should be noted that to facilitate understanding of output from the pipeline, we introduced two new types: "first occurrence" and "consistency" and another subtype: "literal translation", to the original types (substitution and reduction) and subtypes from the official website of the task (WMT24, 2024). First occurrence is assigned to translation of a target where the translation appeared for the first time and thus there is no need to re-

duce redundancy. Consistency refers to situations where the target is translated into the same English expressions across multiple occurrences. The subtype "literal translation" is added to complement "non-literal translation" original included in the examples from the official website.

For this pipeline, we adopted Claude Sonnet 3.5 with few-shot in-context learning prompting techniques. In our prompt, we included the explanation of the shared task and the examples of reduction and substitution from the task's official website. Then we asked the translation pipeline to translate the sentences while trying its best to adopt diversified expressions for the repeated targets tagged in the Japanese sentences, single out the translations for the targets and decide which type and subtype of strategy was used.

For example, for the input sentence:

> RCEP では、7 月 1 日に東京で閣僚 <target id=0 ref=0> 会合 </target> が開かれ、妥結に向け 11 月下旬に シンガポールで首脳 <target id=0 ref=1> 会合 </target> が開催でき るよう、交渉に注力する方針を確認。

the output from the translation pipeline is as follows:

```
{
"en_translation": "For RCEP, a ministerial
```

```
meeting was held in Tokyo on July 1,
confirming the policy to focus on negotiations
with the aim of holding a summit in Singapore
in late November to reach an agreement.",
  "targets":
  [
    {
      "id": "0",
      "ref": "0",
      "ja_element": "会合",
      "en_element": "meeting",
      "type": "f",
      "subtype": "lt"
    },
    {
      "id": "0",
      "ref": "1",
      "ja_element": "会合",
      "en_element": "summit",
      "type": "s",
      "subtype": "syn"
    }
  ]
}
```

### 3.3 Proofreading pipeline

For the proofreading pipeline, we also adopted Claude Sonnet 3.5. Few-shot in-context learning prompt was also designed for this pipeline as for the translation pipeline. The proofreading pipeline receives the Japanese sentence and the translated text, the translations for each occurrence of repeated targets and the types and subtypes of strategies used. It is asked to check if the translation is faithful to the Japanese sentence and if redundancy can be further reduced. If changes are not necessary, it returns

```
{"changed": "No"}
```

Otherwise, it returns

```
{"changed": "Yes"}
```

followed by a revised output in the same format as the translation pipeline. A sample output from the pipeline is shown below:

```
{
  "changed": "Yes",
  "en_translation_updated": "Toshiba stated
  that there is no change to its previous
  projection, as the reversal is already
  incorporated into the full-year earnings
  outlook for the fiscal year ending
  March 2019.",
  "targets_updated":
  [
    {
      "id": "0",
      "ref": "0",
      "ja_element": "予想",
      "en_element": "outlook",
      "type": "f",
```

```
      "subtype": "lt"
    },
    {
      "id": "0",
      "ref": "1",
      "ja_element": "予想",
      "en_element": "projection",
      "type": "s",
      "subtype": "syn"
    }
  ]
}
```

### 3.4 Post-processing pipeline

With the proofread sentence translation and translations of each occurrence of repeated targets together with the types and subtypes of strategies, the post-processing pipeline first appends the types of the translations to the ID and ref number of each occurrence of repeated targets in the Japanese text. Then it replaces the two types we added with the two official types, reduction or substitution. Specifically, in the case of first occurrence, the type of same ID but the following/previous ref number target will be adopted. For example, if for a target A where it appears twice in a sentence, it has two occurrences: ID=0, Ref=0, Type= first occurrence; and ID=0, Ref=1, Type= substitution. The type in ref=0 is replaced by the type in ref=1, substitution. For consistency, which means our system deems it unnecessary to reduce redundancy, we remove the IDs and ref numbers of the targets in the Japanese text. This is based on the official dataset where only the targets that require redundancy reduction are tagged with IDs and ref numbers.

After the post-processing pipeline, the system outputs the Japanese sentences with IDs, ref numbers and types tagged for repeated targets and the raw English translation without any tags.

We also considered situations where our system may fail to identify any repeated targets in the pre-processing pipeline. In such cases, the raw Japanese sentences will be translated into English by Claude Sonnet 3.5 and the system will directly output the raw Japanese sentence and its English translation.

1051

| JPO Adequacy | \<NON-REP\> | \<REP\> | \<INCORRECT\> | Total |
|:---:|:---:|:---:|:---:|:---:|
| [5,5] | 127 | 20 | 0 | 147 |
| [4,5) | 280 | 17 | 3 | 300 |
| [3,4) | 15 | 1 | 7 | 23 |
| [2,3) | 0 | 0 | 0 | 0 |
| [1,2) | 0 | 0 | 0 | 0 |
| Total | 422 | 38 | 10 | 470 |

Table 2: JPO adequacy and translation style

## 4 System implementation and results

Our system was applied to translate 470 sentences from the test dataset provided by the Non-Repetitive Translation Task at WMT24. The output from our system was rigorously evaluated by three human raters assigned by the task organizers. Each rater independently reviewed the translation of the repeated targets and assigned each target a translation quality score ranging from 1 to 5 based on criterion for translated patent documents from the Japanese Patent Office (JPO), with 5 representing the highest quality (JPO adequacy). The raters also labelled each translation target with one of the following translation styles: "c" (consistent/repetitive translation), "s" (substitution), "r" (reduction), or "m" (mistake in translation). The final evaluation score for each sentence was determined by averaging the scores given by the three raters and the label was determined by a majority vote.

In the 470 sentences, there are a total of 489 repeated targets tagged by the organizers in the dataset, meaning that some sentences contain more than 1 repeated target. When there are multiple targets in one sentence, the evaluation of all targets is aggregated to one by the organizers, resulting in 470 evaluation instances for the 470 sentences in total. Results suggest that our system produced 38 instances of repetitive translations, 422 instances of correct non-repetitive translations, and 10 incorrect translations. This shows a correct non-repetitive translation rate of 89.79% for our system.

Table 2 shows the detailed evaluation results

including the instance counts of translation styles (non-repetitive, repetitive and incorrect translation) and translation quality (JPO adequacy).

The average JPO adequacy of our system is 4.60, significantly higher than that of the baseline system at 3.88 (t=14.09, $p<0.00$). To view the balance between translation quality and style, the JPO adequacy score for each instance are converted to 0 if its style is not '\<NON-REP\>', i.e., correct non-repetitive translation. The average of this filtered JPO scores is 4.13. This is significantly higher than the baseline system at 2.13 (t=16.60, $p<0.00$).

For reference purposes, our system's performance was also evaluated using the BLEU metric (Papineni et al., 2002). The BLEU score for our system was 24.4, indicating moderate similarity to the human benchmark translations provided by the organizers. The verbose BLEU score breakdown shows a precision of 58.3% for 1-grams, 30.0% for 2-grams, 18.0% for 3-grams, and 11.3% for 4-grams. No Brevity Penalty (BP) was applied, as the length of the system's output (15,700 words) closely matched that of the benchmark (15,579 words), with a length ratio of 1.008.

## 5 Discussions

Our system demonstrated high performance in the shared task, effectively combining multiple NLP and LLMs pipelines to achieve impressive results. However, some issues were observed that highlight the challenges of this task and the limitations of our approach.

One challenge our system faced was in certain tagging targets deemed repetitive by hu-

man raters, particularly those related to standalone noun-suffix elements. For example:

**Japanese:** 専門家の 1 人は、鑑定した 110 個の遺骨の中で日本人の DNA<target_1> 型 </target_1> は 5 個、フィリピン人の <target_1> 型 </target_1> が 54 個だったと報告。

**System Translation:** One of the experts reported that among the 110 bone samples examined, 5 had DNA patterns matching Japanese individuals, while 54 had patterns matching Filipino individuals.

In this instance, our system failed to identify the target "型" (type), a noun-suffix, as a repeated element and subsequently the translation style was consistent. According to Step 3 in the pre-processing pipeline, noun-suffixes are only identified as repeated targets if their preceding tokens also match. This rule was implemented to reduce noise; however, it led to the negligence of independent suffixes tagged by humans in the benchmark dataset. A similar issue occurred with the word "量" (quantity) in the following sentence:

**Japanese:** 国際貨物 <target_1> 量 </target_1> は 10%減の 16 万 8510 トン。|| ジェット燃料給油 <target_1> 量 </target_1> は 2%減の 37 万 3805 キロリットルだった。

**System Translation:** International cargo volume decreased by 10% to 168,510 tons. Jet fuel supply volume fell by 2% to 373,805 kiloliters.

In this example, "量" was not identified as a repeated target, resulting in the translation being repeated instead of diversified.

Moreover, there were instances where our system correctly identified repeated targets but still opted for consistent translations. This occurred in cases where the system judged consistent translation to be preferable or where it misinterpreted derivatives or inflected forms of a word as non-repetitive. For example:

**Japanese:** 同氏は <target_1> 五輪 </target_1> には 04 年アテネ <target_1> 五輪 </target_1> から 4 大会連続で出場した。

**System Translation:** He competed in four consecutive Olympic Games, starting with the 2004 Athens Olympics.

In this case, our system treated "Olympic Games" and "Olympics" as substitutions. However, human raters considered these terms to be consistent translations, as the word "Olympic" and its inflected form "Olympics" are essentially the same. Though we stated specifically in our prompts to Claude that derivatives and inflected forms are consistent translations, it failed to perform the translation/proofreading task adequately in some cases.

In the 10 instances where the system produced incorrect translations, these errors occur because the translations are too flexible and non-literal, compromising the "faithfulness" of translation. An example of such an error is:

**Japanese:** 一方、66 歳以降も働きたいと答えた人が挙げた理由は、「<target_1> 経済的 </target_1> にゆとりある生活を送りたい」が 28.9%、「働き続けないと生活費が足りないと思う」24.9%などで、<target_1> 経済的 </target_1> な理由が半数を超えた。

**System Translation:** On the other hand, among those who expressed a desire to continue working beyond the age of 66, over half cited financial reasons. The most common reasons were 'wanting to maintain a comfortable lifestyle' (28.9%) and 'believing that living expenses would be insufficient without continued employment' (24.9%).

In this example, "the most common reasons" should be interpreted as "the most common comments given by those who cited financial reasons". However, the translation assumes that readers can infer from the text, but the same word "reasons" makes the sentence confusing.

## 6   Conclusions

In conclusion, the proposed multi-pipeline Japanese-to-English machine translation system successfully addresses the challenge of translating repeated elements from Japanese into fluent and varied English. By integrating MeCab for accurate tokenization and Claude Sonnet 3.5 for translation and proofreading, the system achieved a high rate of correct non-repetitive translations, with a translation quality score that significantly exceeded the baseline. However, certain challenges remain, particularly in identifying and translating standalone noun-suffix elements and in cases where consistent translation is deemed preferable. Additionally, the study highlighted the limitations of current human evaluation processes, where inter-rater reliability was low, affecting the consistency of the evaluation results. Future work could explore more advanced language models and refined evaluation methodologies to further enhance the system's performance and address these challenges.

## Limitations

In our translation pipeline, we compared the performance of Claude with Google Translate and GPT-4 before selecting Claude as the translation model. However, it is important to acknowledge that more effective LLMs may emerge in the future, which could offer improved performance. Besides, one of the inherent issues of relying on commercial LLMs like Claude is the issue of token limits, which can pose challenges in large-scale projects where the tasks requires days to complete.

Furthermore, the inter-rater reliability among the human raters was relatively low. We noticed that one of the raters was conspicuously more severe in their evaluations compared to the other two raters. The inter-rater reliability analysis also revealed only a slight agreement among the raters for JPO Adequacy, with an average Weighted Kappa (Cohen, 1968) of 0.161. The Fleiss' Kappa (Fleiss, 1971) for Style was -0.204, suggesting that the agreement among the raters was not only poor but worse than what would be expected by chance. This means that the evaluation results may have differed if the inter-rater reliability was higher. To illustrate, the raters did not reach consensus on some mistranslations. The following shows an example:

> **Japanese:** JAXA によると、<target_1> クレーター </target_1> は直径 10 メートル規模と推測され、小惑星への人工 <target_1> クレーター </target_1> 生成に成功したのは世界で初めてだという。

> **System Translation:** According to JAXA, the crater is estimated to be about 10 meters in diameter, marking the world's first successful artificial impact on an asteroid.

In this example, the word "impact" can be considered a term referring to "a collision between astronomical objects causing measurable effects" (Rumpf et al., 2017) in planetary science, which usually results in the formation of an impact crater. In this sense, it may be a substitution to "crater". Indeed, one rater regarded it as an appropriate substitution and scored it a 5, while the other two raters considered it an incorrect translation. Such disagreement highlights the importance of a more rigorous and standardized human evaluation process in future tasks.

## Acknowledgement

## References

Anthropic. 2024. Claude 3.5. Large language model, retrieved from https://www.anthropic.com/claude.

Mona Baker. 1992. *In Other Words: A Coursebook on Translation*. Routledge.

Jacob Cohen. 1968. Weighted kappa: nominal scale agreement with provision for scaled disagreement or partial credit. *Psychological Bulletin*, 70(4):213–220.

Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382.

Kayo Fujimura-Wilson. 2007. Japanese exact repetitions involving talk among friends. *Discourse Studies*, 9(3):319–339.

M. A. K. Halliday. 1994. *An Introduction to Functional Grammar*. Edward Arnold.

Eli Hinkel. 2002. *Second Language Writers' Text: Linguistic and Rhetorical Features*. Lawrence Erlbaum Associates.

Toshiki Kawamoto. 2022. Generating repetitions with appropriate repeated words. *Journal of Natural Language Processing*, 29(4):1302–1307.

Kazutaka Kinugawa, Hideya Mino, Isao Goto, and Naoto Shirai. 2024. Findings of the wmt 2024 shared task on non-repetitive translation. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

Taku Kudo. 2005. Mecab: Yet another part-of-speech and morphological analyzer.

Senko K. Maynard. 1997. *Japanese Communication: Language and Thought in Context*. University of Hawai'i Press.

Norie Mogi. 2000. Japanese ways of addressing people. *Investigationes Linguisticae*, VIII. ∗.

Peter Newmark. 1988. *A Textbook of Translation*. Prentice Hall.

Andrew Ng. 2024. Machine learning yearning: Technical strategy for ai engineers, in the era of deep learning. DeepLearning.AI.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 311–318, Morristown, NJ, USA. Association for Computational Linguistics.

Nathaniel R. Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. Chatgpt mt: Competitive for high- (but not low-) resource languages. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, page 392–418. Association for Computational Linguistics.

Clemens M. Rumpf, Hugh G. Lewis, and Peter M. Atkinson. 2017. Asteroid impact effects and their immediate hazards for human populations. *Geophysical Research Letters*, 44(8):3433–3440.

Sai Cheong Siu. 2023. Chatgpt and gpt-4 for professional translators: Exploring the potential of large language models in translation. Available at SSRN: https://ssrn.com/abstract=4448091 or http://dx.doi.org/10.2139/ssrn.4448091.

Natsuko Tsujimura. 2013. *An Introduction to Japanese Linguistics*. Wiley-Blackwell.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *arXiv*, arXiv:1706.03762. Doi: 10.48550/arXiv.1706.03762.

Michael Vincent, Giovanni Sutanto, De Giacomo Gatti, Toshiaki Nakazawa, and Masaru Yamada. 2024. Chatgpt as a translation engine: A case study on japanese-english. In *Proceedings of the Machine Translation Summit XVIII*, pages 95–98.

Haifeng Wang, Hua Wu, Zhongjun He, Liang Huang, and Kenneth Ward Church. 2022. Progress in machine translation. *Engineering*, 18:143–153.

WMT24. 2024. Non-repetitive translation task. Website.

# SYSTRAN @ WMT24 Non-Repetitive Translation Task

**Marko Avila** and **Josep Crego**

SYSTRAN by ChapsVision
5 rue Feydeau,
75002 Paris, France
{mavila,jcrego}@chapsvision.com

## Abstract

Many contemporary NLP systems rely on neural decoders for text generation, which demonstrate an impressive ability to generate text approaching human fluency levels. However, in the case of neural machine translation networks, they often grapple with the production of repetitive content, also known as repetitive diction or word repetition, an aspect they weren't explicitly trained to address. While not inherently negative, this repetition can make writing seem monotonous or awkward if not used intentionally for emphasis or stylistic purposes. This paper presents our submission to the WMT 2024 Non-Repetitive Translation Task, for which we adopt a repetition penalty method applied at learning inspired by the principles of label smoothing. No additional work is needed at inference time. We modify the ground-truth distribution to steer the model towards discouraging repetitions. Experiments show the ability of the proposed methods in reducing repetitions within neural machine translation engines, without compromising efficiency or translation quality.

## 1 Introduction

The Non-Repetitive Translation Task of the ninth Conference on Machine Translation (WMT24) focuses on lexical choice in machine translation, especially choice regarding repeated words in a source sentence. Generally, the repetition of the same words can create a monotonous or awkward impression in English, and it should be appropriately avoided. Typical workarounds in monolingual writing are to

1) remove redundant terms if possible (reduction) or

2) use alternative words such as synonyms as substitutes (substitution).

These techniques are also observed in human translations. The goal of this task is to study how these techniques can be incorporated into machine translation systems to enrich lexical choice capabilities. From a practical standpoint, such capability would be important, for example, in news production, where high quality text that goes beyond robotic word-by-word translation is required.

In addition, repetitions do not always have a negative impact on readability. Without aiming to be exhaustive : i) repetitions play a role when summarizing information or reinforcing a concept ; ii) common expressions are formed using word repetitions, and altering them to eliminate repetition would alter their intended meaning ; iii) in highly specialized domains, expressions convey precise meanings that disallow being reformulated. The following examples illustrate these observations :

  i) once closed, the door stays closed

  ii) over and over ; to be or not to be ; step by step

  iii) the congenital muscular dystrophy in newborns presenting with muscular hypotonia

As previously introduced, finding suitable alternatives without altering the meaning of a sentence can be a challenging task.

Participants are required to control a machine translation system using reduction or substitution so that it does not output the same words for certain repeated words in a source sentence. The translation direction is Japanese to English.

## 2 Related Work

The fluency levels achieved by LLMs are widely acknowledged to be high, primarily owing to the extensive availability of monolingual datasets, which surpasses that of standard neural machine translation (NMT) models trained solely on parallel texts. To the best of our knowledge, no dedicated research has been conducted on addressing the repetition issue tackled in this work within NMT systems.

Closely related, Welleck et al. (2019) describe a method to train neural language models that in addition to maximizing likelihood to model the overall sequence probability distribution, also includes an unlikelihood term in the loss function to correct known biases such as repeated tokens. Li et al. (2020) use the same approach to control copy effect and repetitions observed in dialogue tasks. Su et al. (2022) present a contrastive solution to encourage diversity while maintaining coherence in the generated text.

Various studies have addressed diversity in neural MT systems, which is a closely related topic. Sampling predictions from the output distribution can be an effective decoding strategy for back-translation, as described by Edunov et al. (2018), or sampling from less likely tokens Holtzman et al. (2020). Results show that such techniques enlarge diversity and richness of the generated translations when compared to data generated by beam or greedy search, but introduce semantic inconsistency in translations. In Lin et al. (2022) is proposed a multi-candidate optimization framework for augmenting diversity. The authors propose to guide an NMT model to learn more diverse translations from its candidate translations based on reinforcement learning. During training, the model generates multiple candidate translations, of which rewards are quantified according to their diversity and quality.

A different approach attempts to condition the decoding procedure with diverse signals. Typically, Shu et al. (2019) use syntactic codes to condition the translation process. Lachaux et al. (2020) replace the syntactic codes with latent domain variables derived from target sentences. Similarly, Schioppa et al. (2021) use prefix-based control tokens and vector-based interventions for controlling output translations from a NMT system. In the context of paraphrase generation Vahtola et al. (2023) propose a translation-based guided paraphrase generation model that learns useful features for promoting surface form variation in generated paraphrases.

## 3 Adjusting the ground-truth distribution

Throughout the training process, at every time-step $t$, neural machine translation networks generate predictions over the target-side vocabulary based on the input $x$ and previous predictions $y_{<t}$:

$$p_t^i = p(y_t^i | x, y_{<t}), \ i \in [1, ..., V]$$

where $V$ indicates the size of the target vocabulary.

The loss function evaluates the neural network's capacity to model the training data by comparing its predictions to a reference target vector $r = [r_1, r_2, ..., r_T]$, where $T$ denotes the sequence length. This loss is utilized to update the network's parameters, aiming to minimize the observed error in the model. The loss at time-step $t$ is usually computed as the cross-entropy between the model predictions $p_t = [p_t^1, ..., p_t^V]$ and the ground-truth distribution $q_t = [q_t^1, ..., q_t^V]$:

$$\ell_t = -\sum_{i=1}^{V} q_t^i \, log(p_t^i) \tag{1}$$

Note that the vector $q_t$ is a one-hot encoding representation of $r_t$, with all entries set to $0$ except for the token indicated by $r_t$, which is set to $1$. Addressing the over-fitting risk illustrated by the previous $q_t$ distribution, label smoothing Szegedy et al. (2015); Müller et al. (2019) (LS) is widely employed to achieve a smoother distribution:

$$q_t^{\epsilon LS} = (1 - \epsilon)q_t + \frac{\epsilon}{V} \tag{2}$$

with $\epsilon$ being a commonly small hyper-parameter.[1]

| t | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| r | I | like | cookies | and | cookies | . |

| | | | | | | |
|---|---|---|---|---|---|---|
| . | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 0 | 0 | 0 | 0 | 0 | 0 |
| and | 0 | 0 | 0 | 0 | 0 | 0 |
| like | 0 | 0 | 0 | 0 | 0 | 0 |
| cookies | 0 | 0 | 0 | 0 | 1 | 0 |

Figure 1: Matrix for the ground-truth $r =$'*I like cookies and cookies.*'. Rows $t$ and $r$ represent respectively the time-step and the corresponding ground-truth token. A reduced model vocabulary (matrix rows) is used to facilitate reading.

LS can be interpreted as penalizing the probability of the ground-truth class by a factor of $1 - \epsilon$, while evenly distributing the removed probability mass among all classes, $\epsilon/V$. Building upon a strategy akin to label smoothing, we make additional adjustments to the ground-truth distribution and reduce the likelihood of repeated tokens, with the

---

[1] $\epsilon = 0$ yields the initial distribution $q_t$, whereas $\epsilon = 1$ implies a uniform distribution.

Figure 2: Ground-truth distributions for the $5^{th}$ time-step of our example: the original one-hot encoding $q$; adjusted with label smoothing $q^{\epsilon LS}$; and further adjusted with repetitions $q^{\epsilon L\alpha}$.

goal of enabling the model to learn to predict repetitions with lower probability. We introduce a matrix, denoted as $_{V \times T}$, which indicates whether the ground-truth token $r_t$ is also present in the preceding time-steps.[2] Figure 1 illustrates an example of matrix with ground-truth *I like cookies and cookies.* as translation of the Japanese sentence クッキーとビスケットが好き with a model vocabulary of 5 tokens (matrix rows). Both Japanese terms クッキー *[cookies]* and ビスケット *[biscuits]* are correctly translated into English as *cookies*, yet this choice clearly reduces the fluency and clarity of the translation. As it can be seen, only $_{[i=5,t=5]}$ is set to 1 since only $r_5 = $ '*cookies*' occurs in a preceding time-step ($t = 3$).

We consequently update the ground-truth distribution following:

$$q_t^{\epsilon LS\alpha} = (1-\epsilon)(1-\alpha_t)\, q_t + \frac{\epsilon}{V} \qquad (3)$$

where $\alpha$ is a hyper-parameter, and $\alpha$ is used as a penalty, much like $\epsilon$ in the case of LS. Note that only the label smoothing probabilities discounted are distributed among all classes. As a result, time-steps with repeated tokens (such as $t = 5$ in our example) do not constitute proper probability distributions, as their sum does not add to 1. Figure 2 illustrates ground-truth distributions for our example at time-step $t = 5$: the original one-hot encoding $q$; the original distribution adjusted using label smoothing $q^{\epsilon LS}$, and further adjusted using repetitions $q^{\epsilon LS\alpha}$.[3] A significant challenge with the aforementioned techniques that modify $q$ distribution with repetitions is their limited impact on the training process, primarily caused by the scarcity

of repeated tokens in datasets. In the following section, we present alternative approaches to address this challenge.

## 4 Gathering Examples with Repetitions

As previously depicted, our intention is to instruct the model to minimize certain repetitions while preserving others deemed necessary for an accurate translation. To achieve this, we must compile a relatively large dataset of examples that demonstrate this behavior to the model. We initially focus on repetitions of content words such as *nouns*, *adjectives*, *verbs*, and *adverbs*. Function words, which serve a distinct grammatical role in a sentence, are excluded from this analysis. Current MT networks reliably generate these words based on their understanding of grammatical correctness.

We back-translate the Japanese side of the JiJi corpus (further detailed in Section 5.1) into English and annotate word (or sequence) repetitions of content words based on automatic morpho-syntactic annotations performed by Spacy[4]. We employ word-alignments between Japanese and English words performed by the Giza++ Och and Ney (2003) toolkit[5] in order to consider only repetitions of English content words aligned to Japanese content words (Verbs, Nouns, Adjectives and Adverbs). The resulting set of examples with repetitions from src/tgt training pairs will be regarded as instances that the model needs to learn to discourage. Consequently, we utilize them for training after annotating the repeated target words in their respective matrices.

It's worth noting that the presented approach does not require any alterations to the network architecture and maintains the same training and inference efficiency.

---

[2]Note that repetitions are computed over words while matrix refers to tokens $r \in V$ for each time-step $t \in T$.

[3]As previously discussed, distribution $q^{\epsilon LS\alpha}$ does not form a proper distribution since probabilities do not add to 1 ($0,02 + 0,02 + 0,02 + 0,02 + 0,0092 = 0,0892$). We leave for future experiments the normalization of the output scores in order to allow for a valid probability distribution.

[4]https://spacy.io/

[5]https://github.com/moses-smt/giza-pp.

| | | |
|---|---|---|
| Jap | JEMAの担当者は白物家電について、「機能を絞った低価格製品、高価格な高機能製品とも好調だ」と述べている。 | |
| Eng | "Shipments have been robust for both low-priced models with reduced functions and expensive high-spec products," a JEMA official said. | |
| Jap | JEMAの担当者は白物家電について、「<target id=0 ref=0 type=s>機能<\target>を絞った低価格<target id=1 ref=0 type=s>製品<\target>、高価格な高<target id=0 ref=1 type=s>機能<><target id=1 ref=1 type=s>製品<\target>とも好調だ」と述べている。 | |
| Eng | "Shipments have been robust for both low-priced <target id=1 ref=0 type=s>models<\target> with reduced <target id=0 ref=0 type=s>functions<\target> and expensive <target id=0 ref=1 type=s>high-spec<\target> <target id=1 ref=1 type=s>products<\target>," a JEMA official said. | |

Table 1: example of Japanese-English translation: raw translation is shown at the top, and the tagged translation to annotate repetitions is shown at the bottom.

## 5 Experimental Framework

### 5.1 Datasets

We evaluate the proposed methods in a Japanese-to-English translation task. Thus, we utilize Japanese-English parallel corpora freely obtained from the WMT24 for Non-Repetitive Translation Task website[6]. The corpus is compiled by Jiji Press Ltd in collaboration with the National Institute of Information and Communication Technology (NICT) with various categories, including politics, economy, nation, business, markets, sports, etc., for use in machine translation, in particular for previous the Workshop on Asian Translation (WAT)[7].

Table 2 presents various statistics of the corpora used in this work, including the total number of sentences, vocabularies, words, and average sentence length. Statistics are computed after performing a light tokenization aiming to split-off punctuation. For testing, we use the supplied Japanese-English datasets made available by the task organizers.

| Lang | #Sents | #Vocab | Words | Length |
|---|---|---|---|---|
| *Training-set* | | | | |
| Jap | 200k | $49K$ | 6.9M | 4.46 |
| Eng | | 118K | 4.5M | 24.64 |
| *Repetition-set* | | | | |
| Jap | 470 | $3,297$ | $23,472$ | 4.22 |
| Eng | | $4,341$ | $13,814$ | 11.91 |

Table 2: Corpora statistics. M and K stand for millions and thousands respectively.

Due to the poor alignment quality of the Japanese-English parallel sentences present in the

provided dataset (sentence pairs are coupled using an automatic cross-lingual sentence similarity score) we decided to back-translate the English side using an in-house English-Japanese model. Then, using the resulting Japanese[8]-English dataset we fine-tune our baseline Japanese-English model.

In addition, we use a test set of repetitions also provided by the challenge, consisting of reference English machine translations and their corresponding Japanese machine translations that include at least one word repeated on the target (English) side for a more nuanced analysis of repetition. Among the files corresponding to the test datasets are those containing tagged files in which repeated words and their translations in each sentence pair are marked with tags *<target>* and *</target>*. Marked words indicate that they are evaluated repetitions. Three labels, 'id', 'ref' and 'type' are embedded within the tags. Table 1 illustrates an example, where:

**id** indicates IDs of repated words. In the above example, two tagged repeated words are included, i.e., 機能 (id=0) and 製品 (id=1). The number of instances including multiple id's, such as the above example, are limited.

**ref** indicates IDs of pairs of source/target words, such as 製品/*models* (id=1, ref=0) and 製品/*products* (id=1, ref=1).

**type** indicates whether they are substituted (*s*) or reduced (*r*).

The *Repetition-set* is mainly used to evaluate the performance of our models in handling repetition problems, as well as to assess overall translation accuracy.

## 5.2 NMT Models

Our NMT model is built using an in-house implementation of the state-of-the-art Transformer architecture Vaswani et al. (2017). Details of the network hyper-parameters emplooy for training are given in Table 3.

| | |
|---|---|
| size of word embedding | 512 |
| size of hidden layers | 512 |
| size of inner feed forward layer | 2,048 |
| number of heads | 8 |
| number of layers | 6 |
| batch size | 4,000 (tokens) |
| batch accumulation | 25 (batches) |

Table 3: Network hyperparameters.

For optimization work we use the lazy Adam algorithm Kingma and Ba (2014). We set warmup steps to 4,000 and update learning rate for every 8 iterations. All models are trained using a single NVIDIA V100 GPU.

We limit the source and target sentence lengths to 150 tokens based on BPE Sennrich et al. (2016) preprocessing. A total of $28K$ BPE merge operations are separately computed for each language. We finally use a joint Japanese and English vocabulary of $58K$ tokens. In inference we use a beam size of 5.

Our *baseline* English-to-Japanese model is trained during more than 3 million iterations using all the parallel data available in the Opus website[9].

## 6 Results

To evaluate the method presented in this paper we consider the previous *baseline* model that we update with 15K additional iterations for two different configurations of the ground-truth distribution:

$q^{\epsilon LS}$ follows the same configuration than our *baseline* model with label smoothing set to $\epsilon = 0.1$.

$q^{\epsilon LS\alpha}$ further penalizes the ground-truth distribution with repetition penalties as detailed in Section 3 with $\epsilon = 0.1$ and for different values of $\alpha$.

Note that for both configurations, we use the same training corpus detailed in Table 2 (*Training-set*).

We also assess the effectiveness of two large language models (LLM) with translation capabilities to overcome the repetition issue:

$GPT3.5$ consists of the *GPT3.5-turbo* version of the OpenAI LLM. Built upon the Generative Pre-trained Transformer architecture Radford and Sutskever (2018) which employs only a transformer decoder. Following an auto-regressive approach, the model ensures that the generated text maintains coherence and relevance to the context provided by the input text. Translations are conducted using the OpenAI API, while emphasizing the importance of minimizing word repetitions through the provided prompt: *Translate the following text from English to Japanese, ensuring that the translated output maintains coherence and fluency while minimizing the repetition of words or phrases. Pay attention to using synonyms, varied sentence structures, and appropriate linguistic devices to enhance the overall quality of the translation. Feel free to creatively adapt the language to achieve a natural and engaging tone in the target language. I want you to only reply the translation, do not write explanations.*

$NLLB$ is a family of machine translation models based on the Transformer encoder-decoder architecture, enabling translation between any of the 202 language varieties NLLB Team et al. (2022). We use the *nllb-200-distilled-600M*[10] version and perform translations with the efficient CTranslate2[11] inference toolkit.

To evaluate the presented methods, we report BLEU results computed by `sacrebleu`[12] Post (2018) respectively over test sets. We also report the number of word repetitions that hinder fluency, *Degrading*, after a human evaluation performed on translation hypotheses. Table 4 summarize results obtained by different system configurations.

Models fine-tuned from the *baseline* network exhibit nearly identical quality scores across the *test* set. This suggests that training with the method presented to adjust the ground-truth distribution does not compromise translation quality. On the contrary, unlike Configuration $q^{\epsilon LS}$, Configurations $q^{\epsilon LS\alpha}$ demonstrate a significant decrease in the number of repetitions that degrade fluency over the *Repetition-set*, while retaining most of the acceptable repetitions in the translated output.

---

[9] https://opus.nlpl.eu/

[10] https://huggingface.co/facebook/nllb-200-distilled-600M
[11] https://github.com/OpenNMT/CTranslate2
[12] https://github.com/mjpost/sacrebleu

Results from both LLMs demonstrate a reduced number of repetitions, suggesting an elevated level of diversity and fluency of such models. However, the translation quality scores of LLMs do not align with those achieved by the models presented in this study in either of the test sets, especially translations obtained by GPT-3.5. These findings are consistent with those presented by Bawden and Yvon (2023) where the authors note the challenge of controlling translations performed by BLOOM[13], a multilingual LLM.

| Configuration | BLEU | *Degrading* |
|---|---|---|
| $q^{\epsilon LS}$ | 28.41 | 77 |
| $q^{\epsilon LS\alpha}, 1 - \alpha = 10^{-6}$ | 28.91 | 60 |
| $GPT3.5$ | 19.29 | 64 |
| $NLLB$ | 16.12 | 74 |

Table 4: Translation accuracy results and number of repetitions present in translations performed by models under different configurations. $\epsilon$ is always set to $0.1$.

## 7 Conclusions and Further Work

We presented SYSTRAN submission to the WMT24 Non-Repetitive Translation Task. Our NMT systems introduce a method to reduce the occurrence of repetitions in translation hypotheses, which significantly affects the readability of the generated texts. The method is solely implemented during fine-tuning at the conclusion of the training phase, without any modifications to the inference process. Experiments indicate the ability of our proposed methods in reducing the repetition problem.

We aim to further study the impact of the ratio between the number of reference sentences and synthetic translations that include repetitions during the training process. Additionally, we plan to analyze the influence of the distance (measured in number of words) between repetitions and explore the possibility of replacing the binary penalty in matrix with a softer approach.

## Acknowledgements

---

[13]https://huggingface.co/bigscience/bloom

## References

Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of BLOOM. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 157–170, Tampere, Finland. European Association for Machine Translation.

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 489–500, Brussels, Belgium. Association for Computational Linguistics.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Marie-Anne Lachaux, Armand Joulin, and Guillaume Lample. 2020. Target conditioning for one-to-many generation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2853–2862, Online. Association for Computational Linguistics.

Margaret Li, Stephen Roller, Ilia Kulikov, Sean Welleck, Y-Lan Boureau, Kyunghyun Cho, and Jason Weston. 2020. Don't say that! making inconsistent dialogue unlikely with unlikelihood training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4715–4728, Online. Association for Computational Linguistics.

Huan Lin, Baosong Yang, Liang Yao, Dayiheng Liu, Haibo Zhang, Jun Xie, Min Zhang, and Jinsong Su. 2022. Bridging the gap between training and inference: Multi-candidate optimization for diverse neural machine translation. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2622–2632, Seattle, United States. Association for Computational Linguistics.

Rafael Müller, Simon Kornblith, and Geoffrey Hinton. 2019. *When Does Label Smoothing Help?* Curran Associates Inc., Red Hook, NY, USA.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp

Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Narasimhan Karthic Salimans Tim Radford, Alec and Ilya Sutskever. 2018. Improving language understanding with unsupervised learning. *Technical Report*.

Andrea Schioppa, David Vilar, Artem Sokolov, and Katja Filippova. 2021. Controlling machine translation for multiple attributes with additive interventions. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6676–6696, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Raphael Shu, Hideki Nakayama, and Kyunghyun Cho. 2019. Generating diverse translations with sentence codes. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1823–1827, Florence, Italy. Association for Computational Linguistics.

Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. 2022. A contrastive framework for neural text generation. In *Advances in Neural Information Processing Systems*.

Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2015. Rethinking the inception architecture for computer vision. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826.

Teemu Vahtola, Mathias Creutz, and Jrg Tiedemann. 2023. Guiding zero-shot paraphrase generation with fine-grained control tokens. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)*, pages 323–337, Toronto, Canada. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2019. Neural text generation with unlikelihood training. *ArXiv*, abs/1908.04319.

# Mitigating Metric Bias in Minimum Bayes Risk Decoding

**Geza Kovacs** and **Dan Deutsch** and **Markus Freitag**
Google
{geza, dandeutsch, freitag}@google.com

## Abstract

While Minimum Bayes Risk (MBR) decoding using metrics such as COMET or MetricX has outperformed traditional decoding methods such as greedy or beam search, it introduces a challenge we refer to as *metric bias*. As MBR decoding aims to produce translations that score highly according to a specific utility metric, this very process makes it impossible to use the same metric for both decoding and evaluation, as improvements might simply be due to *reward hacking* rather than reflecting real quality improvements. In this work we find that compared to human ratings, neural metrics not only overestimate the quality of MBR decoding when the same metric is used as the utility metric, but they also overestimate the quality of MBR/QE decoding with other neural utility metrics as well. We also show that the metric bias issue can be mitigated by using an ensemble of utility metrics during MBR decoding: human evaluations show that MBR decoding using an ensemble of utility metrics outperforms a single utility metric.

## 1 Introduction

Minimum bayes risk (MBR) decoding is a decoding approach where $n$ candidate translations are sampled from the MT system, and they are used as pseudoreferences for a reference-based *utility metric*. MBR decoding computes the utility metric for all $O(n^2)$ pairs of candidates and pseudoreferences, selecting the candidate that achieves the best average score across all pseudoreferences. Quality Estimation (QE) decoding[1] selects the candidate that scores best according to a QE utility metric. Previous work on MBR decoding has shown that it results in improvements on the utility metric (Amrhein and Sennrich, 2022; Cheng and Vlachos, 2023; Eikema and Aziz, 2022), however other metrics do not improve as much as the utility metric (Guttmann et al., 2024; Vamvas and Sennrich,

2024). This issue of MBR/QE decoding exhibiting bias towards the utility metric complicates our ability to use automatic metrics to compare the quality of MBR/QE-based MT systems, as we cannot tell whether improvements in automatic metrics from MBR/QE decoding correspond to actual improvements in quality, or if it simply reward hacking. Prior work has assumed that this issue can be avoided by using a different metric for evaluating MBR decoding outputs (Tomani et al., 2023), though this assumption has never been tested.

In this work we compare the results of human vs metric-based evaluation of MBR/QE decoding with a wide variety of metrics to show that the quality of MBR/QE decoding is overestimated by not only the utility metric, but also other similar metrics. While MBR/QE decoding with a single utility metric results in significant gains in automatic metrics, it does not perform better than greedy decoding in our human evaluations. This may be due to MBR decoding preferring fluent yet inaccurate candidates. Using an ensemble of metrics as the utility helps us mitigate the metric bias issue, with human evaluations showing that MBR decoding with an ensemble utility metric results in significantly better translations than greedy decoding or MBR/QE decoding with a single utility metric.

In this paper we contribute:

1. A large-scale analysis of metric bias in MBR and QE decoding with metrics commonly used in MT, showing that this metric bias issue holds across many different metrics and language pairs, and is not resolved by simply using a different metric for evaluation.
2. Mitigation strategies for MBR bias using QE filtering followed by MBR decoding, as well as MBR decoding using an ensemble of metrics as the utility function.
3. A human evaluation showing that MBR decoding with ensembles outperforms MBR decoding with a single metric.

---

[1]Also known as QE reranking or QE filtering.

## 2 Related Work

Cheng and Vlachos (2023); Eikema and Aziz (2022); Guttmann et al. (2024) find that MBR decoding improves automated metrics on various high, medium, and low resource language pairs. Freitag et al. (2023a, 2022); Tomani et al. (2023) find that human raters prefer the outputs of MBR/QE decoding over greedy decoding.

MBR variants achieve speedups via heuristics (Trabelsi et al., 2024; Jinnai and Ariu, 2024), filtering pseudoreferences via a QE metric (Deguchi et al., 2024, 2023) or filtering via another reference-based metric (Vamvas and Sennrich, 2024; Eikema and Aziz, 2022). Quality-aware translation, which incorporates quality estimation into the training process, has been found to improve translation quality over standard MBR (Tomani et al., 2023).

Other techniques for aligning translation models with human preferences include direct preference optimization (Rafailov et al., 2024; Yang et al., 2024), reinforcement learning from human feedback (Christiano et al., 2017), and reinforcement learning from AI feedback (Bai et al., 2022).

Guttmann et al. (2024); Vamvas and Sennrich (2024) show evidence of metric bias in MBR decoding, as they find that neural evaluation metrics favor models using MBR on the metric used as the utility function. However, these papers only cover only 2 metrics, and neither have human evaluations.

Sellam et al. (2020b); Freitag et al. (2023b); Glushkova et al. (2023) find that ensembling metrics can improve their ability to detect critical errors and improve agreement with human preferences, though they do not investigate the effects of ensembling utility metrics on MBR decoding.

Reward hacking (Skalse et al., 2022) is an issue in reinforcement learning where the reward function improves but the system's behavior is not aligned with human preferences. The metric bias problem in MBR decoding can be viewed as an instance of reward hacking, as the utility function improves while not necessarily improving quality.

## 3 Study 1: Metric Bias in MBR Decoding

### 3.1 Methodology

To investigate metric bias in MBR/QE decoding, we perform MBR/QE decoding via various utility metrics and compare how they perform on various evaluation metrics. We investigate MBR decoding using these reference-based utility metrics:

1. MetricX-23 (Juraska et al., 2023)
2. XCOMET-XXL (Guerreiro et al., 2023)
3. XCOMET-XL (Guerreiro et al., 2023)
4. COMET22 (Rei et al., 2022a)
5. AfriCOMET (Wang et al., 2024)
6. IndicCOMET (Sai B et al., 2023)
7. BLEURT (Sellam et al., 2020a)
8. YiSi-1 (Lo, 2019)
9. sentBLEU (Papineni et al., 2002)
10. chrF (Popović, 2015)
11. chrF++ (Popović, 2017)
12. TER (Snover et al., 2006)

We also investigate QE decoding (Fernandes et al., 2022) using the following QE metrics:

1. MetricX-QE (Juraska et al., 2023)
2. CometKiwi23-XXL (Rei et al., 2023)
3. CometKiwi23-XL (Rei et al., 2023)
4. CometKiwi22 (Rei et al., 2022b)
5. AfriCOMET-QE (Wang et al., 2024)

We used a dev set for selecting ensembles, and a test set for reporting final results and human evaluation. The dev datasets and language pairs are:

1. FLORES-200 dev set (Costa-jussà et al., 2022): English-Swahili (en-sw), Igbo (en-ig), Hindi (en-hi), Tamil (en-ta), Somali (en-so), Hausa (en-ha), Malayalam (en-ml), Gujarati (en-gu), Hungarian (en-hu), Vietnamese (en-vi)
2. WMT2022 (Kocmi et al., 2022): English-Chinese (en-zh), Chinese-English (zh-en), English-German (en-de), German-English (de-en)

The test set datasets and language pairs are:

1. FLORES-200 test set: en-sw, en-ig, en-hi, en-ta, en-so, en-ha, en-ml, en-gu, en-hu, en-vi
2. WMT2023 (Kocmi et al., 2023): en-zh, zh-en[2], en-de, de-en

We produced translations using Gemini 1.0 Pro (Gemini Team Google, 2023) with prompts including 5-shot examples. We used epsilon sampling as recommended by Freitag et al. (2023a) with a sample size of 128. See Appendix A for prompts used for generating translations and instructions for computing scores from metrics.

---

[2]Due to errors in the WMT2023 zh-en reference translations, we use the references from Liu et al. (2024) for zh-en.

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.58 | 1.16 | 82.3 | 77.8 | 76.8 | 68.2 | 77.5 | 85.2 | 77.3 | 84.3 | 57.2 | 54.2 | 26.4 | 63.4 |
| MetricX | 0.656‡ | 0.557‡ | 85.5‡ | 79.6‡ | 79.0‡ | 69.4‡ | 77.7‡ | 84.9‡ | 76.6‡ | 81.2‡ | 50.3‡ | 46.9‡ | 18.1‡ | 75.7‡ |
| MetricX-QE | 0.899‡ | 0.349‡ | 84.4‡ | 78.2‡ | 78.3‡ | 68.8‡ | 77.6‡ | 84.4‡ | 75.6‡ | 81.1‡ | 49.3‡ | 45.9‡ | 17.6‡ | 75.3‡ |
| XCOMET-XXL | 1.25‡ | 0.868‡ | 89.9‡ | 80.4‡ | 80.8‡ | 69.9‡ | 78.1‡ | 85.0‡ | 76.6‡ | 81.5‡ | 50.4‡ | 47.0‡ | 18.5‡ | 73.6‡ |
| XCOMET-XL | 1.38‡ | 1.00‡ | 86.4‡ | 85.0‡ | 80.2‡ | 71.5‡ | 78.7‡ | 85.3‡ | 77.6‡ | 82.2‡ | 51.9‡ | 48.7‡ | 20.1‡ | 71.5‡ |
| CometKiwi23-XXL | 1.43‡ | 0.940‡ | 86.6‡ | 80.4‡ | 85.5‡ | 71.4‡ | 78.7‡ | 85.2 | 76.7‡ | 82.2‡ | 51.7‡ | 48.4‡ | 19.9‡ | 71.7‡ |
| CometKiwi23-XL | 1.46‡ | 0.978‡ | 85.0‡ | 81.5‡ | 81.3‡ | 74.8‡ | 78.8‡ | 85.2 | 76.8‡ | 82.1‡ | 51.7‡ | 48.4‡ | 19.8‡ | 72.6‡ |
| CometKiwi22 | 1.57‡ | 1.07‡ | 84.0‡ | 79.6‡ | 79.7‡ | 70.5‡ | 81.9‡ | 85.4‡ | 76.8‡ | 82.3‡ | 51.9‡ | 48.6‡ | 20.1‡ | 71.0‡ |
| COMET22 | 1.40‡ | 1.02‡ | 84.7‡ | 80.0‡ | 79.3‡ | 70.0‡ | 78.7‡ | 87.4‡ | 78.1‡ | 83.5‡ | 55.3‡ | 52.0‡ | 23.2‡ | 67.0‡ |
| BLEURT | 1.35‡ | 0.986‡ | 83.8‡ | 79.1‡ | 78.6‡ | 69.4‡ | 78.1‡ | 85.5‡ | 82.3‡ | 82.6‡ | 53.2‡ | 49.8‡ | 21.0‡ | 71.3‡ |
| YiSi | 1.57 | 1.14† | 82.6‡ | 78.0* | 77.3‡ | 68.7‡ | 77.7‡ | 85.6‡ | 77.7‡ | 85.0‡ | 57.7‡ | 54.5‡ | 26.1* | 62.6 |
| chrF | 1.54‡ | 1.13† | 82.6‡ | 78.0* | 77.6‡ | 68.9‡ | 77.7‡ | 85.7‡ | 77.8‡ | 84.5‡ | 58.6‡ | 55.3‡ | 25.8‡ | 65.1‡ |
| chrF++ | 1.54‡ | 1.13† | 82.6‡ | 78.0† | 77.5‡ | 68.9‡ | 77.7‡ | 85.6‡ | 77.9‡ | 84.6‡ | 58.6‡ | 55.4‡ | 26.2 | 64.6† |
| sentBLEU | 1.61 | 1.18* | 82.2* | 77.8* | 76.8 | 68.2 | 77.5 | 85.2 | 77.3* | 84.3 | 57.0* | 54.1* | 27.1‡ | 62.3 |
| TER | 1.74‡ | 1.27‡ | 81.9‡ | 77.2‡ | 75.9‡ | 67.5‡ | 77.2‡ | 84.7‡ | 76.7‡ | 83.9‡ | 55.7‡ | 52.7‡ | 25.6‡ | 59.7‡ |
| | | | | | | | | | | | | | | |
| rankAvg:all | 1.08‡ | 0.739‡ | 86.5‡ | 81.7‡ | 81.2‡ | 71.4‡ | 79.3‡ | 86.5‡ | 79.3‡ | 84.3 | 57.1 | 53.9 | 25.3‡ | 63.7 |
| rankAvg:qe | 1.04‡ | 0.580‡ | 86.6‡ | 81.8‡ | 83.2‡ | 73.0‡ | 80.3‡ | 85.9‡ | 77.7‡ | 82.6‡ | 52.8‡ | 49.5‡ | 20.8‡ | 70.7‡ |
| rankAvg:top | 0.899‡ | 0.566‡ | 88.2‡ | 83.0‡ | 83.0‡ | 72.7‡ | 78.9‡ | 85.8‡ | 78.1‡ | 82.5‡ | 52.8‡ | 49.5‡ | 20.7‡ | 71.0‡ |
| rankAvg:topQe | 1.00‡ | 0.527‡ | 86.8‡ | 81.7‡ | 83.7‡ | 73.3‡ | 78.9‡ | 85.6‡ | 77.5 | 82.4‡ | 52.3‡ | 48.9‡ | 20.2‡ | 71.7‡ |
| rankAvg:mxmxqe | 0.700‡ | 0.417‡ | 85.6‡ | 79.7‡ | 79.2‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.4‡ | 47.0‡ | 18.2‡ | 75.1‡ |
| rankAvg:noLex | 0.993‡ | 0.657‡ | 87.3‡ | 82.4‡ | 82.0‡ | 72.0‡ | 79.6‡ | 86.6‡ | 79.5‡ | 83.8‡ | 55.6‡ | 52.3‡ | 23.4‡ | 66.7‡ |
| rankAvg:noNC | 1.09‡ | 0.734‡ | 85.2‡ | 80.4‡ | 79.5‡ | 70.1‡ | 78.5‡ | 86.4‡ | 79.2‡ | 84.4‡ | 57.4‡ | 54.1* | 25.7‡ | 63.0* |
| rankAvg:noNCnoLex | 0.968‡ | 0.636‡ | 85.8‡ | 80.8‡ | 80.0‡ | 70.4‡ | 78.6‡ | 86.6‡ | 79.7‡ | 84.0‡ | 56.1‡ | 52.8‡ | 24.0‡ | 66.0‡ |
| allQE(32)allMBR | 1.06‡ | 0.733‡ | 86.7‡ | 81.9‡ | 81.3‡ | 71.4‡ | 79.2‡ | 86.5‡ | 79.2‡ | 84.1‡ | 56.6‡ | 53.4‡ | 24.9‡ | 64.5 |
| allQE(32)nolexMBR | 0.978‡ | 0.680‡ | 87.5‡ | 82.6‡ | 81.6‡ | 71.7‡ | 79.2‡ | 86.6‡ | 79.5‡ | 83.7‡ | 55.6‡ | 52.3‡ | 23.6‡ | 66.6‡ |
| topQE(32)topMBR | 0.861‡ | 0.599‡ | 88.4‡ | 83.3‡ | 82.0‡ | 71.9‡ | 78.8‡ | 85.7‡ | 78.1‡ | 82.4‡ | 52.7‡ | 49.4‡ | 20.7‡ | 70.9‡ |
| noncQE(32)noncMBR | 0.992‡ | 0.629‡ | 85.6‡ | 80.6‡ | 79.8‡ | 70.2‡ | 78.5‡ | 86.3‡ | 78.9‡ | 83.9‡ | 56.1‡ | 52.8‡ | 24.2‡ | 65.2‡ |
| noncQE(32)noncnolexMBR | 0.911‡ | 0.596‡ | 86.0‡ | 81.0‡ | 80.1‡ | 70.4‡ | 78.7‡ | 86.5‡ | 79.4‡ | 83.6‡ | 55.1‡ | 51.7‡ | 22.9‡ | 67.5‡ |
| mxQE(32)mxMBR | 0.662‡ | 0.475‡ | 85.6‡ | 79.8‡ | 79.2‡ | 69.5‡ | 77.8‡ | 85.0‡ | 76.8‡ | 81.5‡ | 50.7‡ | 47.3‡ | 18.5‡ | 74.9‡ |
| ckQE(32)xcMBR | 1.24‡ | 0.847‡ | 89.6‡ | 80.8‡ | 82.8‡ | 70.7‡ | 78.4‡ | 85.2 | 77.0‡ | 81.9‡ | 51.3‡ | 48.0‡ | 19.5‡ | 72.2‡ |
| mxQE(32)xcMBR | 1.03‡ | 0.593‡ | 89.5‡ | 80.6‡ | 80.9‡ | 70.1‡ | 78.2‡ | 85.1 | 76.9‡ | 81.7‡ | 50.7‡ | 47.4‡ | 18.8‡ | 73.1‡ |
| ckQE(32)mxMBR | 0.728‡ | 0.557‡ | 86.5‡ | 80.6‡ | 82.2‡ | 70.7‡ | 78.3‡ | 85.4‡ | 77.3 | 81.9‡ | 51.7‡ | 48.3‡ | 19.5‡ | 73.3‡ |

Table 1: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), averaged across all languages (test datasets). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

| | MetricX-QE | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | AfriCOMET-QE (African only) | MetricX | XCOMET-XXL | XCOMET-XL | COMET22 | AfriCOMET (African only) | IndicCOMET (Indic only) | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| all | ■ | ■ | ■ | ■ | | ■ | ■ | ■ | ■ | | | ■ | ■ | ■ | ■ | ■ | ■ |
| qe | ■ | ■ | ■ | ■ | | | | | | | | | | | | | |
| top | ■ | | ■ | | | ■ | ■ | ■ | | | | | | | | | |
| topQe | ■ | | ■ | | | | | | | | | | | | | | |
| mxmxqe | ■ | | | | | ■ | | | | | | | | | | | |
| noLex | ■ | ■ | ■ | ■ | | ■ | ■ | ■ | ■ | | | ■ | ■ | | | | ■ |
| noNC | ■ | ■ | ■ | ■ | | ■ | ■ | ■ | ■ | | | ■ | ■ | ■ | ■ | ■ | ■ |
| noNCnoLex | ■ | ■ | ■ | ■ | | ■ | ■ | ■ | ■ | | | ■ | ■ | | | | ■ |
| noNCQe | ■ | ■ | ■ | ■ | | | | | | | | | | | | | |
| allQE(N)allMBR | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| allQE(N)nolexMBR | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | | | | |
| topQE(N)topMBR | 1 | 1 | 1 | | | 2 | 2 | 2 | | | | | | | | | |
| noncQE(N)noncMBR | 1 | | | | 1 | 2 | | | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| noncQE(N)noncnolexMBR | 1 | | | | 1 | 2 | | | 2 | 2 | 2 | 2 | 2 | | | | |
| mxQE(N)xcMBR | 1 | | | | | | 2 | | | | | | | | | | |
| ckQE(N)xcMBR | | 1 | | | | | 2 | | | | | | | | | | |
| mxQE(N)mxMBR | 1 | | | | | 2 | | | | | | | | | | | |
| ckQE(N)mxMBR | | 1 | | | | 2 | | | | | | | | | | | |

Table 2: Metrics included in each ensemble. Rows are ensembles, columns are metrics. Black cells indicate that the metric is included in a single-step ensemble. Green cells indicate the metric is used for the 1st step (QE filtering) in a 2-step ensemble. Red cells indicate the metric is used for the 2nd step (MBR decoding) in a 2-step ensemble.

## 3.2 Results

Results are shown in Table 1 for average scores across all language pairs on the test datasets. We observe that for all reference-based metrics, the best-performing system is MBR decoding using the same utility metric. This result also holds for all QE metrics, but that is by definition, because QE decoding picks the sample with the best QE score. These results also hold on individual languages and the dev set (Appendix G and E).

We can also see that MBR decoding outputs for utility metrics which are similar to the evaluation metric tend to score better than when the MBR utility metric is dissimilar to the evaluation metric. For example, MBR/QE decoding with neural metrics (MetricX and COMET families) performs better than greedy when evaluated with other neural metrics, but worse than greedy if evaluated via lexical metrics. Likewise, MBR decoding with lexical metrics (sentBLEU, chrF, chrF++, and TER) and semantic metrics (YiSi) perform highly when evaluated by lexical and semantic metrics, but poorly when evaluated via neural metrics. The pattern also holds for similar metrics within the same family – XCOMET-XXL prefers MBR/QE decoding using CometKiwi23-XXL and XCOMET-XL, and MetricX prefers outputs from MetricX-QE.

These results suggest the existence of metric bias in MBR decoding – that is, they suggest that MBR decoding will result in a disproportionately large improvement in the utility metric and metrics similar to the utility metric, relative to the actual improvement in quality. In order to address this issue, in the next section we will investigate ensembling metrics during MBR decoding as a means of avoiding overfitting to a particular utility metric.

## 4 Study 2: MBR Decoding using Ensembles of Metrics

### 4.1 Methodology

As a mitigation strategy for utility metric bias in MBR decoding, we investigate how using an ensemble of metrics performs for MBR decoding. We explore the following ensembling techniques (see Appendix C for pseudocode for these techniques):

1. rankAvg: For each metric, assigns a rank to each of the 128 samples (where 0 is best and 127 is worst). Select the sample where the average rank across metrics is minimized.
2. rankMed: Select the sample where the median

rank across metrics is minimized.
3. rankMax: Select the sample where the maximum rank across metrics is minimized.
4. rank75q: Select the sample where the 0.75th quartile rank across metrics is minimized.

For each of these ensembling techniques, we compute ensembles with the following groups of metrics (see Table 2 and Appendix B for the complete list of metrics included in each ensemble):

1. all: All metrics
2. qe: All QE metrics
3. top: Top-performing metrics in WMT2023 metrics shared task (Freitag et al., 2023b)
4. topQe: Top-performing QE metrics
5. mxmxqe: MetricX + MetricX-QE ensemble
6. noLex: Non-lexical metrics
7. noNC: Metrics that permit commercial use
8. noNCnoLex: Non-lexical metrics that permit commercial use
9. noNCQe: QE metrics that permit commercial use

In addition to the ensembles above, we also investigate QE filtering followed by MBR decoding (QE filtering selects the top N candidates according to a QE metric, where N can be either 4, 8, 16, 32, 64). This two-step approach is faster than standard MBR decoding, as QE filtering is linear-time whereas MBR decoding is quadratic time. We include the following two-step ensembles:

1. allQE(N)allMBR: QE filter with all QE metrics, then MBR decode with all reference-based metrics
2. allQE(N)nolexMBR: QE filter with all QE metrics, then MBR decode with non-lexical reference-based metrics
3. topQE(N)topMBR: QE filter with top QE metrics, then MBR decode with top reference-based metrics
4. noncQE(N)noncMBR: QE filter with QE metrics that permit commercial use, then MBR decode with reference-based metrics that permit commercial use
5. noncQE(N)noncnolexMBR: QE filter with QE metrics that permit commerical use, then MBR decode with non-lexical reference-based metrics that permit commercial use
6. mxQE(N)xcMBR: QE filter with MetricX-QE, then MBR decode with XCOMET-XXL
7. ckQE(N)xcMBR: QE filter with CometKiwi23-XXL, then MBR decode

with XCOMET-XXL

8. mxQE(N)mxMBR: QE filter with MetricX-QE, then MBR decode with MetricX

9. ckQE(N)mxMBR: QE filter with CometKiwi23-XXL, then MBR decode with MetricX

The metrics included in each ensemble is shown in Table 2 and Appendix B.

## 4.2 Results

Results for a subset of ensembles averaged across all language pairs on the test sets are at Table 1 with additional ensembles shown in Appendix F. Results on the dev sets are shown in Appendix E. Breakdowns per language pair can be found in Appendix G. As expected, ensembles tend to perform better if judged by metrics that are better represented in the ensemble; for example, if judging by MetricX, the best ensembles are mxQE(32)mxMBR and rankAvg:mxmxqe, both of which are ensembles consisting of MetricX and MetricX-QE.

That said, observe that compared to MBR/QE decoding with a single utility metric, ensembles often improve on automated evaluations even according to metrics not included in the ensemble. For example, if we use the XCOMET or CometKiwi families of metrics to evaluate rankAvg:noNCnoLex and noncQE(32)noncnolexMBR (which do not include any metrics from the XCOMET or CometKiwi families), they outperform MBR/QE decoding with any single metric outside the XCOMET or CometKiwi families. Similarly, if lexical metrics are used to evaluate the rankAvg:noLex and allQE(32)nolexMBR ensembles, which do not include any lexical metrics, they still outperform MBR/QE decoding with any single neural metric. This suggests that ensembles help reduce metric bias towards a single metric, which results in improved automated evaluation scores according to other metrics not included in the ensemble.

## 5 Study 3: Human Evaluation

### 5.1 Methodology

For the human evaluation, we chose the following baselines and ensembles to evaluate:

1. Greedy decoding
2. Reference translation
3. MetricX (MBR decoding)
4. MetricX-QE (QE decoding)

5. AfriCOMET for African languages (MBR decoding)

6. AfriCOMET-QE for African languages (QE decoding)

7. IndicCOMET for Indic langauges (MBR decoding)

8. rankAvg:noNC (single-step ensemble)

9. rankAvg:noNCnoLex (single-step ensemble)

10. mxQE(32)mxMBR (multi-step ensemble)

11. noncQE(32)noncnolexMBR (multi-step ensemble)

We evaluated the following conditions only on en-de and zh-en due to budget constraints:

1. XCOMET-XXL (MBR decoding)
2. CometKiwi23-XXL (QE decoding)
3. COMET22 (MBR decoding)
4. rankAvg:all (single-step ensemble)

We chose MetricX, MetricX-QE, AfriCOMET, AfriCOMET-QE, and IndicCOMET because they had shown good performance in previously-published evaluations (Tomani et al., 2023; Wang et al., 2024; Sai B et al., 2023; Freitag et al., 2023b), had good performance in automated evaluations on the dev set (Appendix E), and lacked restrictions on commercial use. In our en-de and zh-en evaluations we also included metrics and ensembles with restrictions on commercial use (XCOMET, CometKiwi, rankAvg:all) for comparison. The 6 language pairs and datasets we evaluate are en-ha en-sw en-ml en-hi (from FLORES200 test) and en-de zh-en (from WMT2023). We chose these languages to have a wide distribution in resource level. For each language pair, we sampled 400 source segments to evaluate. WMT2023 was evaluated with document context, whereas FLORES200 segments were evaluated in isolation. We asked each rater to provide MQM annotations for all translation candidates for each source segment (we evaluted 15 systems on en-de and zh-en and 11 systems on others), and compute scores as described in Freitag et al. (2021). Scores range from 0 to 25, lower is better. To control for variance between raters, the same rater was used to score all candidate translations resulting from each source segment.

### 5.2 Results

Results are shown in Table 3. We observe that overall the best-performing system is rankAvg:noNC, which significantly outperforms greedy (p<0.001 on pairwise t-test). rankAvg:noNC also performs

| | Greedy | Reference | MetricX | MetricX-QE | XCOMET-XXL | CometKiwi23-XXL | COMET22 | AfriCOMET | AfriCOMET-QE | IndicCOMET | rankAvg:all | rankAvg:noNC | rankAvg:noNCnoLex | mxQE32mxMBR | noncQE32noncnolexMBR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| all:total | 1.52 | **1.80†** | 1.59 | **1.77†** | | | | | | | | **1.27‡** | **1.28‡** | 1.53 | **1.27‡** |
| en-de:total | 2.22 | 2.52 | 2.38 | 2.32 | 2.74 | **2.96\*** | 2.07 | | | | 2.07 | 1.89 | 1.83 | 2.13 | **1.69\*** |
| zh-en:total | 2.56 | 2.42 | **3.15†** | 3.05 | **3.04\*** | 2.98 | 2.65 | | | | 2.43 | 2.49 | 2.53 | 2.81 | 2.55 |
| en-sw:total | 1.03 | 1.41 | 1.08 | 0.95 | | | | 0.97 | **1.44\*** | | | **0.75\*** | 0.82 | 0.99 | 0.86 |
| en-ha:total | 1.02 | 1.25 | 1.07 | 1.04 | | | | 1.17 | 1.29 | | | 0.85 | 0.95 | 0.98 | 0.87 |
| en-hi:total | 0.95 | **1.50‡** | 0.70 | 1.09 | | | | | | 0.93 | | 0.78 | 0.71 | 0.86 | **0.70\*** |
| en-ml:total | 1.74 | 1.94 | 1.70 | **2.60‡** | | | | | | **2.29†** | | **1.31\*** | **1.28\*** | 1.84 | **1.39\*** |
| | | | | | | | | | | | | | | | |
| all:fluency | 0.29 | **0.38†** | 0.30 | **0.33†** | | | | | | | | **0.30‡** | **0.32‡** | 0.26 | **0.26‡** |
| en-de:fluency | 0.46 | 0.45 | 0.50 | 0.39 | 0.46 | **0.75\*** | 0.38 | | | | 0.45 | 0.47 | 0.45 | 0.29 | **0.37\*** |
| zh-en:fluency | 0.42 | 0.43 | **0.28†** | 0.24 | **0.27\*** | 0.39 | 0.32 | | | | 0.35 | 0.37 | 0.39 | 0.19 | 0.30 |
| en-sw:fluency | 0.14 | 0.18 | 0.17 | 0.27 | | | | 0.21 | **0.26\*** | | | **0.13\*** | 0.12 | 0.19 | 0.13 |
| en-ha:fluency | 0.37 | 0.49 | 0.38 | 0.36 | | | | 0.48 | 0.47 | | | 0.33 | 0.37 | 0.32 | 0.33 |
| en-hi:fluency | 0.17 | **0.32‡** | 0.24 | 0.30 | | | | | | 0.20 | | 0.26 | 0.26 | 0.24 | **0.16\*** |
| en-ml:fluency | 0.24 | 0.42 | 0.26 | **0.37‡** | | | | | | **0.30†** | | **0.26\*** | **0.33\*** | 0.31 | **0.29\*** |
| | | | | | | | | | | | | | | | |
| all:accuracy | 0.80 | **0.94†** | 0.98 | **1.06†** | | | | | | | | **0.70‡** | **0.70‡** | 0.95 | **0.74‡** |
| en-de:accuracy | 1.06 | 1.45 | 1.24 | 1.42 | 1.62 | **1.53\*** | 1.12 | | | | 1.11 | 0.86 | 0.90 | 1.14 | **0.85\*** |
| zh-en:accuracy | 1.72 | 1.67 | **2.57†** | 2.54 | **2.44\*** | 2.25 | 2.00 | | | | 1.74 | 1.80 | 1.79 | 2.34 | 1.96 |
| en-sw:accuracy | 0.58 | 0.48 | 0.59 | 0.44 | | | | 0.52 | **0.76\*** | | | **0.40\*** | 0.44 | 0.51 | 0.47 |
| en-ha:accuracy | 0.50 | 0.62 | 0.59 | 0.44 | | | | 0.54 | 0.70 | | | 0.45 | 0.45 | 0.58 | 0.46 |
| en-hi:accuracy | 0.32 | **0.65‡** | 0.32 | 0.46 | | | | | | 0.44 | | 0.25 | 0.22 | 0.41 | **0.32\*** |
| en-ml:accuracy | 0.94 | 1.07 | 1.11 | **1.56‡** | | | | | | **1.65†** | | **0.80\*** | **0.75\*** | 1.19 | **0.77\*** |
| | | | | | | | | | | | | | | | |
| all:other | 0.43 | **0.48†** | 0.30 | **0.38†** | | | | | | | | **0.28‡** | **0.26‡** | 0.32 | **0.27‡** |
| en-de:other | 0.69 | 0.62 | 0.64 | 0.51 | 0.66 | **0.68\*** | 0.58 | | | | 0.51 | 0.56 | 0.47 | 0.71 | **0.46\*** |
| zh-en:other | 0.42 | 0.32 | **0.30†** | 0.27 | **0.33\*** | 0.35 | 0.32 | | | | 0.35 | 0.32 | 0.35 | 0.27 | 0.30 |
| en-sw:other | 0.31 | 0.74 | 0.32 | 0.24 | | | | 0.24 | **0.42\*** | | | **0.22\*** | 0.25 | 0.29 | 0.25 |
| en-ha:other | 0.15 | 0.13 | 0.10 | 0.23 | | | | 0.15 | 0.12 | | | 0.06 | 0.13 | 0.08 | 0.08 |
| en-hi:other | 0.46 | **0.52‡** | 0.14 | 0.33 | | | | | | 0.29 | | 0.26 | 0.22 | 0.21 | **0.22\*** |
| en-ml:other | 0.56 | 0.46 | 0.33 | **0.67‡** | | | | | | **0.33†** | | **0.25\*** | **0.20\*** | 0.34 | **0.32\*** |

Table 3: Human evaluation results broken down by language and MQM error type. Columns indicate the system used for MBR/QE decoding; ensembles are defined in Table 2. Rows starting with "all" shows results across all languages. 1st block is total error scores, 2nd is fluency error scores, 3rd is accuracy error scores, 4th is other error scores. For each system, average human evaluation scores across the evaluated segments are shown. Lower scores are better. Colors are relative to greedy, green is better than greedy, red is worse. Black cells were not evaluated. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001.

the best on each language pair except en-hi. Interestingly, rankAvg:noNC and greedy decoding beat the reference translation in all language pairs, suggesting either that the reference translations in WMT2023 and FLORES200 are of poor quality, or that Gemini's translation quality has achieved human parity for these language pairs.

A surprising result from our human evaluation was that although MBR decoding with an ensembles of metrics was judged as having superior quality to greedy decoding, MBR/QE decoding with a single metric (MetricX, MetricX-QE, XCOMET-XXL, CometKiwi23-XXL, COMET22, AfriCOMET, AfriCOMET-QE, IndicCOMET) did not generally improve over greedy decoding (Table 3). In fact, translations from MetricX MBR decoding for zh-en, MetricX-QE decoding for en-ml, AfriCOMET-QE decoding for en-sw, and IndicCOMET MBR decoding for en-ml were rated by humans as significantly worse than greedy decoding (Table 3), even though automatic evaluation with other neural metrics such as MetricX and XCOMET-XXL estimated those translations as being significantly better than greedy (Appendix G). This suggests that evaluation with neutral metrics overestimates the quality of MBR/QE decoding, even if different metrics are used for decoding and evaluation. Our findings contrast with previous

| System | Translation | Fluency MQM | Accuracy MQM | Other MQM | MetricX | MetricX-QE | XCOMET-XXL | CometKiwi23-XXL | COMET22 |
|---|---|---|---|---|---|---|---|---|---|
| Greedy | The seller said not yet, and it **will** be shipped in the afternoon. | 1.0 | 0.0 | 0.0 | 0.659 | 0.88 | 0.999 | 0.83 | 0.74 |
| MetricX /XCOMET-XXL | The seller said that they **don't have it in stock yet**, and **will** be able to ship it out **this afternoon**. | 1.0 | 10.0 | 0.0 | 0.259 | 0.94 | 1.000 | 0.70 | 0.68 |
| MetricX-QE | The seller said he hadn't shipped it, but could ship it that afternoon. | 0.0 | 0.0 | 0.0 | 0.438 | 0.49 | 0.997 | 0.78 | 0.68 |
| CometKiwi23-XXL | The seller said that it was not ready yet and that it would be shipped that afternoon. | 0.0 | 0.0 | 0.0 | 0.264 | 0.67 | 0.998 | 0.87 | 0.73 |
| COMET22 | The seller said not **yet, it** will be sent in the afternoon. | 1.0 | 0.0 | 0.0 | 0.981 | 1.06 | 0.998 | 0.86 | 0.76 |
| noncQE32noncnolexMBR /rankAvg:noNCnoLex | The seller **said no, it** won't be shipped until this afternoon. | 1.0 | 0.0 | 1.0 | 0.552 | 0.60 | 0.998 | 0.76 | 0.77 |
| rankAvg:noNC /rankAvg:all | The seller said not **yet, it will** be shipped in the afternoon. | 2.0 | 0.0 | 0.0 | 0.608 | 0.90 | 0.998 | 0.84 | 0.71 |
| mxQE32mxMBR | The seller said that it is not yet ready, and it will be shipped in the afternoon. | 0.0 | 5.0 | 0.0 | 0.432 | 0.75 | 0.998 | 0.84 | 0.73 |

Table 4: An example where MetricX and XCOMET-XXL MBR decoding result in an inaccurate translation. The source text is 卖家说还没，下午才能发。("Seller says not yet, can ship in the afternoon.") The preceding sentence is 结果，第二天打电话问，发货了吗？("So the next day I called to ask, has it shipped?"). MetricX and XCOMET-XXL MBR decoding, as well as the reference-based MetricX and XCOMET-XXL evaluations, all prefer a translation which inaccurately states the item is out of stock. The other metrics assign a lower score to the inaccurate translation. Lower scores are better for MQM, MetricX, and MetricX-QE, for other metrics higher is better. Green is better than greedy, red is worse. Spans marked as errors by the rater are bolded.

studies which find that MBR decoding with a single metric outperforms greedy decoding in human evaluations (Freitag et al., 2022, 2023a; Tomani et al., 2023).

We hypothesize a few potential causes of the failure of single-metric MBR/QE decoding to outperform greedy decoding: firstly, machine translation quality has improved considerably in recent years. This is reflected by how in our study the greedy decoding outputs achieved better human evaluation results compared to the references generated by professional human translators, especially when looking at fluency scores (Table 3), in contrast with previous work where reference translations were rated as better (Freitag et al., 2022, 2023a). Therefore, it is possible that improvements in greedy translation quality have reduced the quality gains from MBR/QE decoding, and have resulted in the adverse effects of metric bias from MBR/QE decoding with a single utility metric outweighing the benefits to translation quality. For example, in Table 3 we can see that single-metric MBR/QE decoding generally improves fluency on high-resource languages, and reduces errors in style, terminology, and locale convention (labeled "other"). How-

ever, accuracy suffers with single-metric MBR/QE decoding for most language pairs (Table 3). We show an example in Table 4, where MetricX and XCOMET-XXL MBR decoding favor a fluent yet inaccurate translation. Perhaps part of the reason for this decrease in accuracy is that MBR decoding with metrics such as MetricX considers only similarity to the pseudoreferences and does not consider the source sentence, so fluent hallucinations that occur in a large number of pseudoreferences will be favored by MBR decoding. Therefore, we hypothesize that past gains from single-metric MBR/QE decoding might have been driven by improvements in fluency and style, but modern LLMs have become good at producing fluent outputs (as indicated by the low fluency error scores for the greedy condition in Table 3), so we are no longer seeing overall quality improvements from single-metric MQM/QE decoding.

We also considered the effects of domain on the quality of single-metric MBR/QE decoding. Since the WMT2023 datasets which were used include novel domains such as speech transcripts and mastodon posts which are not well-represented in the data that metrics such as MetricX and

| | Greedy | Reference | MetricX | MetricX-QE | XCOMET-XXL | CometKiwi23-XXL | COMET22 | rankAvg:all | rankAvg:noNC | rankAvg:noNCnoLex | mxQE32mxMBR | noncQE32noncnolexMBR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| en-de@news:total | 1.95 | 2.97 | 3.47 | 3.28 | 2.52 | 3.74 | 1.99 | 1.99 | 2.16 | 1.91 | 2.05 | 1.78 |
| en-de@user-review:total | 3.66 | 2.79 | 3.30 | 2.80 | 3.11 | 4.07 | 3.90 | 3.71 | 2.81 | 3.12 | 3.09 | 2.68 |
| en-de@mastodon:total | 1.29 | 1.70 | 1.17 | 1.60 | 1.60 | **1.87***  | 1.13 | 1.19 | 1.04 | 1.03 | 1.65 | 0.98 |
| en-de@speech:total | 3.59 | 3.78 | 3.37 | 2.60 | **5.43*** | 3.83 | 2.97 | 2.97 | 2.88 | 2.65 | 2.61 | 2.48 |
| zh-en@news:total | 3.56 | 4.51 | 3.90 | 3.90 | 3.83 | 4.16 | 3.36 | 2.98 | 3.90 | 3.15 | 3.83 | 4.11 |
| zh-en@user-review:total | 2.28 | 1.73 | **2.93*** | 2.71 | **2.83*** | 2.62 | 2.45 | 2.22 | 2.06 | 2.42 | 2.39 | 2.01 |
| zh-en@manuals:total | 1.70 | 1.32 | **2.60*** | 2.98 | 2.28 | 2.21 | 2.01 | 2.35 | 1.58 | 1.55 | 2.76 | 1.89 |
| | | | | | | | | | | | | |
| en-de@news:fluency | 0.38 | 0.69 | 1.33 | 0.77 | 0.31 | 1.49 | 0.44 | 0.46 | 0.84 | 0.66 | 0.36 | 0.42 |
| en-de@user-review:fluency | 0.57 | 0.65 | 0.37 | 0.70 | 0.52 | 0.89 | 0.88 | 0.53 | 0.49 | 0.79 | 0.18 | 0.60 |
| en-de@mastodon:fluency | 0.15 | 0.21 | 0.17 | 0.13 | 0.25 | **0.30*** | 0.15 | 0.18 | 0.12 | 0.15 | 0.22 | 0.21 |
| en-de@speech:fluency | 1.21 | 0.63 | 0.48 | 0.34 | **1.05*** | 0.90 | 0.47 | 1.03 | 0.87 | 0.70 | 0.48 | 0.55 |
| zh-en@news:fluency | 0.29 | 1.02 | 0.42 | 0.36 | 0.38 | 0.46 | 0.29 | 0.32 | 0.37 | 0.42 | 0.25 | 0.33 |
| zh-en@user-review:fluency | 0.51 | 0.18 | **0.22*** | 0.18 | **0.23*** | 0.31 | 0.32 | 0.33 | 0.34 | 0.37 | 0.10 | 0.22 |
| zh-en@manuals:fluency | 0.20 | 0.47 | **0.29*** | 0.37 | 0.28 | 0.71 | 0.43 | 0.58 | 0.51 | 0.43 | 0.57 | 0.70 |
| | | | | | | | | | | | | |
| en-de@news:accuracy | 0.65 | 1.55 | 1.63 | 2.14 | 1.59 | 1.63 | 0.96 | 0.97 | 0.59 | 0.78 | 1.12 | 0.95 |
| en-de@user-review:accuracy | 2.32 | 1.25 | 0.89 | 1.47 | 1.16 | 1.96 | 1.79 | 2.37 | 1.25 | 1.32 | 1.00 | 0.82 |
| en-de@mastodon:accuracy | 0.54 | 1.06 | 0.68 | 0.97 | 0.95 | **1.10*** | 0.63 | 0.57 | 0.57 | 0.53 | 1.01 | 0.47 |
| en-de@speech:accuracy | 1.77 | 2.41 | 2.44 | 1.66 | **3.65*** | 2.13 | 1.94 | 1.56 | 1.56 | 1.61 | 1.56 | 1.66 |
| zh-en@news:accuracy | 3.03 | 3.10 | 3.21 | 3.39 | 3.14 | 3.36 | 2.79 | 2.40 | 3.21 | 2.39 | 3.39 | 3.56 |
| zh-en@user-review:accuracy | 1.23 | 1.21 | **2.36*** | 2.20 | **2.24*** | 1.92 | 1.75 | 1.48 | 1.35 | 1.66 | 1.99 | 1.46 |
| zh-en@manuals:accuracy | 1.38 | 0.81 | **2.19*** | 2.46 | 1.88 | 1.38 | 1.50 | 1.62 | 0.96 | 1.04 | 1.85 | 0.88 |
| | | | | | | | | | | | | |
| en-de@news:other | 0.92 | 0.73 | 0.51 | 0.37 | 0.62 | 0.62 | 0.59 | 0.55 | 0.73 | 0.46 | 0.58 | 0.41 |
| en-de@user-review:other | 0.77 | 0.89 | 2.04 | 0.63 | 1.44 | 1.21 | 1.23 | 0.81 | 1.07 | 1.02 | 1.91 | 1.26 |
| en-de@mastodon:other | 0.60 | 0.43 | 0.31 | 0.49 | 0.39 | **0.47*** | 0.36 | 0.44 | 0.36 | 0.35 | 0.42 | 0.30 |
| en-de@speech:other | 0.61 | 0.75 | 0.45 | 0.61 | **0.73*** | 0.80 | 0.55 | 0.38 | 0.45 | 0.34 | 0.56 | 0.27 |
| zh-en@news:other | 0.24 | 0.39 | 0.27 | 0.16 | 0.31 | 0.34 | 0.29 | 0.26 | 0.31 | 0.34 | 0.19 | 0.23 |
| zh-en@user-review:other | 0.54 | 0.34 | **0.35*** | 0.33 | **0.36*** | 0.38 | 0.38 | 0.41 | 0.36 | 0.39 | 0.30 | 0.32 |
| zh-en@manuals:other | 0.12 | 0.04 | **0.12*** | 0.15 | 0.12 | 0.12 | 0.08 | 0.15 | 0.12 | 0.08 | 0.35 | 0.31 |

Table 5: Human evaluation results broken down by domain and MQM error type for en-de and zh-en. Columns indicate the system used for MBR/QE decoding; ensembles are defined in Table 2. 1[st] block is total error scores, 2[nd] is fluency error scores, 3[rd] is accuracy error scores, 4[th] is other error scores. For each system, average human evaluation scores across the evaluated segments are shown. Lower scores are better. Colors are relative to greedy, green is better than greedy, red is worse. Significant differences from greedy (pairwise t-test) indicated by * for $p<0.05$, † for $p<0.01$, ‡ for $p<0.001$.

XCOMET-XXL were trained on, we hypothesized that this may adversely impacting MBR quality. However, contrary to our expectations, as we can observe in Table 5 there is no clear effect of the domain on the quality of MBR decoding results. Thus, we do not believe effects of domain to be the primary factor behind our findings.

We also considered whether MBR decoding with other metrics we did not evaluate with human raters, such as BLEURT, would have performed better than the metrics we evaluated. To do so, we looked at the correlation between the MQM scores from our human evaluation, compared to the scores assigned by metrics. We include scores from QE

metrics (to simulate QE decoding), scores from reference-based metrics based on the 128 pseudoreferences (to simulate MBR decoding), as well as scores form reference-based metrics using the actual references (to simulate a reference-based metric oracle). Table 6 shows Kendall-Tau correlation and Table 7 shows Pearson correlation. Note that this an imperfect simulation of what would happen if we actually performed human evaluation with the MBR/QE decoding outputs for these metrics, as we are considering correlations with human judgements only the subset of candidates which were evaluated (which is a biased sample, as they are the results of MBR/QE decoding), not all 128 samples.

We observe that among the individual metrics that we did not evaluate, simulated XCOMET-XL MBR decoding seems to correlate the best with human judgements, and the other metrics are generally worse than MetricX/XCOMET-XXL MBR decoding. We also include some ensembles, finding that they are generally better correlated with human judgements than individual metrics in our simulation. Therefore, we do not expect that changing to another metric for MBR/QE decoding would have resulted in significantly better translation quality.

# 6 Discussion

While previous work has sometimes assumed that MBR decoding outputs can be evaluated by automated metrics so long as a non-utility metric is used (Tomani et al., 2023), we find MBR/QE decoding outputs are often preferred by automated metrics despite the fact that human raters believe they are worse quality. For example, while MetricX-QE decoding outputs are considered by human raters to be of worse quality than greedy decoding (Table 3), they still achieve higher scores when evaluated by XCOMET-XXL, XCOMET-XL, MetricX, CometKiwi22, CometKiwi23-XL, and CometKiwi23-XXL (Table 1 and Appendix G). Thus, the metric bias issue that results from MBR/QE decoding complicates evaluation with automated metrics.

That said, while we have shown that MBR/QE decoding generated translations with higher automated evaluation scores are not always judged as having better quality by humans, this does not mean that automated metrics are no longer useful. In our study, automatic reference-based metrics, QE metrics, and ensembles of metrics are still somewhat correlated with MQM scores, as shown in Table 6. Therefore, while it is advisable to perform a human evaluation when feasible if evaluating systems that make use of MBR/QE decoding, existing metrics still correlate with human preferences. Additionally, using an ensemble of metrics for MBR decoding results in improved translation quality compared to greedy decoding and MBR/QE decoding with a single metric (Table 3).

Why is it that using an ensemble of metrics for MBR decoding improves translation quality compared to just using a single metric (Table 3)? We hypothesize that each metric has its own biases that lead it to prefer bad translations, but different metrics have different biases, so using an ensemble

| | zh-en | en-de | en-ha | en-sw | en-hi | en-ml |
|---|---|---|---|---|---|---|
| XCOMET-XXL | 0.278 | 0.110 | 0.114 | 0.201 | 0.073 | 0.152 |
| XCOMET-XXL:mbr | 0.275 | 0.111 | 0.125 | 0.212 | 0.094 | 0.152 |
| XCOMET-XL | 0.335 | 0.126 | 0.123 | 0.187 | 0.087 | 0.179 |
| XCOMET-XL:mbr | 0.336 | 0.134 | 0.137 | 0.201 | 0.093 | 0.168 |
| MetricX | 0.252 | 0.065 | 0.077 | 0.192 | 0.087 | 0.154 |
| MetricX:mbr | 0.289 | 0.089 | 0.111 | 0.211 | 0.097 | 0.149 |
| MetricX-QE | 0.291 | 0.046 | 0.093 | 0.166 | 0.065 | 0.130 |
| CometKiwi23-XXL | 0.264 | 0.080 | 0.115 | 0.160 | 0.085 | 0.140 |
| CometKiwi23-XL | 0.281 | 0.094 | 0.113 | 0.138 | 0.101 | 0.165 |
| CometKiwi22 | 0.274 | 0.107 | 0.032 | 0.173 | 0.087 | 0.179 |
| COMET22 | 0.271 | 0.100 | 0.062 | 0.179 | 0.076 | 0.166 |
| COMET22:mbr | 0.290 | 0.125 | 0.067 | 0.183 | 0.088 | 0.159 |
| BLEURT | 0.279 | 0.128 | 0.098 | 0.173 | 0.083 | 0.146 |
| BLEURT:mbr | 0.271 | 0.134 | 0.119 | 0.187 | 0.108 | 0.132 |
| YiSi | 0.178 | 0.049 | 0.072 | 0.105 | 0.061 | 0.138 |
| YiSi:mbr | 0.183 | 0.068 | 0.096 | 0.119 | 0.065 | 0.154 |
| chrF | 0.044 | 0.040 | 0.083 | 0.115 | 0.067 | 0.129 |
| chrF:mbr | 0.091 | 0.049 | 0.098 | 0.135 | 0.056 | 0.146 |
| chrF++ | 0.057 | 0.045 | 0.084 | 0.118 | 0.064 | 0.123 |
| chrF++:mbr | 0.103 | 0.052 | 0.098 | 0.135 | 0.057 | 0.141 |
| sentBLEU | 0.102 | 0.059 | 0.072 | 0.106 | 0.052 | 0.083 |
| sentBLEU:mbr | 0.155 | 0.058 | 0.082 | 0.121 | 0.058 | 0.103 |
| TER | 0.129 | 0.061 | 0.084 | 0.086 | 0.077 | 0.087 |
| TER:mbr | 0.114 | 0.060 | 0.088 | 0.097 | 0.067 | 0.116 |
| MetricX +MetricX-QE | 0.287 | 0.055 | 0.084 | 0.196 | 0.088 | 0.155 |
| MetricX +MetricX-QE | 0.304 | 0.070 | 0.107 | 0.203 | 0.097 | 0.151 |
| XCOMET-XXL +XCOMET-XL | 0.326 | 0.124 | 0.121 | 0.210 | 0.088 | 0.186 |
| XCOMET-XXL:mbr +XCOMET-XL:mbr | 0.324 | 0.131 | 0.136 | 0.216 | 0.098 | 0.177 |
| XCOMET-XXL +XCOMET-XL +COMET22 | 0.346 | 0.127 | 0.116 | 0.213 | 0.090 | 0.193 |
| XCOMET-XXL:mbr +XCOMET-XL:mbr +COMET22:mbr | 0.348 | 0.140 | 0.129 | 0.220 | 0.100 | 0.184 |

Table 6: Kendall-Tau correlation between MQM evaluation scores and automated evaluation scores. For reference-based metrics, rows with ":mbr" indicate pseudoreference-based evaluation. Bottom rows are ensembles that take the average between the listed metrics. Higher scores indicate better agreement with human raters. See Table 7 for Pearson correlation.

reduces metric bias. We see an example of this in Table 4 where MetricX and XCOMET-XXL assign high scores to an inaccurate translation, but this translation is rated poorly by CometKiwi23-XXL and COMET22, so the ensemble ends up picking a good translation that is preferred by all metrics.

Techniques other than MBR/QE decoding for making use of human preferences to improve translation quality, such as DPO (direct preference optimization) (Rafailov et al., 2024; Yang et al., 2024)) and RLHF (reinforcement learning from human feedback) (Christiano et al., 2017), might be more resilient to this metric bias issue, as they do not

directly make use of the evaluation metric. However, given that the data used for DPO/RLHF is similar to the data used to train evaluation metrics, and given that the reward hacking issue is prevalent throughout reinforcement learning (Skalse et al., 2022), issues similar to metric bias may still occur with these techniques.

An open question that remains is how to develop new evaluation techniques that are resilient to metric bias in MBR/QE decoding. One potential way is to develop metrics specialized for evaluating MBR/QE decoding outputs from a particular system, by generating MBR/QE decoding outputs from a translation model, obtaining human annotations for those, and training a metric with them. This process is unfortunately costly and time-intensive, and the learned metric might not be able to generalize beyond translations generated by the particular utility metric and translation model it was trained on. Perhaps a better approach would be to view the metric bias problem as an adversarial learning problem, and apply techniques such as generative adversarial training (Yang et al., 2018) to help train metrics resilient to MBR bias.

## 7   Conclusion

In this paper we have explored the problem of metric bias, where MBR or QE decoding with a single utility metric shows improvements on automated evaluation with the utility metric and related metrics, but does not actually improve quality when judged by a human rater. We find that the metric bias issue is most severe when using a single utility metric, and using an ensemble of metrics to perform MBR decoding can help improve quality as judged by a human rater. While we have shown that metric bias can result in overly-optimistic automatic evaluations of systems that make use of MBR/QE decoding, the question of how to resolve this issue and automatically evaluate systems that make use of MBR/QE decoding is still an open problem which we leave to future work.

## Dataset

Dataset is at `https://mbrbias.github.io/`

## Limitations

In this work we compare to only full MBR decoding and QE filtering as baselines, but there are many alternative approaches, such as MBR approximation heuristics (Trabelsi et al., 2024; Jinnai and

Ariu, 2024; Deguchi et al., 2024, 2023; Vamvas and Sennrich, 2024; Eikema and Aziz, 2022), direct preference optimization training (Yang et al., 2024), quality-aware training (Tomani et al., 2023), or training on MBR decoding outputs (Finkelstein and Freitag, 2023), that are more practical to use if translation latency is important. In this work we only look at translations coming from Gemini 1.0 Pro with 5-shot sample prompts and epsilon sampling, and it is possible that results may differ if using a different translation system, different prompts, or a different sampling technique. In this work we only look at using 128 samples due to the computationally expensive $O(n^2)$ cost of running full MBR decoding, but it is possible that using additional samples can achieve further quality improvements. In this work we only looked at segment-level translation, and it is possible that results may differ if performing document-level translation. However, MetricX and the COMET families of models have input token limits – 1024 tokens for MetricX, 512 tokens for COMET – which make it difficult to use them for document-level MBR decoding. Our human evaluation used only a single rater for each translation, which introduces the question of how reliable and consistent the ratings are – using multiple raters and looking at inter-rater agreement is preferable, but was beyond our budget constraints.

## Ethics Statement

MBR decoding is resource-intensive, and using ensembles of multiple metrics increases computational complexity compared to a single utility metric. To mitigate this issue, we presented two-step ensembles that use QE filtering followed by MBR decoding, which reduce the computational cost below the cost of standard MBR decoding with a single metric.

## References

Chantal Amrhein and Rico Sennrich. 2022. Identifying weaknesses in machine translation metrics through minimum Bayes risk decoding: A case study for COMET. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1125–1141, Online only. Association for Computational Linguistics.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini,

Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Julius Cheng and Andreas Vlachos. 2023. Faster minimum Bayes risk decoding with confidence-based pruning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12473–12480, Singapore. Association for Computational Linguistics.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Hiroyuki Deguchi, Kenji Imamura, Yuto Nishida, Yusuke Sakai, Justin Vasselli, and Taro Watanabe. 2023. NAIST-NICT WMT'23 general MT task submission. In *Proceedings of the Eighth Conference on Machine Translation*, pages 110–118, Singapore. Association for Computational Linguistics.

Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, Taro Watanabe, Hideki Tanaka, and Masao Utiyama. 2024. Centroid-based efficient minimum bayes risk decoding. *arXiv preprint arXiv:2402.11197*.

Bryan Eikema and Wilker Aziz. 2022. Sampling-based approximations to minimum Bayes risk decoding for neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10978–10993, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Mara Finkelstein and Markus Freitag. 2023. Mbr and qe finetuning: Training-time distillation of the best and most expensive decoding methods. *arXiv preprint arXiv:2309.10966*.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023a. Epsilon sampling rocks: Investigating sampling strategies for minimum Bayes risk decoding for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9198–9209, Singapore. Association for Computational Linguistics.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023b. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Gemini Team Google. 2023. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Taisiya Glushkova, Chrysoula Zerva, and André F. T. Martins. 2023. BLEU meets COMET: Combining lexical and neural metrics towards robust machine translation evaluation. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 47–58, Tampere, Finland. European Association for Machine Translation.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.

Kamil Guttmann, Mikołaj Pokrywka, Adrian Charkiewicz, and Artur Nowakowski. 2024. Chasing comet: Leveraging minimum bayes risk decoding for self-improving machine translation. *arXiv preprint arXiv:2405.11937*.

Yuu Jinnai and Kaito Ariu. 2024. Hyperparameter-free approach for faster minimum bayes risk decoding. *arXiv preprint arXiv:2401.02749*.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme

Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Zhongtao Liu, Parker Riley, Alison Lui, Daniel Deutsch, Mengmeng Niu, Apurva Shah, and Markus Freitag. 2024. Beyond human-only: Evaluating human-machine collaboration for collecting high-quality translation data. *arXiv preprint*.

Chi-kiu Lo. 2019. YiSi - a unified semantic MT quality evaluation and estimation metric for languages with different levels of available resources. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 507–513, Florence, Italy. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn.

2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ananya Sai B, Tanay Dixit, Vignesh Nagarajan, Anoop Kunchukuttan, Pratyush Kumar, Mitesh M. Khapra, and Raj Dabre. 2023. IndicMT eval: A dataset to meta-evaluate machine translation metrics for Indian languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14210–14228, Toronto, Canada. Association for Computational Linguistics.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020a. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Thibault Sellam, Amy Pu, Hyung Won Chung, Sebastian Gehrmann, Qijun Tan, Markus Freitag, Dipanjan Das, and Ankur Parikh. 2020b. Learning to evaluate translation beyond English: BLEURT submissions to the WMT metrics 2020 shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 921–927, Online. Association for Computational Linguistics.

Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. 2022. Defining and characterizing reward hacking. *Advances in Neural Information Processing Systems*, 35:9460–9471.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Christian Tomani, David Vilar, Markus Freitag, Colin Cherry, Subhajit Naskar, Mara Finkelstein, and Daniel Cremers. 2023. Quality control at your fingertips: Quality-aware translation models. *arXiv preprint arXiv:2310.06707*.

Firas Trabelsi, David Vilar, Mara Finkelstein, and Markus Freitag. 2024. Efficient minimum bayes risk decoding using low-rank matrix completion algorithms. *arXiv preprint arXiv:2406.02832*.

Jannis Vamvas and Rico Sennrich. 2024. Linear-time minimum bayes risk decoding with reference aggregation. *arXiv preprint arXiv:2402.04251*.

Jiayi Wang, David Adelani, Sweta Agrawal, Marek Masiak, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Xuanli He, Sofia Bourhim, Andiswa Bukula, Muhidin Mohamed, Temitayo Olatoye, Tosin Adewumi, Hamam Mokayed, Christine Mwase, Wangui Kimotho, Foutse Yuehgoh, Anuoluwapo Aremu, Jessica Ojo, Shamsuddeen Muhammad, Salomey Osei, Abdul-Hakeem Omotayo, Chiamaka Chukwuneke, Perez Ogayo, Oumaima Hourrane, Salma El Anigri, Lolwethu Ndolela, Thabiso Mangwana, Shafie Mohamed, Hassan Ayinde, Oluwabusayo Awoyomi, Lama Alkhaled, Sana Al-azzawi, Naome Etori, Millicent Ochieng, Clemencia Siro, Njoroge Kiragu, Eric Muchiri, Wangari Kimotho, Toadoum Sari Sakayo, Lyse Naomi Wamba, Daud Abolade, Simbiat Ajao, Iyanuoluwa Shode, Ricky Macharm, Ruqayya Iro, Saheed Abdullahi, Stephen Moore, Bernard Opoku, Zainab Akinjobi, Abeeb Afolabi, Nnaemeka Obiefuna, Onyekachi Ogbu, Sam Ochieng', Verrah Otiende, Chinedu Mbonu, Yao Lu, and Pontus Stenetorp. 2024. AfriMTE and AfriCOMET: Enhancing COMET to embrace under-resourced African languages. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 5997–6023, Mexico City, Mexico. Association for Computational Linguistics.

Guangyu Yang, Jinghong Chen, Weizhe Lin, and Bill Byrne. 2024. Direct preference optimization for neural machine translation with minimum Bayes risk decoding. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 391–398, Mexico City, Mexico. Association for Computational Linguistics.

Zhen Yang, Wei Chen, Feng Wang, and Bo Xu. 2018. Improving neural machine translation with conditional sequence generative adversarial nets. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1346–1355, New Orleans, Louisiana. Association for Computational Linguistics.

## A  Methodology Details

### A.1  Prompts Used for Generating Samples

For each language pair, we obtained 5-shot examples for our prompts from the dev split of FLORES-200 by randomly sampling among those reference pairs that had perfect MetricX QE scores (scores of 0). We used MetricX QE filtering to ensure we used high-quality examples as our 5-shot examples. The sampled examples and prompt text for each language pair is included in our dataset release.

### A.2  Instructions for Computing Metrics

sentBLEU, chrF, chrF++, and TER scores were computed with sacreBLEU 2.4.2 (Post, 2018) on python 3.11.8 with the following parameters:

chrF: -m chrf

chrF++: -m chrf –chrf-word-order 2

sentBLEU: -m bleu –sentence-level

TER: -m ter

For other metrics, we used the publicly released models on HuggingFace, running with the unbabel-comet package version 2.2.1 available on pip, on Python 3.10.14. We ran on an NVIDIA A100 GPU for all metrics except XCOMET-XXL and CometKiwi23-XXL, which required an NVIDIA A100 80GB GPU.

## B  Metrics Included in Each Ensemble

This section presents the same information that is present in Table 2, but in textual format. The following are the groups of metrics included in the single-step ensembles that we include in our study. For each of these metric groups the rankAvg, rankMed, rankMax, and rank75q ensembling techniques are used to generate an ensemble.

1. all: All metrics, both reference-based and QE (MetricX, MetricX-QE, XCOMET-XXL, XCOMET-XL, CometKiwi23-XXL, CometKiwi23-XL, CometKiwi22, COMET22, BLEURT, YiSi, chrF, chrF++, sentBLEU, TER, AfriCOMET and

AfriCOMET-QE for African languages, IndicCOMET for Indic languages)

2. qe: All QE metrics (MetricX-QE, CometKiwi23-XXL, CometKiwi23-XL, CometKiwi22, and AfriCOMET-QE for African languages)

3. top: MetricX, MetricX-QE, XCOMET-XXL, XCOMET-XL, CometKiwi23-XXL, CometKiwi23-XL

4. topQe: MetricX-QE, CometKiwi23-XXL, CometKiwi23-XL

5. mxmxqe: MetricX, MetricX-QE

6. noLex: All non-lexical metrics (MetricX, MetricX-QE, XCOMET-XXL, XCOMET-XL, CometKiwi23-XXL, CometKiwi23-XL, CometKiwi22, COMET22, BLEURT, YiSi, AfriCOMET and AfriCOMET-QE for African languages, IndicCOMET for Indic languages)

7. noNC: All metrics that permit commercial use (MetricX, MetricX-QE, CometKiwi22, COMET22, BLEURT, YiSi, chrF, chrF++, sentBLEU, TER, AfriCOMET and AfriCOMET-QE for African languages, IndicCOMET for Indic languages)

8. noNCnoLex: All non-lexical metrics that permit commercial use (MetricX, MetricX-QE, COMET22, BLEURT, YiSi, AfriCOMET and AfriCOMET-QE for African languages, IndicCOMET for Indic languages)

9. noNCQe: All QE metrics that permit commercial use (MetricX-QE, and AfriCOMET-QE for African languages)

In addition, we also investigate QE filtering followed by MBR decoding (here we define QE filtering as selecting the top N candidates according to a QE metric, where N can be either 4, 8, 16, 32, 64). We include the following ensembles of this form:

1. allQE(N)allMBR: Use QE filtering with an ensemble of all QE metrics (MetricX-QE, CometKiwi23-XXL, CometKiwi23-XL, CometKiwi22, AfriCOMET-QE for African languages), then perform MBR decoding on the N resulting candidates with all reference-based metrics (MetricX, XCOMET-XXL, XCOMET-XL, COMET22, BLEURT, YiSi, chrF, chrF++, sentBLEU, TER, AfriCOMET for African languages, IndicCOMET for Indic languages).

2. allQE(N)nolexMBR: Use QE filtering with an ensemble of all QE metrics (MetricX-QE, CometKiwi23-XXL, CometKiwi23-XL,

CometKiwi22, AfriCOMET-QE for African languages), then perform MBR decoding on the N resulting candidates with all non-lexical reference-based metrics (MetricX, XCOMET-XXL, XCOMET-XL, COMET22, BLEURT, YiSi, AfriCOMET for African languages, IndicCOMET for Indic languages).

3. topQE(N)topMBR: Use QE filtering with an ensemble of top-performing QE metrics (MetricX QE, CometKiwi23-XXL, CometKiwi23-XL), then perform MBR decoding on the N resulting candidates with an ensemble of top-performing reference-based metrics (MetricX, XCOMET-XXL, XCOMET-XL).

4. noncQE(N)noncMBR: Use QE filtering with an ensemble of QE metrics that permit commercial use (MetricX-QE, AfriCOMET-QE for African languages), then perform MBR decoding with an ensemble of reference-based metrics that permit commercial use (MetricX, COMET22, BLEURT, YiSi, chrF, chrF++, sentBLEU, TER, AfriCOMET for African languages, IndicCOMET for Indic languages).

5. noncQE(N)noncnolexMBR: Use QE filtering with an ensemble of QE metrics that permit commercial use (MetricX-QE, AfriCOMET-QE for African languages), then perform MBR decoding with an ensemble of non-lexical reference-based metrics that permit commercial use (MetricX, COMET22, BLEURT, YiSi, AfriCOMET for African languages, IndicCOMET for Indic languages).

6. mxQE(N)xcMBR: Use QE filtering with MetricX-QE, then perform MBR decoding with XCOMET-XXL

7. ckQE(N)xcMBR: Use QE filtering with CometKiwi23-XXL, then perform MBR decoding with XCOMET-XXL

8. mxQE(N)mxMBR: Use QE filtering with MetricX-QE, then perform MBR decoding with MetricX

9. ckQE(N)mxMBR: Use QE filtering with CometKiwi23-XXL, then perform MBR decoding with MetricX

## C  Pseudocode for Ensembles

rankAvg ensembling strategy:

```
def rankAvg(
  sample_list: List[str], metric_list: List[str]
):
  sample_ranks =
      get_ranks_for_samples_by_ensemble(sample_list,
```

```
    metric_list)
  score_list = [np.mean(x) for x in
      sample_ranks]
  return select_samples_by_score(sample_list,
      score_list)
```

### rankMed ensembling strategy:

```
def rankMed(
  sample_list: List[str], metric_list: List[str]
):
  sample_ranks =
      get_ranks_for_samples_by_ensemble(sample_list,
      metric_list)
  score_list = [np.median(x) for x in
      sample_ranks]
  return select_samples_by_score(sample_list,
      score_list)
```

### rankMax ensembling strategy:

```
def rankMax(
  sample_list: List[str], metric_list: List[str]
):
  sample_ranks =
      get_ranks_for_samples_by_ensemble(sample_list,
      metric_list)
  score_list = [np.max(x) for x in sample_ranks]
  return select_samples_by_score(sample_list,
      score_list)
```

### rank75q ensembling strategy:

```
def rank75q(
  sample_list: List[str], metric_list: List[str]
):
  sample_ranks =
      get_ranks_for_samples_by_ensemble(sample_list,
      metric_list)
  score_list = [np.quantile(x, q=[0.75])[0] for
      x in sample_ranks]
  return select_samples_by_score(sample_list,
      score_list)
```

### Here are helper functions that were used:

```
def get_ranks_for_samples_by_ensemble(
  sample_list: List[str], metric_list: List[str]
):
  output = [[None for y in metric_list] for x
      in sample_list]
  for metric_idx, metric in
      enumerate(metric_list):
    sample_to_rank =
        rank_samples_by_metric(sample_list,
        metric)
    for sample_idx, sample in
        enumerate(sample_list):
      output[sample_idx][metric_idx] =
          sample_to_rank[sample]
  return output

def select_samples_by_score(
  sample_list: List[str],
  score_list: List[float]
):
  sample_with_score = zip(sample_list,
      score_list)
```

```
  top_candidate, top_score =
      min(sample_with_score, key=lambda x: x[1])
  return top_candidate
```

## D  Correlation Between Human Evaluation MQM Scores and Metrics

| | zh-en | en-de | en-ha | en-sw | en-hi | en-ml |
|---|---|---|---|---|---|---|
| XCOMET-XXL | 0.391 | 0.084 | 0.146 | 0.139 | 0.111 | 0.202 |
| XCOMET-XXL:mbr | 0.389 | 0.076 | 0.178 | 0.145 | 0.141 | 0.198 |
| XCOMET-XL | 0.543 | 0.126 | 0.154 | 0.160 | 0.141 | 0.208 |
| XCOMET-XL:mbr | 0.550 | 0.124 | 0.170 | 0.174 | 0.156 | 0.194 |
| MetricX | 0.391 | 0.105 | 0.077 | 0.146 | 0.100 | 0.216 |
| MetricX:mbr | 0.431 | 0.127 | 0.173 | 0.150 | 0.153 | 0.200 |
| MetricX-QE | 0.485 | 0.120 | 0.115 | 0.132 | 0.074 | 0.170 |
| CometKiwi23-XXL | 0.241 | 0.088 | 0.116 | 0.121 | 0.128 | 0.208 |
| CometKiwi23-XL | 0.284 | 0.092 | 0.098 | 0.118 | 0.124 | 0.202 |
| CometKiwi22 | 0.277 | 0.148 | 0.050 | 0.156 | 0.116 | 0.235 |
| COMET22 | 0.298 | 0.170 | 0.058 | 0.146 | 0.099 | 0.195 |
| COMET22:mbr | 0.312 | 0.209 | 0.069 | 0.149 | 0.118 | 0.190 |
| BLEURT | 0.308 | 0.152 | 0.134 | 0.143 | 0.115 | 0.205 |
| BLEURT:mbr | 0.322 | 0.170 | 0.143 | 0.150 | 0.149 | 0.191 |
| YiSi | 0.211 | 0.088 | 0.105 | 0.100 | 0.092 | 0.187 |
| YiSi:mbr | 0.214 | 0.124 | 0.138 | 0.105 | 0.100 | 0.202 |
| chrF | 0.054 | 0.055 | 0.106 | 0.106 | 0.086 | 0.164 |
| chrF:mbr | 0.083 | 0.069 | 0.111 | 0.115 | 0.084 | 0.189 |
| chrF++ | 0.062 | 0.058 | 0.109 | 0.108 | 0.087 | 0.157 |
| chrF++:mbr | 0.091 | 0.069 | 0.110 | 0.113 | 0.088 | 0.183 |
| sentBLEU | 0.128 | 0.072 | 0.095 | 0.094 | 0.073 | 0.091 |
| sentBLEU:mbr | 0.160 | 0.072 | 0.098 | 0.111 | 0.088 | 0.113 |
| TER | 0.101 | 0.071 | 0.104 | 0.061 | 0.105 | 0.118 |
| TER:mbr | 0.096 | 0.085 | 0.105 | 0.067 | 0.107 | 0.149 |
| MetricX +MetricX-QE | 0.463 | 0.124 | 0.101 | 0.152 | 0.101 | 0.229 |
| MetricX +MetricX-QE | 0.483 | 0.130 | 0.160 | 0.151 | 0.131 | 0.209 |
| XCOMET-XXL +XCOMET-XL | 0.532 | 0.110 | 0.159 | 0.161 | 0.144 | 0.228 |
| XCOMET-XXL:mbr +XCOMET-XL:mbr | 0.537 | 0.105 | 0.183 | 0.171 | 0.166 | 0.215 |
| XCOMET-XXL +XCOMET-XL +COMET22 | 0.521 | 0.136 | 0.150 | 0.169 | 0.144 | 0.235 |
| XCOMET-XXL:mbr +XCOMET-XL:mbr +COMET22:mbr | 0.529 | 0.136 | 0.173 | 0.176 | 0.167 | 0.223 |

Table 7: Pearson correlation between MQM evaluation scores and automated evaluation scores. For reference-based metrics, rows with ":mbr" indicate pseudoreference-based evaluation. Bottom rows are ensembles that take the average between the listed metrics. Higher scores indicate better agreement with human raters. See Table 6 for Kendall-Tau correlation.

# E   Results on Dev Datasets (WMT2022 and FLORES200 dev)

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.58 | 1.16 | 82.3 | 77.8 | 76.8 | 68.2 | 77.5 | 85.2 | 77.3 | 84.3 | 57.2 | 54.2 | 26.4 | 63.4 |
| rankAvg:all | 1.08‡ | 0.739‡ | 86.5‡ | 81.7‡ | 81.2‡ | 71.4‡ | 79.3‡ | 86.5‡ | 79.3‡ | 84.3 | 57.1 | 53.9 | 25.3‡ | 63.7 |
| rankAvg:qe | 1.04‡ | 0.580‡ | 86.6‡ | 81.8‡ | 83.2‡ | 73.0‡ | 80.3‡ | 85.9‡ | 77.7‡ | 82.6‡ | 52.8‡ | 49.5‡ | 20.8‡ | 70.7‡ |
| rankAvg:top | 0.899‡ | 0.566‡ | 88.2‡ | 83.0‡ | 83.0‡ | 72.7‡ | 78.9‡ | 85.8‡ | 78.1‡ | 82.5‡ | 52.8‡ | 49.5‡ | 20.7‡ | 71.0‡ |
| rankAvg:topQe | 1.00‡ | 0.527‡ | 86.8‡ | 81.7‡ | 83.7‡ | 73.3‡ | 78.9‡ | 85.6‡ | 77.5 | 82.4‡ | 52.3‡ | 48.9‡ | 20.2‡ | 71.7‡ |
| rankAvg:mxmxqe | 0.700‡ | 0.417‡ | 85.6‡ | 79.7‡ | 79.2‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.4‡ | 47.0‡ | 18.2‡ | 75.1‡ |
| rankAvg:noLex | 0.993‡ | 0.657‡ | 87.3‡ | 82.4‡ | 82.0‡ | 72.0‡ | 79.6‡ | 86.6‡ | 79.5‡ | 83.8‡ | 55.6‡ | 52.3‡ | 23.4‡ | 66.7‡ |
| rankAvg:noNC | 1.09‡ | 0.734‡ | 85.2‡ | 80.4‡ | 79.5‡ | 70.1‡ | 78.5‡ | 86.4‡ | 79.2‡ | 84.4‡ | 57.4‡ | 54.1* | 25.7‡ | 63.0* |
| rankAvg:noNCnoLex | 0.968‡ | 0.636‡ | 85.8‡ | 80.8‡ | 80.0‡ | 70.4‡ | 78.6‡ | 86.6‡ | 79.7‡ | 84.0‡ | 56.1‡ | 52.8‡ | 24.0‡ | 66.0‡ |
| rankAvg:noNCQe | 0.934‡ | 0.400‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.7‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.5‡ | 46.1‡ | 17.6‡ | 75.5‡ |
| rankMax:all | 1.16‡ | 0.776‡ | 86.1‡ | 81.0‡ | 80.8‡ | 71.1‡ | 79.2‡ | 86.3‡ | 78.9‡ | 83.9‡ | 56.1‡ | 52.8‡ | 24.3‡ | 64.1 |
| rankMax:qe | 1.06‡ | 0.595‡ | 86.3‡ | 81.5‡ | 82.8‡ | 72.6‡ | 80.2‡ | 85.9‡ | 77.7‡ | 82.7‡ | 53.0‡ | 49.6‡ | 20.9‡ | 70.5‡ |
| rankMax:top | 0.929‡ | 0.586‡ | 88.0‡ | 82.7‡ | 82.7‡ | 71.4‡ | 78.8‡ | 85.7‡ | 78.0‡ | 82.5‡ | 52.8‡ | 49.5‡ | 20.8‡ | 70.6‡ |
| rankMax:topQe | 0.964‡ | 0.480‡ | 86.7‡ | 80.7‡ | 84.0‡ | 71.2‡ | 78.6‡ | 85.4‡ | 77.0‡ | 82.1‡ | 51.7‡ | 48.3‡ | 19.7‡ | 72.0‡ |
| rankMax:mxmxqe | 0.704‡ | 0.420‡ | 85.6‡ | 79.7‡ | 79.3‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.5‡ | 47.1‡ | 18.2‡ | 75.0‡ |
| rankMax:noLex | 1.11‡ | 0.739‡ | 86.6‡ | 81.5‡ | 81.3‡ | 71.4‡ | 79.4‡ | 86.4‡ | 79.1‡ | 83.8‡ | 55.5‡ | 52.2‡ | 23.4‡ | 66.5‡ |
| rankMax:noNC | 1.11‡ | 0.733‡ | 85.1‡ | 80.1‡ | 79.3‡ | 69.9‡ | 78.4‡ | 86.3‡ | 79.1‡ | 84.0‡ | 56.3‡ | 53.1‡ | 24.7‡ | 63.6 |
| rankMax:noNCnoLex | 1.05‡ | 0.685‡ | 85.4‡ | 80.4‡ | 79.6‡ | 70.2‡ | 78.5‡ | 86.4‡ | 79.5‡ | 83.9‡ | 55.9‡ | 52.6‡ | 23.8‡ | 66.0‡ |
| rankMax:noNCQe | 0.937‡ | 0.405‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.6‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.4‡ | 46.0‡ | 17.6‡ | 75.5‡ |
| rankMed:all | 1.06‡ | 0.733‡ | 86.5‡ | 81.9‡ | 81.0‡ | 71.3‡ | 79.1‡ | 86.5‡ | 79.2‡ | 84.1‡ | 56.8‡ | 53.6‡ | 25.1‡ | 64.5* |
| rankMed:qe | 1.14‡ | 0.679‡ | 86.5‡ | 81.7‡ | 83.3‡ | 73.0‡ | 79.9‡ | 85.7‡ | 77.5 | 82.4‡ | 52.3‡ | 49.0‡ | 20.3‡ | 71.6‡ |
| rankMed:top | 0.895‡ | 0.573‡ | 88.2‡ | 83.1‡ | 82.8‡ | 72.5‡ | 78.9‡ | 85.6‡ | 77.9‡ | 82.2‡ | 52.3‡ | 48.9‡ | 20.1‡ | 71.9‡ |
| rankMed:topQe | 1.21‡ | 0.726‡ | 86.5‡ | 81.4‡ | 83.8‡ | 73.2‡ | 78.9‡ | 85.3‡ | 77.1‡ | 82.1‡ | 51.7‡ | 48.3‡ | 19.7‡ | 72.4‡ |
| rankMed:mxmxqe | 0.700‡ | 0.417‡ | 85.6‡ | 79.7‡ | 79.2‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.4‡ | 47.0‡ | 18.2‡ | 75.1‡ |
| rankMed:noLex | 0.935‡ | 0.611‡ | 87.6‡ | 82.8‡ | 82.2‡ | 72.2‡ | 79.4‡ | 86.4‡ | 79.1‡ | 83.1‡ | 54.3‡ | 51.0‡ | 22.1‡ | 69.0‡ |
| rankMed:noNC | 1.28‡ | 0.927‡ | 84.2‡ | 79.6‡ | 78.6‡ | 69.5‡ | 78.2‡ | 86.2‡ | 78.7‡ | 84.6‡ | 57.9‡ | 54.7‡ | 26.3 | 62.6* |
| rankMed:noNCnoLex | 0.910‡ | 0.607‡ | 85.8‡ | 80.9‡ | 80.0‡ | 70.4‡ | 78.6‡ | 86.5‡ | 79.3‡ | 83.5‡ | 55.1‡ | 51.8‡ | 23.0‡ | 67.9‡ |
| rankMed:noNCQe | 0.934‡ | 0.400‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.7‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.5‡ | 46.1‡ | 17.6‡ | 75.5‡ |
| rank75q:all | 1.09‡ | 0.743‡ | 86.5‡ | 81.7‡ | 81.1‡ | 71.3‡ | 79.1‡ | 86.5‡ | 79.2‡ | 84.2 | 56.9‡ | 53.6‡ | 25.0‡ | 64.2 |
| rank75q:qe | 1.06‡ | 0.600‡ | 86.5‡ | 81.7‡ | 83.2‡ | 72.9‡ | 80.0‡ | 85.9‡ | 77.6‡ | 82.6‡ | 52.7‡ | 49.4‡ | 20.7‡ | 70.9‡ |
| rank75q:top | 0.892‡ | 0.564‡ | 88.0‡ | 82.9‡ | 82.8‡ | 72.6‡ | 78.9‡ | 85.7‡ | 78.0‡ | 82.4‡ | 52.7‡ | 49.4‡ | 20.6‡ | 71.2‡ |
| rank75q:topQe | 1.00‡ | 0.526‡ | 86.7‡ | 81.7‡ | 83.6‡ | 73.3‡ | 78.9‡ | 85.6‡ | 77.5 | 82.4‡ | 52.3‡ | 49.0‡ | 20.3‡ | 71.6‡ |
| rank75q:mxmxqe | 0.705‡ | 0.419‡ | 85.6‡ | 79.7‡ | 79.2‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.5‡ | 47.0‡ | 18.2‡ | 75.1‡ |
| rank75q:noLex | 0.990‡ | 0.651‡ | 87.3‡ | 82.5‡ | 82.0‡ | 72.0‡ | 79.5‡ | 86.5‡ | 79.4‡ | 83.5‡ | 55.2‡ | 51.9‡ | 23.0‡ | 67.4‡ |
| rank75q:noNC | 1.13‡ | 0.780‡ | 85.0‡ | 80.2‡ | 79.3‡ | 69.9‡ | 78.4‡ | 86.4‡ | 79.0‡ | 84.3 | 57.3* | 54.1 | 25.6‡ | 63.3 |
| rank75q:noNCnoLex | 0.955‡ | 0.628‡ | 85.8‡ | 80.9‡ | 80.0‡ | 70.4‡ | 78.6‡ | 86.6‡ | 79.6‡ | 83.7‡ | 55.6‡ | 52.2‡ | 23.4‡ | 67.0‡ |
| rank75q:noNCQe | 0.937‡ | 0.403‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.6‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.4‡ | 46.1‡ | 17.6‡ | 75.5‡ |

Table 8: Reference-based and QE evaluation scores for greedy, MBR, and QE decoding using a single-step ensemble utility metric, averaged across all languages (test datasets). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001.

# F Results for Additional Ensembles

## F.1 Additional Single-Step Ensembles on Test Datasets

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.58 | 1.16 | 82.3 | 77.8 | 76.8 | 68.2 | 77.5 | 85.2 | 77.3 | 84.3 | 57.2 | 54.2 | 26.4 | 63.4 |
| rankAvg:all | 1.08‡ | 0.739‡ | 86.5‡ | 81.7‡ | 81.2‡ | 71.4‡ | 79.3‡ | 86.5‡ | 79.3‡ | 84.3 | 57.1 | 53.9 | 25.3‡ | 63.7 |
| rankAvg:qe | 1.04‡ | 0.580‡ | 86.6‡ | 81.8‡ | 83.2‡ | 73.0‡ | 80.3‡ | 85.9‡ | 77.7‡ | 82.6‡ | 52.8‡ | 49.5‡ | 20.8‡ | 70.7‡ |
| rankAvg:top | 0.899‡ | 0.566‡ | 88.2‡ | 83.0‡ | 83.0‡ | 72.7‡ | 78.9‡ | 85.8‡ | 78.1‡ | 82.5‡ | 52.8‡ | 49.5‡ | 20.7‡ | 71.0‡ |
| rankAvg:topQe | 1.00‡ | 0.527‡ | 86.8‡ | 81.7‡ | 83.7‡ | 73.3‡ | 78.9‡ | 85.6‡ | 77.5 | 82.4‡ | 52.3‡ | 48.9‡ | 20.2‡ | 71.7‡ |
| rankAvg:mxmxqe | 0.700‡ | 0.417‡ | 85.6‡ | 79.7‡ | 79.2‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.4‡ | 47.0‡ | 18.2‡ | 75.1‡ |
| rankAvg:noLex | 0.993‡ | 0.657‡ | 87.3‡ | 82.4‡ | 82.0‡ | 72.0‡ | 79.6‡ | 86.6‡ | 79.5‡ | 83.8‡ | 55.6‡ | 52.3‡ | 23.4‡ | 66.7‡ |
| rankAvg:noNC | 1.09‡ | 0.734‡ | 85.2‡ | 80.4‡ | 79.5‡ | 70.1‡ | 78.5‡ | 86.4‡ | 79.2‡ | 84.4‡ | 57.4‡ | 54.1* | 25.7‡ | 63.0* |
| rankAvg:noNCnoLex | 0.968‡ | 0.636‡ | 85.8‡ | 80.8‡ | 80.0‡ | 70.4‡ | 78.6‡ | 86.6‡ | 79.7‡ | 84.0‡ | 56.1‡ | 52.8‡ | 24.0‡ | 66.0‡ |
| rankAvg:noNCQe | 0.934‡ | 0.400‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.7‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.5‡ | 46.1‡ | 17.6‡ | 75.5‡ |
| rankMax:all | 1.16‡ | 0.776‡ | 86.1‡ | 81.0‡ | 80.8‡ | 71.1‡ | 79.2‡ | 86.3‡ | 78.9‡ | 83.9‡ | 56.1‡ | 52.8‡ | 24.3‡ | 64.1 |
| rankMax:qe | 1.06‡ | 0.595‡ | 86.3‡ | 81.5‡ | 82.8‡ | 72.6‡ | 80.2‡ | 85.9‡ | 77.7‡ | 82.7‡ | 53.0‡ | 49.6‡ | 20.9‡ | 70.5‡ |
| rankMax:top | 0.929‡ | 0.586‡ | 88.0‡ | 82.7‡ | 82.7‡ | 71.4‡ | 78.8‡ | 85.7‡ | 78.0‡ | 82.5‡ | 52.8‡ | 49.5‡ | 20.8‡ | 70.6‡ |
| rankMax:topQe | 0.964‡ | 0.480‡ | 86.7‡ | 80.7‡ | 84.0‡ | 71.2‡ | 78.6‡ | 85.4‡ | 77.0‡ | 82.1‡ | 51.7‡ | 48.3‡ | 19.7‡ | 72.0‡ |
| rankMax:mxmxqe | 0.704‡ | 0.420‡ | 85.6‡ | 79.7‡ | 79.3‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.5‡ | 47.1‡ | 18.2‡ | 75.0‡ |
| rankMax:noLex | 1.11‡ | 0.739‡ | 86.6‡ | 81.5‡ | 81.3‡ | 71.4‡ | 79.4‡ | 86.4‡ | 79.1‡ | 83.8‡ | 55.5‡ | 52.2‡ | 23.4‡ | 66.5‡ |
| rankMax:noNC | 1.11‡ | 0.733‡ | 85.1‡ | 80.1‡ | 79.3‡ | 69.9‡ | 78.4‡ | 86.3‡ | 79.1‡ | 84.0‡ | 56.3‡ | 53.1‡ | 24.7‡ | 63.6 |
| rankMax:noNCnoLex | 1.05‡ | 0.685‡ | 85.4‡ | 80.4‡ | 79.6‡ | 70.2‡ | 78.5‡ | 86.4‡ | 79.5‡ | 83.9‡ | 55.9‡ | 52.6‡ | 23.8‡ | 66.0‡ |
| rankMax:noNCQe | 0.937‡ | 0.405‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.6‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.4‡ | 46.0‡ | 17.6‡ | 75.5‡ |
| rankMed:all | 1.06‡ | 0.733‡ | 86.5‡ | 81.9‡ | 81.0‡ | 71.3‡ | 79.1‡ | 86.5‡ | 79.2‡ | 84.1‡ | 56.8‡ | 53.6‡ | 25.1‡ | 64.5* |
| rankMed:qe | 1.14‡ | 0.679‡ | 86.5‡ | 81.7‡ | 83.3‡ | 73.0‡ | 79.9‡ | 85.7‡ | 77.5 | 82.4‡ | 52.3‡ | 49.0‡ | 20.3‡ | 71.6‡ |
| rankMed:top | 0.895‡ | 0.573‡ | 88.2‡ | 83.1‡ | 82.8‡ | 72.5‡ | 78.9‡ | 85.6‡ | 77.9‡ | 82.2‡ | 52.3‡ | 48.9‡ | 20.1‡ | 71.9‡ |
| rankMed:topQe | 1.21‡ | 0.726‡ | 86.5‡ | 81.4‡ | 83.8‡ | 73.2‡ | 78.9‡ | 85.3‡ | 77.1‡ | 82.1‡ | 51.7‡ | 48.3‡ | 19.7‡ | 72.4‡ |
| rankMed:mxmxqe | 0.700‡ | 0.417‡ | 85.6‡ | 79.7‡ | 79.2‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.4‡ | 47.0‡ | 18.2‡ | 75.1‡ |
| rankMed:noLex | 0.935‡ | 0.611‡ | 87.6‡ | 82.8‡ | 82.2‡ | 72.2‡ | 79.4‡ | 86.4‡ | 79.1‡ | 83.1‡ | 54.3‡ | 51.0‡ | 22.1‡ | 69.0‡ |
| rankMed:noNC | 1.28‡ | 0.927‡ | 84.2‡ | 79.6‡ | 78.6‡ | 69.5‡ | 78.2‡ | 86.2‡ | 78.7‡ | 84.6‡ | 57.9‡ | 54.7‡ | 26.3 | 62.6* |
| rankMed:noNCnoLex | 0.910‡ | 0.607‡ | 85.8‡ | 80.9‡ | 80.0‡ | 70.4‡ | 78.6‡ | 86.5‡ | 79.3‡ | 83.5‡ | 55.1‡ | 51.8‡ | 23.0‡ | 67.9‡ |
| rankMed:noNCQe | 0.934‡ | 0.400‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.7‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.5‡ | 46.1‡ | 17.6‡ | 75.5‡ |
| rank75q:all | 1.09‡ | 0.743‡ | 86.5‡ | 81.7‡ | 81.1‡ | 71.3‡ | 79.1‡ | 86.5‡ | 79.2‡ | 84.2 | 56.9‡ | 53.6‡ | 25.0‡ | 64.2 |
| rank75q:qe | 1.06‡ | 0.600‡ | 86.5‡ | 81.7‡ | 83.2‡ | 72.9‡ | 80.0‡ | 85.9‡ | 77.6‡ | 82.6‡ | 52.7‡ | 49.4‡ | 20.7‡ | 70.9‡ |
| rank75q:top | 0.892‡ | 0.564‡ | 88.0‡ | 82.9‡ | 82.8‡ | 72.6‡ | 78.9‡ | 85.7‡ | 78.0‡ | 82.4‡ | 52.7‡ | 49.4‡ | 20.6‡ | 71.2‡ |
| rank75q:topQe | 1.00‡ | 0.526‡ | 86.7‡ | 81.7‡ | 83.6‡ | 73.3‡ | 78.9‡ | 85.6‡ | 77.5 | 82.4‡ | 52.3‡ | 49.0‡ | 20.3‡ | 71.6‡ |
| rank75q:mxmxqe | 0.705‡ | 0.419‡ | 85.6‡ | 79.7‡ | 79.2‡ | 69.6‡ | 77.8‡ | 84.9‡ | 76.7‡ | 81.3‡ | 50.5‡ | 47.0‡ | 18.2‡ | 75.1‡ |
| rank75q:noLex | 0.990‡ | 0.651‡ | 87.3‡ | 82.5‡ | 82.0‡ | 72.0‡ | 79.5‡ | 86.5‡ | 79.4‡ | 83.5‡ | 55.2‡ | 51.9‡ | 23.0‡ | 67.4‡ |
| rank75q:noNC | 1.13‡ | 0.780‡ | 85.0‡ | 80.2‡ | 79.3‡ | 69.9‡ | 78.4‡ | 86.4‡ | 79.0‡ | 84.3‡ | 57.3* | 54.1 | 25.6‡ | 63.3 |
| rank75q:noNCnoLex | 0.955‡ | 0.628‡ | 85.8‡ | 80.9‡ | 80.0‡ | 70.4‡ | 78.6‡ | 86.6‡ | 79.6‡ | 83.7‡ | 55.6‡ | 52.2‡ | 23.4‡ | 67.0‡ |
| rank75q:noNCQe | 0.937‡ | 0.403‡ | 84.5‡ | 78.3‡ | 78.5‡ | 69.0‡ | 77.6‡ | 84.6‡ | 75.6‡ | 81.1‡ | 49.4‡ | 46.1‡ | 17.6‡ | 75.5‡ |

Table 9: Reference-based and QE evaluation scores for greedy, MBR, and QE decoding using a single-step ensemble utility metric, averaged across all languages (test datasets). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for $p<0.05$, † for $p<0.01$, ‡ for $p<0.001$.

## F.2 Additional Two-Step Ensembles on Test Datasets

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.58 | 1.16 | 82.3 | 77.8 | 76.8 | 68.2 | 77.5 | 85.2 | 77.3 | 84.3 | 57.2 | 54.2 | 26.4 | 63.4 |
| allQE(64)allMBR | 1.09‡ | 0.781‡ | 86.5‡ | 81.7‡ | 80.6‡ | 71.0‡ | 78.9‡ | 86.5‡ | 79.3‡ | 84.3 | 57.1 | 53.9 | 25.4‡ | 63.6 |
| allQE(32)allMBR | 1.06‡ | 0.733‡ | 86.7‡ | 81.9‡ | 81.3‡ | 71.4‡ | 79.2‡ | 86.5‡ | 79.2‡ | 84.1‡ | 56.6‡ | 53.4‡ | 24.9‡ | 64.5 |
| allQE(16)allMBR | 1.04‡ | 0.688‡ | 86.8‡ | 82.0‡ | 81.8‡ | 71.9‡ | 79.4‡ | 86.4‡ | 79.1‡ | 83.9‡ | 56.2‡ | 52.9‡ | 24.3‡ | 65.4‡ |
| allQE(8)allMBR | 1.04‡ | 0.654‡ | 86.8‡ | 82.0‡ | 82.2‡ | 72.2‡ | 79.7‡ | 86.3‡ | 78.7‡ | 83.6‡ | 55.4‡ | 52.1‡ | 23.4‡ | 66.8‡ |
| allQE(4)allMBR | 1.04‡ | 0.629‡ | 86.8‡ | 82.0‡ | 82.7‡ | 72.5‡ | 79.9‡ | 86.2‡ | 78.4‡ | 83.2‡ | 54.5‡ | 51.2‡ | 22.3‡ | 68.3‡ |
| allQE(64)nolexMBR | 0.991‡ | 0.708‡ | 87.4‡ | 82.4‡ | 81.1‡ | 71.3‡ | 79.0‡ | 86.6‡ | 79.7‡ | 83.9‡ | 55.9‡ | 52.6‡ | 23.9‡ | 66.0‡ |
| allQE(32)nolexMBR | 0.978‡ | 0.680‡ | 87.5‡ | 82.6‡ | 81.6‡ | 71.7‡ | 79.2‡ | 86.6‡ | 79.5‡ | 83.7‡ | 55.6‡ | 52.3‡ | 23.6‡ | 66.6‡ |
| allQE(16)nolexMBR | 0.972‡ | 0.647‡ | 87.5‡ | 82.6‡ | 82.1‡ | 72.0‡ | 79.4‡ | 86.5‡ | 79.3‡ | 83.5‡ | 55.2‡ | 51.9‡ | 23.2‡ | 67.2‡ |
| allQE(8)nolexMBR | 0.977‡ | 0.625‡ | 87.3‡ | 82.5‡ | 82.4‡ | 72.3‡ | 79.7‡ | 86.4‡ | 79.0‡ | 83.3‡ | 54.6‡ | 51.3‡ | 22.5‡ | 68.3‡ |
| allQE(4)nolexMBR | 0.988‡ | 0.608‡ | 87.2‡ | 82.4‡ | 82.8‡ | 72.6‡ | 79.9‡ | 86.2‡ | 78.6‡ | 83.0‡ | 53.9‡ | 50.5‡ | 21.7‡ | 69.4‡ |
| topQE(64)topMBR | 0.868‡ | 0.621‡ | 88.5‡ | 83.3‡ | 81.5‡ | 71.5‡ | 78.7‡ | 85.6‡ | 78.1‡ | 82.3‡ | 52.4‡ | 49.1‡ | 20.4‡ | 71.2‡ |
| topQE(32)topMBR | 0.861‡ | 0.599‡ | 88.4‡ | 83.3‡ | 82.0‡ | 71.9‡ | 78.8‡ | 85.7‡ | 78.1‡ | 82.4‡ | 52.7‡ | 49.4‡ | 20.7‡ | 70.9‡ |
| topQE(16)topMBR | 0.879‡ | 0.585‡ | 88.3‡ | 83.2‡ | 82.4‡ | 72.2‡ | 78.9‡ | 85.7‡ | 78.1‡ | 82.5‡ | 52.8‡ | 49.4‡ | 20.8‡ | 70.8‡ |
| topQE(8)topMBR | 0.897‡ | 0.567‡ | 88.1‡ | 82.9‡ | 82.8‡ | 72.6‡ | 78.9‡ | 85.7‡ | 78.0‡ | 82.5‡ | 52.8‡ | 49.5‡ | 20.7‡ | 71.0‡ |
| topQE(4)topMBR | 0.925‡ | 0.548‡ | 87.7‡ | 82.6‡ | 83.2‡ | 72.9‡ | 78.9‡ | 85.7‡ | 77.8‡ | 82.4‡ | 52.6‡ | 49.2‡ | 20.5‡ | 71.3‡ |
| noncQE(64)noncnolexMBR | 0.955‡ | 0.668‡ | 85.9‡ | 81.0‡ | 80.0‡ | 70.4‡ | 78.7‡ | 86.6‡ | 79.8‡ | 83.9‡ | 55.9‡ | 52.6‡ | 23.8‡ | 66.3‡ |
| noncQE(32)noncnolexMBR | 0.911‡ | 0.596‡ | 86.0‡ | 81.0‡ | 80.1‡ | 70.4‡ | 78.7‡ | 86.5‡ | 79.4‡ | 83.6‡ | 55.1‡ | 51.7‡ | 22.9‡ | 67.5‡ |
| noncQE(16)noncnolexMBR | 0.883‡ | 0.533‡ | 86.0‡ | 80.8‡ | 80.0‡ | 70.3‡ | 78.6‡ | 86.2‡ | 78.8‡ | 83.1‡ | 54.1‡ | 50.7‡ | 21.8‡ | 69.2‡ |
| noncQE(8)noncnolexMBR | 0.877‡ | 0.487‡ | 85.8‡ | 80.5‡ | 79.8‡ | 70.2‡ | 78.4‡ | 85.9‡ | 78.2‡ | 82.6‡ | 53.0‡ | 49.6‡ | 20.7‡ | 70.7‡ |
| noncQE(4)noncnolexMBR | 0.890‡ | 0.450‡ | 85.4‡ | 79.8‡ | 79.5‡ | 69.9‡ | 78.2‡ | 85.5‡ | 77.4 | 82.0‡ | 51.5‡ | 48.1‡ | 19.4‡ | 72.9‡ |
| noncQE(64)noncMBR | 1.06‡ | 0.728‡ | 85.3‡ | 80.6‡ | 79.6‡ | 70.1‡ | 78.5‡ | 86.4‡ | 79.2‡ | 84.3 | 57.0 | 53.8‡ | 25.3‡ | 63.7 |
| noncQE(32)noncMBR | 0.992‡ | 0.629‡ | 85.6‡ | 80.6‡ | 79.8‡ | 70.2‡ | 78.5‡ | 86.3‡ | 78.9‡ | 83.9‡ | 56.1‡ | 52.8‡ | 24.2‡ | 65.2‡ |
| noncQE(16)noncMBR | 0.960‡ | 0.559‡ | 85.6‡ | 80.4‡ | 79.7‡ | 70.1‡ | 78.5‡ | 86.0‡ | 78.4‡ | 83.5‡ | 54.9‡ | 51.7‡ | 22.9‡ | 67.1‡ |
| noncQE(8)noncMBR | 0.942‡ | 0.506‡ | 85.5‡ | 80.1‡ | 79.6‡ | 69.9‡ | 78.3‡ | 85.8‡ | 77.8‡ | 82.9‡ | 53.7‡ | 50.3‡ | 21.5‡ | 69.2‡ |
| noncQE(4)noncMBR | 0.931‡ | 0.461‡ | 85.2‡ | 79.5‡ | 79.3‡ | 69.7‡ | 78.1‡ | 85.4‡ | 77.0‡ | 82.2‡ | 52.1‡ | 48.7‡ | 19.9‡ | 71.8‡ |
| mxQE(64)xcMBR | 1.11‡ | 0.690‡ | 89.8‡ | 80.6‡ | 80.9‡ | 70.1‡ | 78.2‡ | 85.1 | 76.9‡ | 81.7‡ | 50.7‡ | 47.3‡ | 18.8‡ | 73.1‡ |
| mxQE(32)xcMBR | 1.03‡ | 0.593‡ | 89.5‡ | 80.6‡ | 80.9‡ | 70.1‡ | 78.2‡ | 85.1 | 76.9‡ | 81.7‡ | 50.7‡ | 47.4‡ | 18.8‡ | 73.1‡ |
| mxQE(16)xcMBR | 0.965‡ | 0.517‡ | 89.1‡ | 80.5‡ | 80.7‡ | 70.0‡ | 78.2‡ | 85.1* | 76.9‡ | 81.6‡ | 50.6‡ | 47.2‡ | 18.7‡ | 73.4‡ |
| mxQE(8)xcMBR | 0.924‡ | 0.459‡ | 88.4‡ | 80.3‡ | 80.3‡ | 69.9‡ | 78.1‡ | 85.0‡ | 76.7‡ | 81.6‡ | 50.4‡ | 47.1‡ | 18.6‡ | 73.4‡ |
| mxQE(4)xcMBR | 0.904‡ | 0.411‡ | 87.5‡ | 79.8‡ | 79.9‡ | 69.8‡ | 78.0‡ | 84.9‡ | 76.4‡ | 81.4‡ | 50.1‡ | 46.8‡ | 18.3‡ | 73.9‡ |
| ckQE(64)xcMBR | 1.23‡ | 0.851‡ | 89.8‡ | 80.7‡ | 81.9‡ | 70.4‡ | 78.3‡ | 85.1 | 76.9‡ | 81.8‡ | 50.9‡ | 47.6‡ | 19.1‡ | 72.8‡ |
| ckQE(32)xcMBR | 1.24‡ | 0.847‡ | 89.6‡ | 80.8‡ | 82.8‡ | 70.7‡ | 78.4‡ | 85.2 | 77.0‡ | 81.9‡ | 51.3‡ | 48.0‡ | 19.5‡ | 72.2‡ |
| ckQE(16)xcMBR | 1.25‡ | 0.850‡ | 89.3‡ | 81.0‡ | 83.5‡ | 71.0‡ | 78.6‡ | 85.3‡ | 77.1‡ | 82.1‡ | 51.6‡ | 48.3‡ | 19.9‡ | 71.6‡ |
| ckQE(8)xcMBR | 1.30‡ | 0.870‡ | 88.9‡ | 80.9‡ | 84.1‡ | 71.2‡ | 78.7‡ | 85.3‡ | 77.0‡ | 82.2‡ | 51.8‡ | 48.5‡ | 19.9‡ | 71.5‡ |
| ckQE(4)xcMBR | 1.33‡ | 0.883‡ | 88.3‡ | 80.8‡ | 84.7‡ | 71.4‡ | 78.7‡ | 85.3‡ | 76.9‡ | 82.2‡ | 51.8‡ | 48.5‡ | 20.0‡ | 71.5‡ |
| mxQE(64)mxMBR | 0.653‡ | 0.508‡ | 85.6‡ | 79.8‡ | 79.2‡ | 69.5‡ | 77.8‡ | 85.0‡ | 76.8‡ | 81.4‡ | 50.6‡ | 47.2‡ | 18.4‡ | 75.2‡ |
| mxQE(32)mxMBR | 0.662‡ | 0.475‡ | 85.6‡ | 79.8‡ | 79.2‡ | 69.5‡ | 77.8‡ | 85.0‡ | 76.8‡ | 81.5‡ | 50.7‡ | 47.3‡ | 18.5‡ | 74.9‡ |
| mxQE(16)mxMBR | 0.681‡ | 0.450‡ | 85.5‡ | 79.6‡ | 79.1‡ | 69.4‡ | 77.8‡ | 85.0‡ | 76.7‡ | 81.5‡ | 50.5‡ | 47.1‡ | 18.4‡ | 74.9‡ |
| mxQE(8)mxMBR | 0.712‡ | 0.421‡ | 85.3‡ | 79.5‡ | 79.0‡ | 69.4‡ | 77.8‡ | 84.9‡ | 76.4‡ | 81.4‡ | 50.3‡ | 46.9‡ | 18.3‡ | 74.9‡ |
| mxQE(4)mxMBR | 0.762‡ | 0.395‡ | 85.0‡ | 79.0‡ | 78.8‡ | 69.3‡ | 77.7‡ | 84.7‡ | 76.2‡ | 81.3‡ | 50.1‡ | 46.7‡ | 18.1‡ | 75.1‡ |
| ckQE(64)mxMBR | 0.687‡ | 0.553‡ | 86.1‡ | 80.3‡ | 81.0‡ | 70.2‡ | 78.1‡ | 85.2 | 77.1‡ | 81.7‡ | 51.2‡ | 47.7‡ | 19.0‡ | 74.2‡ |
| ckQE(32)mxMBR | 0.728‡ | 0.557‡ | 86.5‡ | 80.6‡ | 82.2‡ | 70.7‡ | 78.3‡ | 85.4‡ | 77.3 | 81.9‡ | 51.7‡ | 48.3‡ | 19.5‡ | 73.3‡ |
| ckQE(16)mxMBR | 0.798‡ | 0.594‡ | 86.8‡ | 80.9‡ | 83.2‡ | 71.1‡ | 78.5‡ | 85.4‡ | 77.4 | 82.1‡ | 51.9‡ | 48.5‡ | 19.8‡ | 72.7‡ |
| ckQE(8)mxMBR | 0.892‡ | 0.644‡ | 87.0‡ | 81.0‡ | 84.0‡ | 71.3‡ | 78.7‡ | 85.5‡ | 77.4 | 82.2‡ | 52.1‡ | 48.8‡ | 20.1‡ | 72.0‡ |
| ckQE(4)mxMBR | 1.01‡ | 0.714‡ | 86.9‡ | 80.9‡ | 84.6‡ | 71.4‡ | 78.7‡ | 85.4‡ | 77.2† | 82.2‡ | 52.0‡ | 48.7‡ | 20.0‡ | 71.9‡ |

Table 10: Reference-based and QE evaluation scores for greedy, MBR, and QE decoding using a two-step ensemble (QE filtering followed by MBR) utility metric, averaged across all languages (test datasets). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001.

# G  Breakdown of Results on Individual Language Pairs

## G.1  Results for English-Swahili (en-sw) on FLORES200 test dataset

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | AfriCOMET | AfriCOMET-QE | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.70 | 1.28 | 83.7 | 85.0 | 84.6 | 73.2 | 83.8 | 86.0 | 85.7 | 76.9 | 77.5 | 86.3 | 65.0 | 62.6 | 34.9 | 51.7 |
| MetricX | **0.598‡** | **0.477‡** | **88.9‡** | **87.7‡** | **89.0‡** | **75.5‡** | **84.9‡** | 85.9 | **87.4‡** | **79.7‡** | **76.3‡** | **83.0‡** | **58.1‡** | **55.1‡** | **24.4‡** | **61.1*** |
| MetricX-QE | **0.811‡** | **0.293‡** | **87.5‡** | **86.6‡** | **88.4‡** | **75.0‡** | **84.7‡** | **85.2‡** | **86.6‡** | **79.3‡** | **75.3‡** | **82.7‡** | **57.2‡** | **54.1‡** | **23.8‡** | **61.4*** |
| XCOMET-XXL | **1.03‡** | **0.698‡** | **94.2‡** | **89.1‡** | **91.3‡** | **76.4‡** | **85.3‡** | **86.3*** | **87.6‡** | **79.5‡** | **76.9*** | **83.5‡** | **58.8‡** | **56.0‡** | **25.6‡** | 59.3 |
| XCOMET-XL | **1.08‡** | **0.788‡** | **89.9‡** | **92.4‡** | **89.9‡** | **77.9‡** | **85.4‡** | **86.4†** | **87.9‡** | **79.5‡** | 77.8 | **83.7‡** | **59.3‡** | **56.4‡** | **26.1‡** | 58.5 |
| CometKiwi23-XXL | **1.23‡** | **0.784‡** | **90.6‡** | **88.3‡** | **93.6‡** | **76.9‡** | **85.5‡** | 86.0 | **87.1‡** | **79.4‡** | **76.6‡** | **83.5‡** | **58.7‡** | **55.8‡** | **25.6‡** | 58.8 |
| CometKiwi23-XL | **1.26‡** | **0.816‡** | **88.3‡** | **90.2‡** | **90.1‡** | **79.9‡** | **85.5‡** | 86.2 | **87.3‡** | **79.3‡** | 77.1 | **83.7‡** | **59.5‡** | **56.7‡** | **26.4‡** | 58.4 |
| CometKiwi22 | **1.26‡** | **0.852‡** | **87.6‡** | **87.8‡** | **89.1‡** | **76.4‡** | **87.3‡** | **86.7‡** | **87.9‡** | **79.8‡** | 77.0 | **83.7‡** | **59.4‡** | **56.4‡** | **25.8‡** | 58.3 |
| COMET22 | **1.25‡** | **0.927‡** | **87.5‡** | **88.0‡** | **88.4‡** | **75.9‡** | **85.4‡** | **88.3‡** | **87.8‡** | **79.7‡** | **78.5‡** | **85.2‡** | **62.8‡** | **60.0‡** | **30.2‡** | 52.7 |
| AfriCOMET | **1.10‡** | **0.769‡** | **88.7‡** | **88.4‡** | **89.2‡** | **75.9‡** | **85.7‡** | **86.7‡** | **90.0‡** | **80.7‡** | 77.5 | **84.2‡** | **60.7‡** | **57.8‡** | **27.5‡** | 56.3 |
| AfriCOMET-QE | **1.42‡** | **0.964‡** | **85.1‡** | 85.3 | **87.1‡** | **74.5‡** | **84.6‡** | 85.8 | **87.6‡** | **83.3‡** | **74.6‡** | **82.9‡** | **57.6‡** | **54.5‡** | **23.5‡** | **61.8*** |
| BLEURT | **1.37‡** | **1.05‡** | **86.3‡** | **86.9‡** | **87.3‡** | **75.0‡** | **84.8‡** | 86.2 | **86.7‡** | **78.5‡** | **82.9‡** | **84.0‡** | **60.0‡** | **57.0‡** | **25.8‡** | 58.0 |
| YiSi | 1.62 | 1.26 | 84.2 | 85.3 | 84.9 | 73.4 | 83.9 | **86.2*** | 85.7 | 76.9 | 77.9 | **86.9‡** | **65.7†** | **63.2†** | 35.1 | 46.7 |
| chrF | **1.57†** | 1.23 | **84.7‡** | **85.7‡** | **85.5‡** | **74.0‡** | **84.1*** | **86.5‡** | **86.2‡** | **77.3‡** | **78.3‡** | 86.4 | **66.4‡** | **63.7‡** | **34.3*** | 49.3 |
| chrF++ | **1.57†** | 1.22 | **84.8‡** | **85.8‡** | **85.5‡** | **74.0‡** | **84.1*** | **86.5‡** | **86.2‡** | **77.2‡** | **78.4‡** | **86.5‡** | **66.4‡** | **63.9‡** | 34.7 | 48.8 |
| sentBLEU | 1.64 | **1.29** | 84.1 | 85.4 | 84.7 | 73.4 | 83.9 | 86.1 | 85.7 | 76.8 | 77.6 | **86.5*** | 65.3 | 63.0 | **35.8†** | 46.3 |
| TER | 1.73 | **1.36*** | 83.2 | **84.2‡** | **83.7†** | **72.7*** | **83.6*** | 85.8 | **85.0‡** | **76.3‡** | 77.3 | 86.2 | **64.4†** | **62.0†** | 34.6 | 45.1 |
| | | | | | | | | | | | | | | | | |
| rankAvg:all | **1.01‡** | **0.711‡** | **89.8‡** | **89.7‡** | **90.1‡** | **77.1‡** | **85.9‡** | **87.5‡** | **88.5‡** | **79.8‡** | **79.4‡** | **85.9‡** | **64.5*** | **61.9†** | **33.0‡** | 49.7 |
| rankAvg:qe | **0.893‡** | **0.506‡** | **90.3‡** | **89.8‡** | **91.9‡** | **78.3‡** | **86.5‡** | **86.9‡** | **88.5‡** | **81.3‡** | 77.4 | **83.8‡** | **59.7‡** | **56.8‡** | **26.2‡** | 58.6 |
| rankAvg:top | **0.781‡** | **0.484‡** | **92.2‡** | **90.9‡** | **92.0‡** | **78.4‡** | **85.8‡** | **86.7‡** | **88.2‡** | **80.0‡** | 77.8 | **83.8‡** | **59.8‡** | **57.0‡** | **26.7‡** | 58.3 |
| rankAvg:topQe | **0.900‡** | **0.455‡** | **90.7‡** | **89.8‡** | **92.3‡** | **78.7‡** | **85.7‡** | **86.4†** | **87.9‡** | **79.9‡** | 77.4 | **83.5‡** | **59.0‡** | **56.0‡** | **25.7‡** | 59.4 |
| rankAvg:mxmxqe | **0.638‡** | **0.347‡** | **88.8‡** | **87.4‡** | **89.0‡** | **75.7‡** | **85.0‡** | 85.8 | **87.3‡** | **79.7‡** | **76.2‡** | **83.0‡** | **58.2‡** | **55.1‡** | **24.4‡** | **60.8*** |
| rankAvg:noLex | **0.899‡** | **0.606‡** | **91.2‡** | **90.5‡** | **91.0‡** | **77.7‡** | **86.2‡** | **87.5‡** | **88.9‡** | **80.4‡** | **79.7‡** | **85.2‡** | **62.7‡** | **60.0‡** | **30.2‡** | 53.4 |
| rankAvg:noNC | **1.06‡** | **0.724‡** | **88.0‡** | **88.4‡** | **88.4‡** | **75.8‡** | **85.3‡** | **87.4‡** | **88.3‡** | **79.7‡** | **79.5‡** | 86.2 | 64.8 | 62.3 | **33.7‡** | 48.6 |
| rankAvg:noNCnoLex | **0.919‡** | **0.597‡** | **89.2‡** | **89.0‡** | **89.4‡** | **76.5‡** | **85.7‡** | **87.6‡** | **88.7‡** | **80.6‡** | **80.1‡** | **85.7‡** | **63.6‡** | **60.9‡** | **31.4‡** | 51.8 |
| allQE(32)allMBR | **0.992‡** | **0.705‡** | **90.2‡** | **89.8‡** | **90.9‡** | **77.1‡** | **85.8‡** | **87.4‡** | **88.5‡** | **79.8‡** | **79.4‡** | **85.8‡** | **64.1‡** | **61.5‡** | **32.6‡** | 50.4 |
| allQE(32)nolexMBR | **0.904‡** | **0.636‡** | **91.2‡** | **90.6‡** | **90.7‡** | **77.5‡** | **86.0‡** | **87.5‡** | **88.8‡** | **80.1‡** | **79.5‡** | **85.2‡** | **62.7‡** | **60.0‡** | **30.6‡** | 53.3 |
| topQE(32)topMBR | **0.761‡** | **0.552‡** | **92.4‡** | **90.9‡** | **91.2‡** | **77.7‡** | **85.7‡** | **86.6‡** | **88.1‡** | **79.9‡** | 78.0 | **83.7‡** | **59.6‡** | **56.8‡** | **26.9‡** | 58.1 |
| noncQE(32)noncMBR | **0.968‡** | **0.648‡** | **88.7‡** | **88.7‡** | **89.0‡** | **76.1‡** | **85.5‡** | **87.3‡** | **88.5‡** | **80.1‡** | **79.4‡** | **85.8‡** | **64.3†** | **61.7‡** | **32.8‡** | 49.9 |
| noncQE(32)noncnolexMBR | **0.885‡** | **0.604‡** | **89.3‡** | **89.1‡** | **89.5‡** | **76.5‡** | **85.7‡** | **87.5‡** | **88.8‡** | **80.4‡** | **79.7‡** | **85.4‡** | **63.1‡** | **60.3‡** | **30.7‡** | 52.7 |
| mxQE(32)mxMBR | **0.628‡** | **0.434‡** | **89.2‡** | **87.7‡** | **88.9‡** | **75.6‡** | **85.0‡** | 86.0 | **87.4‡** | **79.7‡** | **76.5‡** | **83.2‡** | **58.5‡** | **55.5‡** | **25.0‡** | **60.3*** |
| ckQE(32)xcMBR | **1.05‡** | **0.696‡** | **94.0‡** | **89.1‡** | **91.8‡** | **76.6‡** | **85.4‡** | **86.3*** | **87.6‡** | **79.6‡** | 77.2 | **83.6‡** | **59.1‡** | **56.2‡** | **25.9‡** | 59.1 |
| mxQE(32)xcMBR | **0.928‡** | **0.538‡** | **93.6‡** | **89.0‡** | **91.2‡** | **76.5‡** | **85.4‡** | 86.2 | **87.6‡** | **79.7‡** | **77.0*** | **83.5‡** | **58.7‡** | **55.8‡** | **25.4‡** | 59.4 |
| ckQE(32)mxMBR | **0.657‡** | **0.499‡** | **90.0‡** | **88.3‡** | **90.9‡** | **76.2‡** | **85.3‡** | 86.2 | **87.5‡** | **79.8‡** | **76.8*** | **83.4‡** | **58.7‡** | **55.8‡** | **25.4‡** | 59.6 |

Table 11: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-sw (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

1081

## G.2 Results for English-Hausa (en-ha) on FLORES200 test dataset

| MBR/QE Method | Evaluated Metric | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | AfriCOMET | AfriCOMET-QE | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | | 2.31 | 1.46 | 74.7 | 70.3 | 79.9 | 66.2 | 60.1 | 81.5 | 79.3 | 72.3 | 84.1 | 79.5 | 52.5 | 49.8 | 21.3 | 65.9 |
| MetricX | | 0.818‡ | 0.515‡ | 76.7‡ | 70.5 | 82.8‡ | 67.5‡ | 60.0 | 81.5 | 81.3‡ | 76.2‡ | 81.2‡ | 77.0‡ | 47.5‡ | 44.3‡ | 14.3‡ | 80.2‡ |
| MetricX-QE | | 1.19‡ | 0.278‡ | 76.5‡ | 69.6 | 82.4‡ | 67.1† | 59.5 | 81.0† | 80.7‡ | 75.6‡ | 80.0‡ | 76.9‡ | 46.6‡ | 43.3‡ | 13.7‡ | 78.7‡ |
| XCOMET-XXL | | 1.65‡ | 0.909‡ | 87.5‡ | 73.9‡ | 86.9‡ | 69.8‡ | 61.4‡ | 82.5‡ | 81.5‡ | 75.4‡ | 82.6‡ | 77.7‡ | 48.6‡ | 45.5‡ | 15.3‡ | 75.9‡ |
| XCOMET-XL | | 1.75‡ | 0.990‡ | 80.9‡ | 80.6‡ | 85.1‡ | 71.8‡ | 62.8‡ | 82.5‡ | 81.7‡ | 75.4‡ | 83.2† | 77.8‡ | 48.8‡ | 45.8‡ | 15.9‡ | 74.2‡ |
| CometKiwi23-XXL | | 1.97‡ | 0.989‡ | 82.2‡ | 72.9‡ | 90.2‡ | 70.5‡ | 62.2‡ | 82.2‡ | 81.1‡ | 75.6‡ | 81.8‡ | 77.6‡ | 48.0‡ | 45.0‡ | 15.4‡ | 75.4‡ |
| CometKiwi23-XL | | 1.99‡ | 1.12‡ | 78.5‡ | 75.4‡ | 85.4‡ | 74.5‡ | 62.6‡ | 82.3‡ | 80.6‡ | 74.9‡ | 82.0‡ | 77.7‡ | 48.3‡ | 45.3‡ | 15.5‡ | 76.1‡ |
| CometKiwi22 | | 2.55† | 1.51 | 73.5* | 69.0† | 81.2‡ | 67.8‡ | 74.8‡ | 82.0† | 78.8* | 73.1‡ | 81.6‡ | 77.5‡ | 47.5‡ | 44.4‡ | 14.8‡ | 75.6‡ |
| COMET22 | | 1.94‡ | 1.14‡ | 77.7‡ | 72.4‡ | 84.0‡ | 69.0‡ | 63.4‡ | 84.9‡ | 81.6‡ | 75.9‡ | 83.6 | 78.4‡ | 50.0‡ | 46.9‡ | 16.4‡ | 74.3‡ |
| AfriCOMET | | 1.71‡ | 0.956‡ | 77.8‡ | 71.3* | 83.2‡ | 68.0‡ | 60.2 | 82.0* | 85.4‡ | 77.5‡ | 81.8‡ | 77.8‡ | 48.9‡ | 45.7‡ | 15.7‡ | 75.6‡ |
| AfriCOMET-QE | | 1.93‡ | 1.07‡ | 73.4* | 67.6‡ | 81.3‡ | 66.5 | 59.0‡ | 81.3 | 82.2‡ | 80.3‡ | 78.9‡ | 76.8‡ | 46.6‡ | 43.3‡ | 13.4‡ | 81.1‡ |
| BLEURT | | 2.10‡ | 1.34* | 75.8* | 71.2* | 81.3‡ | 67.5‡ | 61.4‡ | 82.2‡ | 80.2‡ | 73.7‡ | 89.8‡ | 78.3‡ | 49.9‡ | 46.9‡ | 16.8‡ | 73.3‡ |
| YiSi | | 2.37 | 1.46 | 73.7* | 69.2† | 79.4 | 66.2 | 60.5 | 81.7 | 79.1 | 72.5 | 83.6 | 80.4‡ | 52.3 | 49.6 | 20.2‡ | 66.6 |
| chrF | | 2.30 | 1.40 | 74.7 | 69.7 | 80.3 | 67.0‡ | 60.6* | 82.1‡ | 80.0† | 73.5‡ | 84.3 | 79.4 | 53.7‡ | 50.6‡ | 19.4‡ | 71.9‡ |
| chrF++ | | 2.34 | 1.43 | 74.4 | 69.9 | 80.2 | 66.9† | 60.7* | 82.1‡ | 79.9† | 73.4‡ | 84.4 | 79.5 | 53.5‡ | 50.6‡ | 19.8‡ | 71.4‡ |
| sentBLEU | | 2.36 | 1.50 | 73.6† | 69.8 | 78.8‡ | 65.7 | 60.1 | 81.5 | 79.0 | 72.3 | 83.7 | 79.6 | 52.3 | 49.7 | 21.2 | 65.1* |
| TER | | 2.66‡ | 1.69‡ | 72.8‡ | 68.8‡ | 77.3‡ | 64.4‡ | 59.6* | 80.8‡ | 77.9‡ | 71.3‡ | 82.9‡ | 79.6 | 51.3‡ | 48.8‡ | 21.1 | 61.3‡ |
| | | | | | | | | | | | | | | | | | |
| rankAvg:all | | 1.47‡ | 0.782‡ | 81.8‡ | 75.8‡ | 85.9‡ | 70.8‡ | 64.8‡ | 83.6‡ | 82.9‡ | 76.1‡ | 85.7‡ | 79.3 | 51.9* | 49.0‡ | 19.3‡ | 68.8‡ |
| rankAvg:qe | | 1.40‡ | 0.581‡ | 81.1‡ | 75.1‡ | 87.5‡ | 71.9‡ | 68.1‡ | 83.0‡ | 82.7‡ | 77.5‡ | 82.8‡ | 77.7‡ | 48.9‡ | 45.7‡ | 15.5‡ | 76.7‡ |
| rankAvg:top | | 1.15‡ | 0.524‡ | 84.1‡ | 77.5‡ | 88.0‡ | 72.4‡ | 62.4‡ | 82.6‡ | 82.4‡ | 76.3‡ | 83.3* | 77.8‡ | 49.1‡ | 46.0‡ | 16.1‡ | 76.0‡ |
| rankAvg:topQe | | 1.36‡ | 0.491‡ | 81.8‡ | 75.8‡ | 88.4‡ | 72.8‡ | 62.7‡ | 82.5‡ | 82.0‡ | 76.0‡ | 82.5‡ | 77.6‡ | 48.4‡ | 45.4‡ | 15.7‡ | 75.9‡ |
| rankAvg:mxmxqe | | 0.891‡ | 0.336‡ | 77.4‡ | 70.9 | 82.8‡ | 67.6‡ | 60.1 | 81.4 | 81.5‡ | 76.1‡ | 81.1‡ | 77.0‡ | 47.3‡ | 44.1‡ | 14.2‡ | 79.4‡ |
| rankAvg:noLex | | 1.34‡ | 0.657‡ | 82.8‡ | 76.4‡ | 86.9‡ | 71.6‡ | 65.8‡ | 83.7‡ | 83.3‡ | 76.9‡ | 85.1‡ | 78.7‡ | 50.5‡ | 47.5‡ | 17.5‡ | 72.8‡ |
| rankAvg:noNC | | 1.43‡ | 0.736‡ | 79.4‡ | 73.4‡ | 83.6‡ | 68.9‡ | 62.0‡ | 83.4‡ | 82.9‡ | 76.2‡ | 85.5‡ | 79.5 | 52.1 | 49.2† | 19.8‡ | 68.0‡ |
| rankAvg:noNCnoLex | | 1.27‡ | 0.625‡ | 80.1‡ | 73.9‡ | 84.4‡ | 69.3‡ | 62.3‡ | 83.7‡ | 83.5‡ | 77.2‡ | 85.5‡ | 79.0‡ | 51.0‡ | 48.0‡ | 17.9‡ | 72.4‡ |
| allQE(32)allMBR | | 1.45‡ | 0.802‡ | 82.0‡ | 76.1‡ | 85.8‡ | 70.6‡ | 63.7‡ | 83.6‡ | 83.0‡ | 76.1‡ | 85.6‡ | 79.1‡ | 51.6‡ | 48.7‡ | 19.1‡ | 70.0‡ |
| allQE(32)nolexMBR | | 1.31‡ | 0.715‡ | 83.4‡ | 76.9‡ | 86.4‡ | 71.2‡ | 63.7‡ | 83.7‡ | 83.4‡ | 76.6‡ | 85.3‡ | 78.6‡ | 50.5‡ | 47.6‡ | 17.7‡ | 73.1‡ |
| topQE(32)topMBR | | 1.12‡ | 0.595‡ | 84.6‡ | 78.0‡ | 87.0‡ | 71.6‡ | 62.4‡ | 82.6‡ | 82.4‡ | 76.1‡ | 83.5 | 77.8‡ | 49.1‡ | 46.1‡ | 16.2‡ | 75.6‡ |
| noncQE(32)noncMBR | | 1.35‡ | 0.704‡ | 79.4‡ | 73.6‡ | 83.7‡ | 69.0‡ | 61.9‡ | 83.4‡ | 83.0‡ | 76.5‡ | 84.9* | 79.0‡ | 51.5‡ | 48.5‡ | 18.7‡ | 70.2‡ |
| noncQE(32)noncnolexMBR | | 1.22‡ | 0.653‡ | 80.4‡ | 74.3‡ | 84.6‡ | 69.3‡ | 62.3‡ | 83.7‡ | 83.6‡ | 76.9‡ | 85.1† | 78.8‡ | 50.6‡ | 47.6‡ | 17.7‡ | 73.0‡ |
| mxQE(32)mxMBR | | 0.859‡ | 0.444‡ | 76.7‡ | 70.7 | 82.8‡ | 67.6‡ | 60.3 | 81.6 | 81.3‡ | 76.1‡ | 81.1‡ | 77.1‡ | 47.7‡ | 44.5‡ | 14.6‡ | 79.5‡ |
| ckQE(32)xcMBR | | 1.60‡ | 0.865‡ | 87.2‡ | 74.2‡ | 87.7‡ | 70.1‡ | 61.3‡ | 82.4‡ | 81.6‡ | 75.7‡ | 82.4‡ | 77.7‡ | 48.6‡ | 45.5‡ | 15.3‡ | 76.6‡ |
| mxQE(32)xcMBR | | 1.39‡ | 0.613‡ | 86.4‡ | 73.9‡ | 86.9‡ | 69.9‡ | 61.5‡ | 82.3‡ | 81.9‡ | 75.8‡ | 82.2‡ | 77.6‡ | 48.5‡ | 45.3‡ | 15.3‡ | 76.8‡ |
| ckQE(32)mxMBR | | 0.908‡ | 0.516‡ | 79.6‡ | 72.8‡ | 86.6‡ | 69.3‡ | 61.1‡ | 82.1‡ | 82.0‡ | 76.3‡ | 81.7‡ | 77.3‡ | 48.0‡ | 44.9‡ | 14.7‡ | 78.9‡ |

Table 12: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-ha (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

## G.3 Results for English-Igbo (en-ig) on FLORES200 test dataset

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | AfriCOMET | AfriCOMET-QE | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 3.73 | 2.67 | 27.3 | 19.8 | 16.9 | 16.8 | 30.2 | 72.1 | 68.9 | 66.4 | 33.5 | 78.3 | 42.7 | 40.7 | 15.2 | 74.2 |
| MetricX | **1.68‡** | **1.31‡** | **26.7‡** | **18.9‡** | **18.0‡** | **17.8‡** | 30.3 | 72.1 | **71.7‡** | **69.7‡** | 28.8‡ | 77.2‡ | 40.7‡ | 38.6‡ | 12.8‡ | 79.6‡ |
| MetricX-QE | **2.15‡** | **0.937‡** | **26.7‡** | **19.2*** | 17.4 | **17.7‡** | **30.6*** | 72.1 | **71.2‡** | **69.3‡** | 29.0‡ | 77.3‡ | 40.7‡ | 38.5‡ | 12.8‡ | 78.8‡ |
| XCOMET-XXL | 5.12‡ | 3.55‡ | **31.4‡** | **21.0‡** | **20.3‡** | **18.9‡** | **32.1‡** | 72.2 | 66.5‡ | 65.6‡ | 31.2‡ | 77.5‡ | 40.1‡ | 38.0‡ | 12.8‡ | 74.1 |
| XCOMET-XL | 6.81‡ | 5.12‡ | 27.2 | **26.7‡** | **20.0‡** | **19.6‡** | **33.1‡** | 69.1‡ | 62.6‡ | 64.3‡ | 29.2‡ | 75.9‡ | 37.2‡ | 35.2‡ | 11.8‡ | 79.5‡ |
| CometKiwi23-XXL | 6.08‡ | 4.17‡ | 26.4‡ | 20.1 | **36.6‡** | **21.7‡** | **32.3‡** | 70.2‡ | 64.5‡ | 64.9‡ | 27.1‡ | 76.0‡ | 37.9‡ | 35.8‡ | 11.4‡ | 83.2‡ |
| CometKiwi23-XL | 5.94‡ | 3.99‡ | 26.2‡ | **20.7†** | **24.8‡** | **29.3‡** | **32.6‡** | 70.2‡ | 64.5‡ | 65.5‡ | 27.4‡ | 75.8‡ | 38.0‡ | 35.7‡ | 10.9‡ | 83.6‡ |
| CometKiwi22 | 6.13‡ | 4.33‡ | **26.7‡** | **21.9‡** | **21.0‡** | **20.1‡** | **40.6‡** | 69.6‡ | 63.9‡ | 65.3‡ | 28.8‡ | 75.9‡ | 37.0‡ | 35.0‡ | 11.3‡ | 77.6‡ |
| COMET22 | 5.28‡ | 3.80‡ | 27.2 | **18.5‡** | 17.4 | **17.3†** | **30.7†** | **76.3‡** | 66.7‡ | 66.1 | 29.1‡ | 77.7‡ | 41.0‡ | 38.9‡ | 13.0‡ | 78.9‡ |
| AfriCOMET | **3.26‡** | **2.26‡** | **26.8†** | **19.0‡** | **18.4‡** | **17.6‡** | **30.6*** | 72.3 | **77.3‡** | **72.6‡** | 28.8‡ | 77.4‡ | 40.9‡ | 38.6‡ | 12.7‡ | 79.7‡ |
| AfriCOMET-QE | 4.16‡ | **2.87*** | 26.4‡ | **19.0†** | **17.9†** | **18.0‡** | **31.0‡** | 71.8 | **73.4‡** | **75.3‡** | 24.7‡ | 76.6‡ | 38.8‡ | 36.6‡ | 11.4‡ | 81.8‡ |
| BLEURT | 4.20‡ | 3.05‡ | 27.2 | 19.9 | **19.0‡** | **17.6‡** | **30.7†** | 71.9 | 69.0 | 66.3 | **40.9‡** | 77.8‡ | 41.8‡ | 39.8‡ | **14.1‡** | 74.9 |
| YiSi | 4.03‡ | **2.87†** | **27.9‡** | **19.2*** | 17.1 | **17.3‡** | 30.0 | 73.2‡ | 69.6† | 66.9‡ | 33.9 | **79.2‡** | **43.5‡** | **41.5‡** | 15.3 | 73.7 |
| chrF | **4.01†** | **2.86†** | 27.3 | **18.8‡** | **17.7†** | **17.8‡** | **29.9†** | 73.0‡ | 69.9‡ | 66.9‡ | 33.1 | **78.5*** | **44.5‡** | **42.3‡** | 15.3 | 77.2‡ |
| chrF++ | **3.95*** | **2.84†** | 27.3 | **18.9‡** | **17.6‡** | **17.7‡** | **30.0*** | 72.8‡ | 69.7‡ | 66.8* | 33.7 | **78.5‡** | **44.4‡** | **42.3‡** | 15.4 | 76.8‡ |
| sentBLEU | 4.05‡ | 2.94‡ | **27.7*** | 20.0 | **18.2‡** | **17.5‡** | **30.9‡** | 72.2 | 68.8 | 66.1 | 34.6† | 78.3 | 42.9 | **41.1*** | **15.8‡** | **72.4†** |
| TER | 4.30‡ | 3.07‡ | **28.2‡** | **20.7‡** | **18.1‡** | **17.2*** | **31.5‡** | **71.7†** | 68.2† | 65.9† | 34.8‡ | 78.2 | 42.1† | **40.3*** | 15.3 | **69.5‡** |
| | | | | | | | | | | | | | | | | |
| rankAvg:all | **3.05‡** | **2.05‡** | **28.6‡** | **20.8‡** | **22.1‡** | **20.6‡** | **32.5‡** | **73.6‡** | **72.4‡** | **69.3‡** | 34.8† | **78.5*** | 43.2† | **41.1*** | 15.4 | 73.0 |
| rankAvg:qe | **3.34‡** | **1.94‡** | 27.5 | **21.1‡** | **27.2‡** | **24.1‡** | **35.1‡** | 71.6† | **71.7‡** | **71.2‡** | 29.5‡ | 76.9‡ | 40.2‡ | 37.9‡ | 12.3‡ | 78.2‡ |
| rankAvg:top | **2.86‡** | **1.77‡** | **29.1‡** | **22.0‡** | **26.2‡** | **23.1‡** | **32.5‡** | 72.0 | 70.4‡ | 68.4‡ | 31.6‡ | 77.4‡ | 41.1‡ | 38.9‡ | 13.4‡ | **75.5*** |
| rankAvg:topQe | **3.17‡** | **1.70‡** | 27.2 | **20.4*** | **28.9‡** | **24.8‡** | **32.0‡** | 71.6† | 69.8‡ | 68.4‡ | 29.4‡ | 76.8‡ | 40.4‡ | 38.1‡ | 12.4‡ | 80.1‡ |
| rankAvg:mxmxqe | **1.83‡** | **1.03‡** | **26.8†** | **19.2*** | **17.9‡** | **17.9‡** | 30.4 | 72.1 | **71.9‡** | **69.6‡** | 29.5‡ | 77.3‡ | 41.0‡ | 38.8‡ | 12.9‡ | 79.0‡ |
| rankAvg:noLex | **2.98‡** | **1.93‡** | **28.8‡** | **21.2‡** | **23.4‡** | **21.3‡** | **33.1‡** | **73.8‡** | **72.9‡** | **70.0‡** | 34.1 | 78.3 | 42.7 | 40.5 | 14.8 | 73.7 |
| rankAvg:noNC | **2.74‡** | **1.85‡** | **28.0‡** | 19.7 | **18.5‡** | **18.1‡** | **30.8‡** | **73.6‡** | **73.1‡** | **69.7‡** | 34.6† | **78.6‡** | **43.4‡** | **41.3‡** | 15.4 | 73.4 |
| rankAvg:noNCnoLex | **2.60‡** | **1.67‡** | **27.7‡** | 19.5 | **17.9‡** | **18.0‡** | **30.7‡** | **74.1‡** | **73.9‡** | **70.8‡** | 34.6† | **78.6†** | 42.9 | 40.7 | **14.7*** | 74.9 |
| allQE(32)allMBR | **3.06‡** | **2.06‡** | **28.6‡** | **20.9‡** | **21.9‡** | **20.4‡** | **32.5‡** | **73.2‡** | **72.3‡** | **69.3‡** | 34.3 | 78.3 | 42.9 | 40.8 | 15.1 | 73.3 |
| allQE(32)nolexMBR | **2.90‡** | **1.97‡** | **28.8‡** | **21.3‡** | **22.1‡** | **20.3‡** | **32.5‡** | **73.5‡** | **72.9‡** | **69.7‡** | 34.6* | 78.3 | 42.7 | 40.5 | **14.7*** | 73.6 |
| topQE(32)topMBR | **2.71‡** | **1.80‡** | **29.3‡** | **22.3‡** | **23.5‡** | **21.5‡** | **32.6‡** | 71.9 | 70.6‡ | 68.4‡ | 32.3† | 77.6‡ | 41.3‡ | 39.1‡ | 13.7‡ | 74.3 |
| noncQE(32)noncMBR | **2.64‡** | **1.76‡** | **27.8†** | 19.7 | **18.4‡** | **18.0‡** | **30.7‡** | **73.3‡** | **73.6‡** | **70.1‡** | 34.2 | **78.6‡** | **43.4‡** | **41.3‡** | 15.1 | 73.9 |
| noncQE(32)noncnolexMBR | **2.51‡** | **1.69‡** | **27.9‡** | 19.6 | **18.1‡** | **17.9‡** | **30.7‡** | **73.9‡** | **74.2‡** | **70.6‡** | 34.3 | **78.5‡** | 43.0 | 40.8 | **14.6*** | 75.0 |
| mxQE(32)mxMBR | **1.73‡** | **1.24‡** | **26.7‡** | **19.0†** | **18.0‡** | **17.6‡** | 30.4 | 72.1 | **71.8‡** | **69.6‡** | 29.4‡ | 77.3‡ | 40.9‡ | 38.8‡ | 13.0‡ | 78.9‡ |
| ckQE(32)xcMBR | 5.07‡ | 3.53‡ | **30.4‡** | **21.3‡** | **26.6‡** | **20.4‡** | **32.4‡** | 71.9 | 66.6‡ | 66.0* | 30.2‡ | 77.2‡ | 39.8‡ | 37.7‡ | 12.7‡ | **75.8*** |
| mxQE(32)xcMBR | **3.41‡** | **1.90‡** | **30.5‡** | **21.0‡** | **20.6‡** | **18.7‡** | **32.0‡** | 72.5* | 69.3 | 67.6‡ | 31.6‡ | 77.7‡ | 41.0‡ | 38.8‡ | 13.3‡ | 74.2 |
| ckQE(32)mxMBR | **2.11‡** | **1.56‡** | 27.0 | 19.4 | **25.4‡** | **19.6‡** | **30.9‡** | 72.0 | **71.2‡** | **69.2‡** | 29.0‡ | 77.1‡ | 40.7‡ | 38.5‡ | 12.7‡ | 80.2‡ |

Table 13: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-ig (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

## G.4 Results for English-Somali (en-so) on FLORES200 test dataset

| MBR/QE Method | Evaluated Metric | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | AfriCOMET | AfriCOMET-QE | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | | 2.66 | 1.89 | 67.9 | 66.0 | 78.3 | 65.8 | 70.0 | 80.7 | 75.9 | 73.4 | 108. | 77.7 | 46.7 | 42.4 | 11.3 | 86.7 |
| MetricX | | 0.996‡ | 0.635‡ | 71.5‡ | 69.2‡ | 83.8‡ | 69.5‡ | 70.5* | 81.9‡ | 79.8‡ | 78.4‡ | 110.‡ | 76.2‡ | 44.2‡ | 39.5‡ | 8.89‡ | 91.9 |
| MetricX-QE | | 1.34‡ | 0.396‡ | 70.6‡ | 68.3‡ | 83.0‡ | 68.6‡ | 70.3 | 81.5‡ | 79.0‡ | 77.8‡ | 109.† | 76.3‡ | 43.8‡ | 39.2‡ | 8.94‡ | 89.2 |
| XCOMET-XXL | | 1.69‡ | 0.998‡ | 81.0‡ | 72.3‡ | 87.2‡ | 70.5‡ | 71.5‡ | 82.1‡ | 79.8‡ | 78.0‡ | 110.‡ | 76.5‡ | 44.2‡ | 39.7‡ | 9.14‡ | 89.0 |
| XCOMET-XL | | 1.82‡ | 1.08‡ | 74.7‡ | 78.7‡ | 85.4‡ | 72.7‡ | 72.6‡ | 82.5‡ | 80.1‡ | 77.7‡ | 111.‡ | 76.8‡ | 45.0‡ | 40.5‡ | 9.85‡ | 88.2 |
| CometKiwi23-XXL | | 1.89‡ | 1.09‡ | 75.6‡ | 71.6‡ | 90.4‡ | 71.4‡ | 72.3‡ | 82.1‡ | 79.6‡ | 78.0‡ | 110.‡ | 76.4‡ | 44.1‡ | 39.6‡ | 9.34‡ | 89.0 |
| CometKiwi23-XL | | 2.00‡ | 1.17‡ | 72.8‡ | 74.6‡ | 85.9‡ | 75.1‡ | 72.7‡ | 82.2‡ | 79.5‡ | 77.6‡ | 110.‡ | 76.5‡ | 44.4‡ | 39.9‡ | 9.64‡ | 89.8 |
| CometKiwi22 | | 2.56 | 1.65† | 68.0 | 66.9* | 81.4‡ | 68.4‡ | 80.3‡ | 81.7‡ | 77.1‡ | 75.7‡ | 109.* | 76.7‡ | 44.3‡ | 39.7‡ | 9.49‡ | 87.8 |
| COMET22 | | 2.01‡ | 1.30‡ | 72.1‡ | 69.9‡ | 83.5‡ | 69.6‡ | 73.0‡ | 84.6‡ | 79.8‡ | 78.3‡ | 111.‡ | 77.2‡ | 45.8‡ | 41.2‡ | 9.93‡ | 87.7 |
| AfriCOMET | | 1.78‡ | 1.08‡ | 73.7‡ | 71.2‡ | 84.7‡ | 70.0‡ | 71.5‡ | 82.4‡ | 83.9‡ | 80.2‡ | 110.‡ | 76.8‡ | 45.1‡ | 40.5‡ | 9.66‡ | 88.5 |
| AfriCOMET-QE | | 2.01‡ | 1.25‡ | 69.3† | 67.3‡ | 82.1‡ | 68.6‡ | 70.2 | 82.2‡ | 81.1‡ | 82.6‡ | 108. | 76.1‡ | 43.9‡ | 39.1‡ | 8.55‡ | 92.1 |
| BLEURT | | 2.23‡ | 1.52‡ | 68.4 | 67.5† | 81.7‡ | 68.2‡ | 71.0‡ | 81.5‡ | 77.6‡ | 76.6‡ | 120.‡ | 76.0‡ | 43.9‡ | 39.2‡ | 8.29‡ | 96.9‡ |
| YiSi | | 2.66 | 1.71* | 68.0 | 65.8 | 79.4† | 66.4* | 70.5* | 81.3‡ | 76.6† | 74.3‡ | 108. | 78.6‡ | 47.2* | 42.7 | 11.2 | 82.9 |
| chrF | | 2.48* | 1.66† | 68.2 | 66.4 | 80.1‡ | 67.1‡ | 70.3 | 81.5‡ | 77.4‡ | 74.9‡ | 109.‡ | 77.9 | 48.3‡ | 43.6‡ | 11.2 | 87.1 |
| chrF++ | | 2.52 | 1.67† | 68.5 | 66.5 | 79.9‡ | 66.9‡ | 70.4* | 81.4‡ | 77.1‡ | 74.6‡ | 109.‡ | 77.9* | 48.2‡ | 43.7‡ | 11.2 | 87.1 |
| sentBLEU | | 2.65 | 1.76 | 67.9 | 66.4 | 79.0 | 66.2 | 69.8 | 80.8 | 76.3 | 73.7 | 108. | 77.8 | 46.9 | 42.6 | 11.8† | 81.9 |
| TER | | 2.85* | 1.95 | 67.9 | 65.2 | 77.4* | 65.1† | 69.8 | 80.5 | 75.3* | 73.0* | 106.‡ | 77.7 | 45.9‡ | 41.7† | 11.6 | 77.1‡ |
| | | | | | | | | | | | | | | | | | |
| rankAvg:all | | 1.63‡ | 0.909‡ | 75.8‡ | 73.5‡ | 85.9‡ | 71.4‡ | 74.2‡ | 83.1‡ | 81.1‡ | 78.4‡ | 112.‡ | 77.8 | 47.0 | 42.5 | 11.0 | 84.1 |
| rankAvg:qe | | 1.55‡ | 0.724‡ | 74.9‡ | 73.5‡ | 87.8‡ | 72.7‡ | 76.4‡ | 82.9‡ | 81.1‡ | 80.0‡ | 111.‡ | 76.7‡ | 44.7‡ | 40.1‡ | 9.39‡ | 89.5 |
| rankAvg:top | | 1.34‡ | 0.658‡ | 77.7‡ | 75.6‡ | 88.1‡ | 73.1‡ | 72.5‡ | 82.7‡ | 81.0‡ | 78.7‡ | 112.‡ | 76.6‡ | 44.8‡ | 40.3‡ | 9.54‡ | 89.5 |
| rankAvg:topQe | | 1.42‡ | 0.623‡ | 75.8‡ | 74.0‡ | 88.6‡ | 73.4‡ | 72.3‡ | 82.6‡ | 80.6‡ | 78.6‡ | 111.‡ | 76.6‡ | 44.7‡ | 40.2‡ | 9.60‡ | 89.6 |
| rankAvg:mxmxqe | | 1.08‡ | 0.458‡ | 71.8‡ | 69.4‡ | 83.9‡ | 69.4‡ | 70.7† | 81.9‡ | 79.7‡ | 78.3‡ | 110.‡ | 76.3‡ | 44.1‡ | 39.4‡ | 8.78‡ | 90.9 |
| rankAvg:noLex | | 1.45‡ | 0.786‡ | 76.8‡ | 74.7‡ | 87.1‡ | 72.3‡ | 75.0‡ | 83.3‡ | 81.7‡ | 79.3‡ | 113.‡ | 77.4* | 46.0† | 41.4‡ | 10.1‡ | 87.7 |
| rankAvg:noNC | | 1.66‡ | 0.928‡ | 73.5‡ | 71.1‡ | 83.9‡ | 69.8‡ | 72.3‡ | 82.9‡ | 80.8‡ | 78.3‡ | 112.‡ | 77.9 | 47.2* | 42.7 | 11.2 | 83.8 |
| rankAvg:noNCnoLex | | 1.41‡ | 0.764‡ | 74.4‡ | 71.9‡ | 85.1‡ | 70.5‡ | 72.5‡ | 83.3‡ | 81.7‡ | 79.5‡ | 114.‡ | 77.6 | 46.4 | 41.8* | 10.2‡ | 87.9 |
| allQE(32)allMBR | | 1.55‡ | 0.903‡ | 76.1‡ | 74.0‡ | 85.8‡ | 71.4‡ | 73.2‡ | 83.1‡ | 81.2‡ | 78.4‡ | 112.‡ | 77.5 | 46.6 | 42.1 | 10.7† | 85.7 |
| allQE(32)nolexMBR | | 1.43‡ | 0.841‡ | 77.4‡ | 75.0‡ | 86.8‡ | 71.9‡ | 73.5‡ | 83.3‡ | 81.9‡ | 79.1‡ | 114.‡ | 77.3† | 46.1† | 41.5‡ | 10.1‡ | 87.9 |
| topQE(32)topMBR | | 1.25‡ | 0.740‡ | 78.4‡ | 76.1‡ | 87.2‡ | 72.3‡ | 72.2‡ | 82.5‡ | 80.9‡ | 78.5‡ | 111.‡ | 76.6‡ | 44.7‡ | 40.2‡ | 9.68‡ | 89.4 |
| noncQE(32)noncMBR | | 1.54‡ | 0.887‡ | 74.4‡ | 71.7‡ | 84.5‡ | 70.1‡ | 72.3‡ | 83.0‡ | 81.2‡ | 78.8‡ | 112.‡ | 77.6 | 46.7 | 42.2 | 10.7† | 85.5 |
| noncQE(32)noncnolexMBR | | 1.41‡ | 0.809‡ | 74.7‡ | 72.1‡ | 85.2‡ | 70.7‡ | 72.7‡ | 83.3‡ | 81.9‡ | 79.3‡ | 114.‡ | 77.4* | 46.2* | 41.6‡ | 9.95‡ | 88.4 |
| mxQE(32)mxMBR | | 1.03‡ | 0.590‡ | 71.6‡ | 69.1‡ | 83.9‡ | 69.3‡ | 70.5* | 81.9‡ | 79.9‡ | 78.3‡ | 110.‡ | 76.3‡ | 44.3‡ | 39.6‡ | 9.03‡ | 91.2 |
| ckQE(32)xcMBR | | 1.68‡ | 0.974‡ | 80.4‡ | 72.2‡ | 88.0‡ | 70.9‡ | 71.8‡ | 82.2‡ | 79.9‡ | 78.1‡ | 110.‡ | 76.3‡ | 44.0‡ | 39.6‡ | 9.26‡ | 89.4 |
| mxQE(32)xcMBR | | 1.51‡ | 0.770‡ | 80.3‡ | 72.4‡ | 87.1‡ | 70.7‡ | 71.5‡ | 82.2‡ | 80.1‡ | 78.2‡ | 110.‡ | 76.5‡ | 44.4‡ | 39.9‡ | 9.40‡ | 89.3 |
| ckQE(32)mxMBR | | 1.09‡ | 0.650‡ | 73.8‡ | 70.6‡ | 86.6‡ | 70.3‡ | 71.1‡ | 82.1‡ | 80.2‡ | 78.6‡ | 111.‡ | 76.4‡ | 44.4‡ | 39.8‡ | 8.90‡ | 91.4 |

Table 14: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-so (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

## G.5 Results for English-Hindi (en-hi) on FLORES200 test dataset

| MBR/QE Method | Evaluated Metric | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | IndicCOMET | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | | 0.815 | 0.455 | 84.9 | 75.9 | 77.5 | 68.7 | 85.5 | 82.4 | 80.5 | 74.3 | 86.6 | 59.8 | 57.8 | 32.9 | 51.4 |
| MetricX | | 0.257‡ | 0.0928‡ | 92.1‡ | 81.4‡ | 78.7‡ | 68.9 | 86.0‡ | 82.0† | 80.8 | 73.1‡ | 83.1‡ | 50.9‡ | 48.7‡ | 21.5‡ | 66.8‡ |
| MetricX-QE | | 0.463‡ | 0.0413‡ | 89.6‡ | 77.3‡ | 76.9* | 67.6‡ | 85.6 | 81.1‡ | 79.8‡ | 72.1‡ | 82.9‡ | 49.4‡ | 47.2‡ | 20.7‡ | 65.9‡ |
| XCOMET-XXL | | 0.462‡ | 0.158‡ | 96.4‡ | 80.5‡ | 79.5‡ | 68.5 | 85.9‡ | 81.6‡ | 80.2 | 72.8‡ | 83.0‡ | 49.8‡ | 47.6‡ | 20.8‡ | 65.5‡ |
| XCOMET-XL | | 0.455‡ | 0.176‡ | 91.9‡ | 87.7‡ | 79.4‡ | 70.1‡ | 86.3‡ | 82.6 | 81.2‡ | 73.9† | 83.8‡ | 52.3‡ | 50.2‡ | 23.5‡ | 62.8‡ |
| CometKiwi23-XXL | | 0.531‡ | 0.195‡ | 91.7‡ | 80.7‡ | 84.3‡ | 70.4‡ | 86.4‡ | 82.3 | 80.6 | 73.7‡ | 84.2‡ | 53.0‡ | 50.8‡ | 23.8‡ | 61.7‡ |
| CometKiwi23-XL | | 0.590‡ | 0.238‡ | 89.4‡ | 81.4‡ | 80.3‡ | 73.3‡ | 86.4‡ | 82.4 | 80.4 | 73.5‡ | 84.3‡ | 53.1‡ | 50.9‡ | 23.8‡ | 61.7‡ |
| CometKiwi22 | | 0.581‡ | 0.202‡ | 90.1‡ | 80.5‡ | 79.7‡ | 69.9‡ | 87.6‡ | 82.8† | 81.1‡ | 74.0* | 83.9‡ | 52.5‡ | 50.4‡ | 23.2‡ | 62.8‡ |
| COMET22 | | 0.563‡ | 0.244‡ | 89.3‡ | 80.4‡ | 79.2‡ | 69.8‡ | 86.5‡ | 84.7‡ | 81.9‡ | 75.0‡ | 85.6‡ | 57.1‡ | 55.0‡ | 28.9‡ | 55.6‡ |
| IndicCOMET | | 0.641‡ | 0.306‡ | 88.7‡ | 78.4‡ | 77.8 | 68.4* | 86.0‡ | 82.6 | 85.4‡ | 73.7‡ | 84.4‡ | 53.7‡ | 51.5‡ | 24.8‡ | 59.7‡ |
| BLEURT | | 0.579‡ | 0.259‡ | 89.1‡ | 80.1‡ | 79.0‡ | 69.4‡ | 86.3‡ | 83.1‡ | 81.2‡ | 76.6‡ | 85.3‡ | 56.5‡ | 54.5‡ | 28.3‡ | 55.7‡ |
| YiSi | | 0.772* | 0.403† | 85.7* | 76.4 | 77.9‡ | 69.0* | 85.8‡ | 82.8‡ | 80.6 | 74.8‡ | 86.9‡ | 59.8 | 57.8 | 32.7 | 51.4 |
| chrF | | 0.746‡ | 0.397‡ | 85.9‡ | 76.6* | 78.0‡ | 69.1‡ | 85.8‡ | 82.9‡ | 80.8* | 74.8‡ | 86.8 | 60.6‡ | 58.5‡ | 32.6 | 52.4† |
| chrF++ | | 0.752† | 0.404† | 85.9‡ | 76.6† | 78.0† | 69.0† | 85.8‡ | 82.9‡ | 80.7 | 74.9‡ | 86.8* | 60.6‡ | 58.7‡ | 33.1 | 52.0 |
| sentBLEU | | 0.779 | 0.419* | 85.2 | 76.1 | 77.4 | 68.6 | 85.5 | 82.6* | 80.5 | 74.5 | 86.6 | 59.8 | 57.9 | 33.1 | 50.9 |
| TER | | 0.803 | 0.431 | 85.4 | 75.8 | 77.0‡ | 68.1‡ | 85.3* | 82.4 | 80.4 | 74.3 | 86.5 | 58.8‡ | 56.9‡ | 32.1† | 50.0‡ |
| | | | | | | | | | | | | | | | | |
| rankAvg:all | | 0.477‡ | 0.175‡ | 91.4‡ | 82.2‡ | 80.7‡ | 70.7‡ | 86.7‡ | 83.9‡ | 82.4‡ | 75.5‡ | 86.3‡ | 58.7‡ | 56.7‡ | 30.7‡ | 53.1‡ |
| rankAvg:qe | | 0.442‡ | 0.0972‡ | 91.8‡ | 83.1‡ | 82.5‡ | 72.0‡ | 87.1‡ | 83.0‡ | 81.4‡ | 74.3 | 84.4‡ | 53.8‡ | 51.5‡ | 24.3‡ | 61.5‡ |
| rankAvg:top | | 0.349‡ | 0.0980‡ | 94.2‡ | 85.0‡ | 82.3‡ | 71.7‡ | 86.7‡ | 82.9‡ | 81.5‡ | 74.2 | 84.2‡ | 53.4‡ | 51.1‡ | 23.9‡ | 62.4‡ |
| rankAvg:topQe | | 0.442‡ | 0.0875‡ | 92.1‡ | 82.6‡ | 82.9‡ | 72.2‡ | 86.6‡ | 82.7* | 81.2‡ | 73.9* | 84.3‡ | 53.3‡ | 51.0‡ | 24.0‡ | 62.0‡ |
| rankAvg:mxmxqe | | 0.257‡ | 0.0579‡ | 92.1‡ | 81.4‡ | 78.8‡ | 68.9 | 86.0‡ | 82.0† | 80.8 | 73.1‡ | 83.0‡ | 50.7‡ | 48.4‡ | 21.5‡ | 66.7‡ |
| rankAvg:noLex | | 0.420‡ | 0.141‡ | 92.8‡ | 83.7‡ | 81.5‡ | 71.2‡ | 86.8‡ | 83.9‡ | 82.8‡ | 75.5‡ | 85.8‡ | 57.3‡ | 55.2‡ | 28.7‡ | 55.8‡ |
| rankAvg:noNC | | 0.491‡ | 0.186‡ | 89.8‡ | 80.4‡ | 79.5‡ | 69.8‡ | 86.3‡ | 83.8‡ | 82.4‡ | 75.5‡ | 86.4* | 59.2* | 57.2† | 31.5‡ | 52.3† |
| rankAvg:noNCnoLex | | 0.431‡ | 0.144‡ | 90.7‡ | 81.7‡ | 79.9‡ | 70.1‡ | 86.5‡ | 84.0‡ | 83.0‡ | 75.6‡ | 86.1‡ | 57.9‡ | 55.8‡ | 29.7‡ | 54.6‡ |
| allQE(32)allMBR | | 0.464‡ | 0.171‡ | 91.6‡ | 82.7‡ | 80.8‡ | 70.8‡ | 86.7‡ | 83.7‡ | 82.3‡ | 75.4‡ | 86.1‡ | 58.1‡ | 56.0‡ | 29.9‡ | 54.0‡ |
| allQE(32)nolexMBR | | 0.421‡ | 0.146‡ | 92.6‡ | 83.8‡ | 81.1‡ | 70.9‡ | 86.8‡ | 83.8‡ | 82.8‡ | 75.4‡ | 85.6‡ | 56.9‡ | 54.8‡ | 28.4‡ | 56.2‡ |
| topQE(32)topMBR | | 0.323‡ | 0.110‡ | 94.7‡ | 85.7‡ | 81.3‡ | 70.8‡ | 86.5‡ | 82.8† | 81.5‡ | 74.2 | 84.0‡ | 53.0‡ | 50.8‡ | 23.5‡ | 62.9‡ |
| noncQE(32)noncMBR | | 0.446‡ | 0.121‡ | 90.4‡ | 81.2‡ | 79.4‡ | 69.8‡ | 86.4‡ | 83.5‡ | 82.1‡ | 75.1‡ | 85.8‡ | 57.3‡ | 55.2‡ | 29.2‡ | 54.8‡ |
| noncQE(32)noncnolexMBR | | 0.398‡ | 0.109‡ | 91.0‡ | 82.2‡ | 79.6‡ | 69.9‡ | 86.5‡ | 83.6‡ | 82.7‡ | 75.1‡ | 85.4‡ | 56.2‡ | 54.1‡ | 27.6‡ | 56.9‡ |
| mxQE(32)mxMBR | | 0.266‡ | 0.0795‡ | 92.1‡ | 81.4‡ | 78.9‡ | 68.9 | 86.1‡ | 82.2 | 80.8 | 73.3‡ | 83.4‡ | 51.5‡ | 49.3‡ | 22.2‡ | 65.7‡ |
| ckQE(32)xcMBR | | 0.445‡ | 0.157‡ | 95.8‡ | 81.4‡ | 81.9‡ | 69.6‡ | 86.3‡ | 82.1* | 80.6 | 73.5‡ | 83.6‡ | 51.3‡ | 49.1‡ | 22.2‡ | 63.6‡ |
| mxQE(32)xcMBR | | 0.412‡ | 0.104‡ | 95.9‡ | 80.8‡ | 79.5‡ | 68.5 | 85.9‡ | 81.7‡ | 80.5 | 73.0‡ | 83.2‡ | 50.0‡ | 47.9‡ | 21.1‡ | 65.0‡ |
| ckQE(32)mxMBR | | 0.282‡ | 0.0967‡ | 92.8‡ | 82.7‡ | 81.7‡ | 70.0‡ | 86.5‡ | 82.5 | 81.0† | 73.9† | 83.7‡ | 52.5‡ | 50.2‡ | 22.9‡ | 64.3‡ |

Table 15: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-hi (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | IndicCOMET | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 0.789 | 0.491 | 80.6 | 65.9 | 83.4 | 76.6 | 86.7 | 90.2 | 71.9 | 80.9 | 83.7 | 60.2 | 54.4 | 17.6 | 73.4 |
| MetricX | 0.215‡ | 0.100‡ | 87.5‡ | 67.9‡ | 85.5‡ | 77.6‡ | 86.8 | 90.2 | 72.1 | 82.1‡ | 80.8‡ | 53.8‡ | 47.3‡ | 11.0‡ | 87.1‡ |
| MetricX-QE | 0.424‡ | 0.0432‡ | 84.1‡ | 65.3 | 84.4‡ | 76.7 | 86.7 | 89.7‡ | 71.1‡ | 81.0 | 80.7‡ | 53.0‡ | 46.7‡ | 10.9‡ | 85.6‡ |
| XCOMET-XXL | 0.442‡ | 0.185‡ | 93.6‡ | 69.1‡ | 86.7‡ | 77.9‡ | 87.0† | 90.3 | 72.0 | 82.1‡ | 81.0‡ | 53.4‡ | 47.1‡ | 11.4‡ | 84.9‡ |
| XCOMET-XL | 0.471‡ | 0.225‡ | 87.4‡ | 76.7‡ | 86.1‡ | 79.1‡ | 87.3‡ | 90.6‡ | 72.6‡ | 82.5‡ | 81.7‡ | 55.3‡ | 49.0‡ | 12.5‡ | 82.4‡ |
| CometKiwi23-XXL | 0.492‡ | 0.215‡ | 88.4‡ | 69.4‡ | 90.1‡ | 79.4‡ | 87.5‡ | 90.6‡ | 72.3* | 82.4‡ | 81.7‡ | 55.0‡ | 48.6‡ | 12.0‡ | 82.0‡ |
| CometKiwi23-XL | 0.549‡ | 0.256‡ | 85.4‡ | 70.1‡ | 87.2‡ | 82.3‡ | 87.5‡ | 90.6‡ | 72.1 | 82.8‡ | 81.5‡ | 54.9‡ | 48.4‡ | 11.5‡ | 84.6‡ |
| CometKiwi22 | 0.578‡ | 0.263‡ | 83.9‡ | 68.0‡ | 86.2‡ | 78.9‡ | 88.5‡ | 90.8‡ | 72.4† | 81.6‡ | 81.9‡ | 55.1‡ | 48.7‡ | 12.0‡ | 80.8‡ |
| COMET22 | 0.509‡ | 0.274‡ | 84.8‡ | 69.3‡ | 86.0‡ | 78.9‡ | 87.6‡ | 92.1‡ | 73.7‡ | 82.7‡ | 82.9‡ | 58.2‡ | 52.0‡ | 14.5‡ | 78.8‡ |
| IndicCOMET | 0.628‡ | 0.356‡ | 82.3‡ | 66.4 | 83.9* | 77.1* | 87.0† | 90.7‡ | 76.9‡ | 82.1‡ | 82.0‡ | 55.7‡ | 49.3‡ | 12.5‡ | 81.5‡ |
| BLEURT | 0.636‡ | 0.347‡ | 82.2‡ | 65.0* | 84.6‡ | 77.9‡ | 86.7 | 90.0 | 71.4† | 87.0‡ | 80.7‡ | 53.6‡ | 46.9‡ | 10.3‡ | 89.7‡ |
| YiSi | 0.698* | 0.416* | 80.9 | 66.3 | 84.0† | 77.4‡ | 87.0‡ | 90.6‡ | 72.1 | 81.3† | 84.2‡ | 60.4 | 54.4 | 16.7† | 73.6 |
| chrF | 0.679† | 0.418* | 80.9 | 66.3 | 84.1‡ | 77.7‡ | 86.9* | 90.6‡ | 72.1 | 81.8‡ | 83.9 | 61.4‡ | 55.0‡ | 16.4‡ | 76.3* |
| chrF++ | 0.686* | 0.439 | 80.9 | 66.2 | 83.9* | 77.6‡ | 86.9‡ | 90.5‡ | 72.2† | 81.6‡ | 83.9* | 61.3‡ | 55.2‡ | 17.1* | 75.2 |
| sentBLEU | 0.772 | 0.490 | 80.4 | 66.0 | 83.2 | 76.7 | 86.8 | 90.3 | 71.9 | 80.8 | 83.8 | 59.9 | 54.2 | 17.9 | 71.5 |
| TER | 0.835 | 0.495 | 79.8* | 65.2* | 82.5‡ | 75.9‡ | 86.5 | 89.9* | 71.6‡ | 80.2‡ | 83.6 | 59.0‡ | 53.2‡ | 17.2 | 69.5‡ |
| | | | | | | | | | | | | | | | |
| rankAvg:all | 0.419‡ | 0.175‡ | 87.2‡ | 71.3‡ | 87.5‡ | 79.7‡ | 87.8‡ | 91.4‡ | 73.8‡ | 83.2‡ | 83.5* | 59.6† | 53.5‡ | 16.0‡ | 75.6 |
| rankAvg:qe | 0.388‡ | 0.105‡ | 87.8‡ | 71.3‡ | 88.8‡ | 81.0‡ | 88.1‡ | 91.1‡ | 72.8‡ | 82.8‡ | 82.2‡ | 56.1‡ | 49.7‡ | 12.5‡ | 81.3‡ |
| rankAvg:top | 0.321‡ | 0.0985‡ | 90.8‡ | 73.1‡ | 88.7‡ | 80.8‡ | 87.7‡ | 91.0‡ | 72.9‡ | 83.1‡ | 81.8‡ | 55.6‡ | 49.1‡ | 12.2‡ | 83.3‡ |
| rankAvg:topQe | 0.385‡ | 0.0913‡ | 88.4‡ | 70.9‡ | 89.1‡ | 81.3‡ | 87.7‡ | 90.8‡ | 72.7‡ | 82.8‡ | 81.8‡ | 55.2‡ | 48.8‡ | 12.0‡ | 83.3‡ |
| rankAvg:mxmxqe | 0.229‡ | 0.0563‡ | 87.4‡ | 67.9‡ | 85.8‡ | 77.8‡ | 86.9 | 90.2 | 72.1 | 82.2‡ | 80.8‡ | 53.9‡ | 47.4‡ | 11.0‡ | 86.3‡ |
| rankAvg:noLex | 0.368‡ | 0.136‡ | 88.8‡ | 72.5‡ | 88.1‡ | 80.2‡ | 87.9‡ | 91.5‡ | 74.1‡ | 83.7‡ | 83.0‡ | 58.0‡ | 51.7‡ | 13.9‡ | 79.3‡ |
| rankAvg:noNC | 0.441‡ | 0.181‡ | 85.2‡ | 69.5‡ | 86.1‡ | 78.7‡ | 87.5‡ | 91.4‡ | 74.0‡ | 83.2‡ | 83.7 | 60.1 | 54.0 | 16.6‡ | 74.0 |
| rankAvg:noNCnoLex | 0.366‡ | 0.136‡ | 86.2‡ | 70.1‡ | 86.5‡ | 79.0‡ | 87.5‡ | 91.5‡ | 74.5‡ | 83.9‡ | 83.2‡ | 58.5‡ | 52.2‡ | 14.5‡ | 78.4‡ |
| allQE(32)allMBR | 0.414‡ | 0.184‡ | 87.6‡ | 71.3‡ | 87.3‡ | 79.6‡ | 87.7‡ | 91.4‡ | 73.8‡ | 83.0‡ | 83.4‡ | 59.3‡ | 53.1‡ | 15.4‡ | 76.3* |
| allQE(32)nolexMBR | 0.370‡ | 0.151‡ | 89.2‡ | 72.5‡ | 87.8‡ | 79.9‡ | 87.8‡ | 91.5‡ | 74.2‡ | 83.7‡ | 82.9‡ | 57.9‡ | 51.7‡ | 14.2‡ | 79.3‡ |
| topQE(32)topMBR | 0.308‡ | 0.121‡ | 91.0‡ | 73.6‡ | 88.1‡ | 79.8‡ | 87.5‡ | 90.8‡ | 72.8‡ | 82.7‡ | 81.8‡ | 55.5‡ | 49.2‡ | 12.7‡ | 82.5‡ |
| noncQE(32)noncMBR | 0.399‡ | 0.138‡ | 85.6‡ | 69.1‡ | 86.3‡ | 78.4‡ | 87.5‡ | 91.3‡ | 73.6‡ | 82.8‡ | 83.4‡ | 59.1‡ | 53.0‡ | 15.9‡ | 75.8* |
| noncQE(32)noncnolexMBR | 0.343‡ | 0.117‡ | 86.8‡ | 69.9‡ | 86.7‡ | 78.9‡ | 87.5‡ | 91.4‡ | 74.3‡ | 83.7‡ | 82.9‡ | 57.8‡ | 51.5‡ | 14.2‡ | 78.9‡ |
| mxQE(32)mxMBR | 0.230‡ | 0.0857‡ | 87.5‡ | 68.2‡ | 85.7‡ | 77.8‡ | 86.9* | 90.3 | 72.2 | 82.3‡ | 81.0‡ | 54.2‡ | 47.7‡ | 11.2‡ | 85.9‡ |
| ckQE(32)xcMBR | 0.432‡ | 0.176‡ | 93.2‡ | 69.5‡ | 88.3‡ | 78.7‡ | 87.3‡ | 90.4* | 72.1 | 82.3‡ | 81.2‡ | 54.1‡ | 47.7‡ | 11.6‡ | 84.1‡ |
| mxQE(32)xcMBR | 0.387‡ | 0.115‡ | 93.2‡ | 68.9‡ | 86.9‡ | 77.9‡ | 87.1‡ | 90.3 | 72.0 | 82.3‡ | 81.0‡ | 53.4‡ | 47.1‡ | 11.4‡ | 84.4‡ |
| ckQE(32)mxMBR | 0.254‡ | 0.102‡ | 88.2‡ | 69.8‡ | 88.0‡ | 78.9‡ | 87.3‡ | 90.6‡ | 72.5‡ | 82.7‡ | 81.5‡ | 54.9‡ | 48.4‡ | 11.8‡ | 84.6‡ |

Table 16: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-ta (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | IndicCOMET | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 0.794 | 0.263 | 94.3 | 86.2 | 76.2 | 69.1 | 87.7 | 89.1 | 95.2 | 83.2 | 84.7 | 54.9 | 51.7 | 21.8 | 65.0 |
| MetricX | 0.196‡ | 0.0327‡ | 96.2‡ | 87.9‡ | 76.5 | 68.6* | 87.7 | 89.1 | 95.3 | 84.0‡ | 81.8‡ | 48.0‡ | 44.5‡ | 14.7‡ | 78.9‡ |
| MetricX-QE | 0.583‡ | 0.00737‡ | 94.5 | 84.1‡ | 73.5‡ | 67.0‡ | 87.2‡ | 88.0‡ | 93.6‡ | 81.8‡ | 81.3‡ | 45.6‡ | 42.4‡ | 14.0‡ | 78.3‡ |
| XCOMET-XXL | 0.547‡ | 0.107‡ | 98.5‡ | 87.4‡ | 77.6‡ | 68.3‡ | 87.6 | 88.7‡ | 94.3‡ | 82.8 | 81.5‡ | 46.5‡ | 43.1‡ | 14.1‡ | 77.9‡ |
| XCOMET-XL | 0.468‡ | 0.0935‡ | 96.5‡ | 92.9‡ | 78.3‡ | 70.5‡ | 88.2‡ | 89.5‡ | 95.2 | 84.4‡ | 82.7‡ | 49.3‡ | 45.9‡ | 16.0‡ | 74.7‡ |
| CometKiwi23-XXL | 0.609‡ | 0.133‡ | 96.2‡ | 87.9‡ | 86.2‡ | 71.7‡ | 88.2‡ | 89.1 | 94.8‡ | 83.3 | 82.6‡ | 49.2‡ | 45.8‡ | 15.7‡ | 74.9‡ |
| CometKiwi23-XL | 0.713† | 0.186‡ | 94.7* | 88.0‡ | 79.8‡ | 75.8‡ | 88.1‡ | 88.8† | 94.3‡ | 82.7* | 82.3‡ | 48.3‡ | 45.0‡ | 15.0‡ | 77.7‡ |
| CometKiwi22 | 0.631‡ | 0.132‡ | 95.2‡ | 87.9‡ | 79.2‡ | 70.8‡ | 89.3‡ | 89.6‡ | 95.4* | 83.6‡ | 82.6‡ | 48.5‡ | 45.2‡ | 15.4‡ | 74.5‡ |
| COMET22 | 0.578‡ | 0.134‡ | 95.6‡ | 88.7‡ | 78.5‡ | 70.5‡ | 88.5‡ | 91.0‡ | 96.5‡ | 84.2‡ | 83.8‡ | 51.9‡ | 48.6‡ | 17.9‡ | 70.3‡ |
| IndicCOMET | 0.745 | 0.233 | 94.4 | 85.4* | 75.4* | 68.5† | 87.8 | 89.3* | 99.4‡ | 83.2 | 82.8‡ | 49.5‡ | 46.1‡ | 16.1‡ | 73.4‡ |
| BLEURT | 0.608‡ | 0.168‡ | 94.9† | 87.3‡ | 77.0* | 69.2 | 87.9‡ | 89.1 | 95.1 | 88.0‡ | 82.3‡ | 48.1‡ | 44.6‡ | 14.4‡ | 77.7‡ |
| YiSi | 0.762 | 0.252 | 94.3 | 86.4 | 77.6‡ | 69.9‡ | 87.8 | 89.5‡ | 95.4* | 83.6† | 85.3‡ | 55.1 | 51.8 | 21.0† | 65.1 |
| chrF | 0.772 | 0.269 | 94.2 | 86.5 | 77.4‡ | 70.2‡ | 87.8 | 89.4‡ | 95.4† | 83.6† | 84.9‡ | 55.8‡ | 52.4‡ | 21.2* | 66.9‡ |
| chrF++ | 0.767 | 0.266 | 94.2 | 86.4 | 77.4‡ | 70.2‡ | 87.8 | 89.4‡ | 95.4* | 83.6† | 85.0† | 55.8‡ | 52.5‡ | 21.5 | 66.7‡ |
| sentBLEU | 0.813 | 0.290 | 94.0 | 85.7 | 75.8 | 68.7* | 87.5† | 89.1 | 95.0* | 82.9 | 84.7 | 54.4* | 51.4 | 21.8 | 64.1* |
| TER | 0.827 | 0.298 | 94.0 | 85.2‡ | 74.6‡ | 67.6‡ | 87.4‡ | 88.9† | 94.7‡ | 82.5‡ | 84.5‡ | 53.6‡ | 50.5‡ | 21.4 | 62.4‡ |
| | | | | | | | | | | | | | | | |
| rankAvg:all | 0.449‡ | 0.0762‡ | 96.7‡ | 90.0‡ | 81.3‡ | 72.0‡ | 88.6‡ | 90.3‡ | 96.7‡ | 84.9‡ | 84.6 | 54.2† | 50.8‡ | 20.1‡ | 67.1‡ |
| rankAvg:qe | 0.466‡ | 0.0402‡ | 96.4‡ | 89.7‡ | 83.7‡ | 73.8‡ | 88.9‡ | 89.8‡ | 95.6‡ | 83.9‡ | 83.0‡ | 49.9‡ | 46.5‡ | 16.3‡ | 74.6‡ |
| rankAvg:top | 0.346‡ | 0.0339‡ | 97.4‡ | 91.1‡ | 83.3‡ | 73.3‡ | 88.4‡ | 89.7‡ | 95.7‡ | 84.4‡ | 83.1‡ | 50.5‡ | 47.1‡ | 16.5‡ | 74.8‡ |
| rankAvg:topQe | 0.487‡ | 0.0313‡ | 96.4‡ | 89.5‡ | 84.3‡ | 74.3‡ | 88.4‡ | 89.4† | 95.1 | 83.5 | 82.8‡ | 49.7‡ | 46.3‡ | 16.1‡ | 75.3‡ |
| rankAvg:mxmxqe | 0.201‡ | 0.0128‡ | 96.1‡ | 87.9‡ | 76.5 | 68.5* | 87.7 | 89.0 | 95.3 | 84.0‡ | 81.8‡ | 47.8‡ | 44.3‡ | 14.5‡ | 79.0‡ |
| rankAvg:noLex | 0.397‡ | 0.0527‡ | 97.0‡ | 90.6‡ | 82.2‡ | 72.6‡ | 88.7‡ | 90.3‡ | 97.0‡ | 85.4‡ | 84.1‡ | 52.6‡ | 49.2‡ | 18.3‡ | 70.6‡ |
| rankAvg:noNC | 0.467‡ | 0.0805‡ | 95.7‡ | 88.8‡ | 78.9‡ | 70.2‡ | 88.2‡ | 90.2‡ | 96.7‡ | 84.9‡ | 84.6 | 54.3† | 51.0‡ | 20.6‡ | 66.2† |
| rankAvg:noNCnoLex | 0.381‡ | 0.0535‡ | 96.2‡ | 89.3‡ | 79.1‡ | 70.4‡ | 88.3‡ | 90.3‡ | 97.3‡ | 85.7‡ | 84.4‡ | 53.2‡ | 49.8‡ | 18.8‡ | 69.3‡ |
| allQE(32)allMBR | 0.450‡ | 0.0803‡ | 96.7‡ | 90.1‡ | 81.4‡ | 71.8‡ | 88.6‡ | 90.3‡ | 96.6‡ | 84.9‡ | 84.5‡ | 54.0‡ | 50.7‡ | 20.0‡ | 67.6‡ |
| allQE(32)nolexMBR | 0.384‡ | 0.0615‡ | 97.2‡ | 90.9‡ | 81.6‡ | 71.8‡ | 88.6‡ | 90.3‡ | 97.0‡ | 85.5‡ | 84.0‡ | 52.7‡ | 49.3‡ | 18.6‡ | 70.5‡ |
| topQE(32)topMBR | 0.312‡ | 0.0438‡ | 97.7‡ | 91.5‡ | 81.8‡ | 71.8‡ | 88.4‡ | 89.7‡ | 95.7‡ | 84.6‡ | 82.8‡ | 50.0‡ | 46.5‡ | 16.3‡ | 75.0‡ |
| noncQE(32)noncMBR | 0.417‡ | 0.0465‡ | 95.9‡ | 88.6‡ | 78.7‡ | 69.9‡ | 88.2‡ | 90.0‡ | 96.3‡ | 84.7‡ | 84.2‡ | 53.2‡ | 49.9‡ | 19.5‡ | 68.2‡ |
| noncQE(32)noncnolexMBR | 0.361‡ | 0.0395‡ | 96.2‡ | 89.0‡ | 79.1‡ | 70.2‡ | 88.3‡ | 90.1‡ | 96.8‡ | 85.4‡ | 83.7‡ | 51.6‡ | 48.2‡ | 17.7‡ | 71.9‡ |
| mxQE(32)mxMBR | 0.217‡ | 0.0204‡ | 96.3‡ | 88.1‡ | 76.7 | 68.6* | 87.7 | 89.1 | 95.1 | 83.9‡ | 82.1‡ | 48.2‡ | 44.8‡ | 15.1‡ | 77.4‡ |
| ckQE(32)xcMBR | 0.550‡ | 0.0946‡ | 98.2‡ | 88.1‡ | 82.1‡ | 70.0‡ | 88.0‡ | 89.1 | 94.7‡ | 83.4 | 82.3‡ | 48.2‡ | 44.8‡ | 15.0‡ | 76.0‡ |
| mxQE(32)xcMBR | 0.455‡ | 0.0397‡ | 98.3‡ | 87.8‡ | 77.7‡ | 68.3‡ | 87.7 | 88.8* | 94.5‡ | 83.2 | 81.7‡ | 46.7‡ | 43.3‡ | 14.0‡ | 77.5‡ |
| ckQE(32)mxMBR | 0.244‡ | 0.0348‡ | 96.7‡ | 89.0‡ | 82.1‡ | 70.8‡ | 88.2‡ | 89.5‡ | 95.5* | 84.4‡ | 82.5‡ | 49.5‡ | 46.0‡ | 15.7‡ | 76.6‡ |

Table 17: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-gu (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

## G.8 Results for English-Malayalam (en-ml) on FLORES200 test dataset

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | IndicCOMET | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 0.849 | 0.393 | 91.7 | 83.6 | 82.6 | 77.1 | 86.9 | 88.5 | 94.3 | 80.9 | 82.4 | 56.5 | 50.9 | 15.9 | 74.9 |
| MetricX | **0.245‡** | **0.0652‡** | 95.8‡ | **86.3‡** | **84.6‡** | 77.8‡ | 86.9 | **89.0†** | 94.5 | **81.8‡** | 79.4‡ | 50.4‡ | 44.3‡ | 10.3‡ | 87.7‡ |
| MetricX-QE | **0.572‡** | **0.0197‡** | 93.3‡ | 82.5† | 82.5 | 76.4‡ | 86.4‡ | 87.7‡ | 92.8‡ | 79.6‡ | 78.8‡ | 48.2‡ | 42.2‡ | 9.63‡ | 89.1‡ |
| XCOMET-XXL | **0.526‡** | **0.150‡** | **98.5‡** | 86.3‡ | 85.1‡ | 77.8‡ | 87.1* | 88.9 | 93.7† | 81.2 | 79.7‡ | 50.6‡ | 44.5‡ | 10.1‡ | 86.9‡ |
| XCOMET-XL | **0.464‡** | **0.129‡** | 96.0‡ | **91.9‡** | 85.8‡ | 79.6‡ | 87.7‡ | 89.8‡ | 95.0‡ | 83.0‡ | 80.7‡ | 52.6‡ | 46.7‡ | 12.2‡ | 82.4‡ |
| CometKiwi23-XXL | **0.553‡** | **0.134‡** | 96.5‡ | 87.4‡ | **89.3‡** | 79.5‡ | 87.9‡ | 89.5‡ | 94.6 | 82.0‡ | 80.6‡ | 51.8‡ | 45.8‡ | 11.1‡ | 82.9‡ |
| CometKiwi23-XL | **0.577‡** | **0.165‡** | 95.0‡ | 88.9‡ | 86.4‡ | **81.6‡** | 87.9‡ | 89.5‡ | 94.5 | 82.1‡ | 80.5‡ | 52.1‡ | 46.1‡ | 11.4‡ | 84.2‡ |
| CometKiwi22 | **0.623‡** | **0.208‡** | 94.0‡ | 86.7‡ | 85.3‡ | 78.9‡ | **88.9‡** | 89.6‡ | 94.7† | 81.7‡ | 80.7‡ | 52.0‡ | 46.1‡ | 11.6‡ | 82.5‡ |
| COMET22 | **0.519‡** | **0.181‡** | 94.7‡ | 87.4‡ | 85.5‡ | 78.9‡ | 87.9‡ | **91.2‡** | 96.2‡ | 83.0‡ | 81.7‡ | 55.1‡ | 49.0‡ | 13.0‡ | 79.4‡ |
| IndicCOMET | **0.738†** | **0.295†** | 92.3 | 83.7 | 82.9 | 77.1 | 87.2‡ | 89.5‡ | **99.1‡** | 81.6‡ | 80.5‡ | 52.0‡ | 45.9‡ | 11.1‡ | 83.3‡ |
| BLEURT | **0.621‡** | **0.237‡** | 93.3‡ | 84.8‡ | 84.0‡ | 77.8‡ | 87.4‡ | 89.2‡ | 94.3 | **87.0‡** | 79.9‡ | 50.9‡ | 44.7‡ | 9.60‡ | 89.5‡ |
| YiSi | **0.754*** | **0.332*** | 92.1 | 84.1 | 83.5‡ | 77.7‡ | 87.3‡ | 89.5‡ | 95.1‡ | 81.9‡ | **83.3‡** | 57.9‡ | 52.0‡ | 15.4 | 74.3 |
| chrF | 0.790 | 0.348 | 91.2 | 83.4 | 83.3‡ | 77.5† | 87.2‡ | 89.5‡ | 95.1‡ | 81.9‡ | 83.0‡ | **59.1‡** | 52.9‡ | **14.9‡** | 76.9‡ |
| chrF++ | 0.789 | 0.384 | 91.3 | 83.6 | 83.1* | 77.5† | 87.2‡ | 89.4‡ | 94.8‡ | 81.9‡ | 82.9‡ | 58.7‡ | **52.8‡** | **15.2*** | 76.5† |
| sentBLEU | 0.877 | 0.415 | 91.1* | 83.0 | 82.1‡ | 76.7* | 86.8 | 89.0† | 94.4 | 81.1 | 82.7* | 56.8 | 51.5 | **16.9‡** | 72.4‡ |
| TER | 0.914 | 0.454 | 91.1* | 82.7* | 80.9‡ | 75.9‡ | 86.4‡ | 88.5 | 94.1 | 80.5 | 82.4 | 55.5† | 50.2* | 15.9 | **70.0‡** |
| | | | | | | | | | | | | | | | |
| rankAvg:all | **0.424‡** | **0.114‡** | 96.1‡ | 88.9‡ | 86.7‡ | 79.7‡ | 88.2‡ | 90.5‡ | 96.5‡ | 83.6‡ | 82.6 | 57.1 | 51.2 | **15.1*** | 75.2 |
| rankAvg:qe | **0.435‡** | **0.0638‡** | 96.1‡ | 89.1‡ | 88.1‡ | 80.6‡ | 88.5‡ | 90.1‡ | 95.5‡ | 83.0‡ | 81.1‡ | 53.3‡ | 47.2‡ | 12.0‡ | 82.1‡ |
| rankAvg:top | **0.345‡** | **0.0578‡** | 97.6‡ | 90.1‡ | 87.9‡ | 80.4‡ | 88.0‡ | 90.1‡ | 95.4‡ | 83.3‡ | 81.0‡ | 53.3‡ | 47.3‡ | 12.5‡ | 82.1‡ |
| rankAvg:topQe | **0.422‡** | **0.0521‡** | 96.2‡ | 89.0‡ | 88.4‡ | 80.9‡ | 88.0‡ | 89.9‡ | 95.1‡ | 82.7‡ | 80.9‡ | 52.8‡ | 46.7‡ | 11.6‡ | 82.9‡ |
| rankAvg:mxmxqe | **0.264‡** | **0.0280‡** | 95.9‡ | 86.2‡ | 84.8‡ | 77.9‡ | 86.9 | **89.0*** | 94.5 | 81.6† | 79.4‡ | 50.3‡ | 44.2‡ | 10.3‡ | 87.6‡ |
| rankAvg:noLex | **0.375‡** | **0.0845‡** | 96.8‡ | 89.8‡ | 87.3‡ | 80.1‡ | 88.3‡ | 90.6‡ | 96.9‡ | 84.2‡ | 82.2 | 55.8 | **49.7‡** | 13.6‡ | 79.2‡ |
| rankAvg:noNC | **0.464‡** | **0.127‡** | 94.7‡ | 87.4‡ | 85.5‡ | 78.8‡ | 87.9‡ | 90.5‡ | 96.5‡ | 83.6‡ | 82.7* | 57.4† | 51.5 | 15.3 | 74.4 |
| rankAvg:noNCnoLex | **0.385‡** | **0.0801‡** | 95.6‡ | 88.2‡ | 86.1‡ | 79.2‡ | 88.0‡ | 90.6‡ | 97.1‡ | 84.5‡ | 82.2 | 55.9 | **49.7‡** | 13.5‡ | 79.1‡ |
| allQE(32)allMBR | **0.413‡** | **0.112‡** | 96.2‡ | 89.2‡ | 86.6‡ | 79.6‡ | 88.1‡ | 90.5‡ | 96.4‡ | 83.7‡ | 82.5 | 56.8 | 50.9 | **14.8‡** | **76.1*** |
| allQE(32)nolexMBR | **0.367‡** | **0.0948‡** | 96.8‡ | 90.0‡ | 86.8‡ | 79.8‡ | 88.1‡ | 90.6‡ | 96.9‡ | 84.3‡ | **82.1‡** | 55.6* | 49.5‡ | 13.7‡ | 79.1‡ |
| topQE(32)topMBR | **0.325‡** | **0.0696‡** | 97.7‡ | 90.5‡ | 86.9‡ | 79.7‡ | 87.8‡ | 89.9‡ | 95.3‡ | 83.1‡ | 80.8‡ | 53.0‡ | 46.9‡ | 12.2‡ | 82.7‡ |
| noncQE(32)noncMBR | **0.427‡** | **0.0816‡** | 95.0‡ | 87.5‡ | 85.5‡ | 78.7‡ | 87.8‡ | 90.1‡ | 96.2‡ | 83.0‡ | 82.2 | 56.2 | 50.3 | 14.5‡ | 76.8† |
| noncQE(32)noncnolexMBR | **0.376‡** | **0.0669‡** | 95.7‡ | 88.0‡ | 85.8‡ | 79.0‡ | 87.9‡ | 90.4‡ | 96.7‡ | 84.0‡ | 81.8‡ | 54.8‡ | 48.8‡ | 13.1‡ | 80.1‡ |
| mxQE(32)mxMBR | **0.259‡** | **0.0452‡** | 95.9‡ | 86.6‡ | 84.9‡ | 78.0‡ | 87.0 | **89.0*** | 94.5 | 81.8‡ | 79.7‡ | 50.8‡ | 44.7‡ | 10.5‡ | 87.1‡ |
| ckQE(32)xcMBR | **0.487‡** | **0.131‡** | **98.4‡** | 87.2‡ | 86.9‡ | 78.5‡ | 87.5‡ | 89.4‡ | 94.5 | 81.9‡ | 80.3‡ | 51.6‡ | 45.6‡ | 11.2‡ | 83.7‡ |
| mxQE(32)xcMBR | **0.453‡** | **0.0756‡** | **98.4‡** | 87.0‡ | 85.3‡ | 77.9‡ | 87.2* | 89.1† | 94.1 | 81.3 | 79.9‡ | 50.6‡ | 44.6‡ | 10.7‡ | 85.8‡ |
| ckQE(32)mxMBR | **0.266‡** | **0.0570‡** | 96.4‡ | 87.6‡ | 87.0‡ | 78.9‡ | 87.5‡ | 89.5‡ | 95.1‡ | 82.5‡ | 80.1‡ | 51.6‡ | 45.4‡ | 11.0‡ | 85.1‡ |

Table 18: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-ml (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

## G.9 Results for English-Vietnamese (en-vi) on FLORES200 test dataset

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.16 | 0.555 | 93.7 | 92.7 | 91.6 | 79.6 | 85.9 | 90.4 | 76.8 | 89.2 | 62.6 | 62.6 | 42.7 | 41.3 |
| MetricX | 0.486‡ | 0.211‡ | 96.6‡ | 94.7‡ | 93.8‡ | 81.0‡ | 86.4‡ | 90.2† | 75.9‡ | 86.8‡ | 56.4‡ | 56.2‡ | 32.7‡ | 53.4‡ |
| MetricX-QE | 0.673‡ | 0.117‡ | 95.7‡ | 93.8‡ | 93.3‡ | 80.5‡ | 86.2‡ | 89.8‡ | 74.7‡ | 86.4‡ | 54.5‡ | 54.3‡ | 31.6‡ | 53.5‡ |
| XCOMET-XXL | 0.755‡ | 0.298‡ | 98.4‡ | 94.8‡ | 94.6‡ | 81.1‡ | 86.4‡ | 90.1† | 75.4‡ | 86.9‡ | 56.2‡ | 56.0‡ | 33.1‡ | 52.5‡ |
| XCOMET-XL | 0.725‡ | 0.303‡ | 96.9‡ | 96.6‡ | 94.1‡ | 82.1‡ | 86.6‡ | 90.5 | 76.6 | 87.6‡ | 57.9‡ | 57.8‡ | 35.3‡ | 49.4‡ |
| CometKiwi23-XXL | 0.809‡ | 0.310‡ | 97.0‡ | 94.7‡ | 96.0‡ | 81.7‡ | 86.7‡ | 90.4 | 76.0‡ | 87.5‡ | 57.3‡ | 57.1‡ | 34.7‡ | 50.2‡ |
| CometKiwi23-XL | 0.832‡ | 0.328‡ | 96.1‡ | 95.3‡ | 94.4‡ | 83.6‡ | 86.7‡ | 90.5 | 76.1‡ | 87.4‡ | 57.3‡ | 57.2‡ | 34.5‡ | 50.5‡ |
| CometKiwi22 | 0.912‡ | 0.366‡ | 95.4‡ | 94.3‡ | 93.8‡ | 81.4‡ | 87.6‡ | 90.4 | 75.9‡ | 87.3‡ | 57.1‡ | 56.9‡ | 34.2‡ | 50.7‡ |
| COMET22 | 0.907‡ | 0.410‡ | 95.4‡ | 94.3‡ | 93.3‡ | 81.0‡ | 86.6‡ | 91.5‡ | 77.2‡ | 88.6‡ | 60.9‡ | 60.9‡ | 39.2‡ | 44.5‡ |
| BLEURT | 0.923‡ | 0.420‡ | 95.0‡ | 94.2‡ | 93.1‡ | 80.7‡ | 86.5‡ | 90.7‡ | 78.9‡ | 88.6‡ | 60.4‡ | 60.4‡ | 39.2‡ | 44.8‡ |
| YiSi | 1.08* | 0.519† | 94.0 | 93.1* | 92.1‡ | 80.0‡ | 86.1‡ | 90.5† | 76.8 | 89.3* | 62.5 | 62.5 | 42.1* | 42.1† |
| chrF | 1.09 | 0.531 | 94.1* | 92.9 | 92.1‡ | 80.0‡ | 86.0† | 90.5 | 76.7 | 89.1 | 63.1† | 63.0* | 41.4‡ | 43.4‡ |
| chrF++ | 1.09* | 0.527* | 94.1† | 92.9 | 92.1‡ | 79.9‡ | 86.0† | 90.5 | 76.7 | 89.1 | 63.1† | 63.1‡ | 41.5‡ | 43.4‡ |
| sentBLEU | 1.11 | 0.546 | 93.9 | 92.9 | 91.7 | 79.8 | 86.0 | 90.4 | 76.8 | 89.2 | 62.6 | 62.6 | 42.5 | 41.3 |
| TER | 1.21 | 0.592* | 93.6 | 92.5 | 91.3* | 79.4* | 85.8* | 90.4 | 76.5 | 89.2 | 62.1* | 62.2* | 42.5 | 39.7‡ |
| | | | | | | | | | | | | | | |
| rankAvg:all | 0.759‡ | 0.302‡ | 96.4‡ | 95.3‡ | 94.4‡ | 81.7‡ | 86.8‡ | 91.1‡ | 77.8‡ | 89.0* | 62.3 | 62.2 | 41.2‡ | 43.1‡ |
| rankAvg:qe | 0.688‡ | 0.198‡ | 97.0‡ | 95.4‡ | 95.4‡ | 82.7‡ | 87.3‡ | 90.6† | 76.5 | 87.6‡ | 57.8‡ | 57.7‡ | 35.2‡ | 49.7‡ |
| rankAvg:top | 0.608‡ | 0.203‡ | 97.8‡ | 95.9‡ | 95.3‡ | 82.5‡ | 86.8‡ | 90.6† | 76.7 | 87.5‡ | 58.1‡ | 57.9‡ | 35.3‡ | 50.4‡ |
| rankAvg:topQe | 0.670‡ | 0.185‡ | 97.2‡ | 95.4‡ | 95.5‡ | 82.8‡ | 86.9‡ | 90.5 | 76.3* | 87.4‡ | 57.6‡ | 57.4‡ | 34.8‡ | 50.7‡ |
| rankAvg:mxmxqe | 0.518‡ | 0.144‡ | 96.7‡ | 94.7‡ | 94.0‡ | 81.1‡ | 86.5‡ | 90.1† | 75.6‡ | 86.7‡ | 56.0‡ | 55.8‡ | 32.3‡ | 53.8‡ |
| rankAvg:noLex | 0.680‡ | 0.254‡ | 97.2‡ | 95.8‡ | 94.9‡ | 82.2‡ | 87.0‡ | 91.1‡ | 77.9‡ | 88.6‡ | 60.9‡ | 60.8‡ | 39.1‡ | 45.4‡ |
| rankAvg:noNC | 0.795‡ | 0.324‡ | 95.7‡ | 94.5‡ | 93.4‡ | 81.0‡ | 86.5‡ | 91.0‡ | 77.9‡ | 89.2 | 62.5 | 62.5 | 41.6‡ | 42.3‡ |
| rankAvg:noNCnoLex | 0.715‡ | 0.264‡ | 96.1‡ | 94.9‡ | 93.8‡ | 81.3‡ | 86.7‡ | 91.2‡ | 78.1‡ | 89.0† | 61.8‡ | 61.8‡ | 40.5‡ | 43.7‡ |
| allQE(32)allMBR | 0.735‡ | 0.287‡ | 96.7‡ | 95.4‡ | 94.6‡ | 81.8‡ | 86.9‡ | 91.0‡ | 77.8‡ | 88.8‡ | 61.5‡ | 61.5‡ | 40.4‡ | 44.1‡ |
| allQE(32)nolexMBR | 0.681‡ | 0.261‡ | 97.2‡ | 95.8‡ | 94.8‡ | 82.0‡ | 86.9‡ | 91.1‡ | 77.9‡ | 88.6‡ | 60.7‡ | 60.7‡ | 39.1‡ | 45.4‡ |
| topQE(32)topMBR | 0.589‡ | 0.220‡ | 97.8‡ | 96.0‡ | 94.9‡ | 82.2‡ | 86.8‡ | 90.5 | 76.5 | 87.4‡ | 57.8‡ | 57.7‡ | 35.1‡ | 50.8‡ |
| noncQE(32)noncMBR | 0.717‡ | 0.242‡ | 96.2‡ | 94.8‡ | 93.8‡ | 81.2‡ | 86.6‡ | 91.0‡ | 77.7‡ | 88.8‡ | 61.3‡ | 61.2‡ | 40.1‡ | 44.1‡ |
| noncQE(32)noncnolexMBR | 0.650‡ | 0.226‡ | 96.3‡ | 95.0‡ | 94.0‡ | 81.4‡ | 86.7‡ | 91.1‡ | 77.8‡ | 88.6‡ | 60.9‡ | 60.8‡ | 39.5‡ | 44.8‡ |
| mxQE(32)mxMBR | 0.520‡ | 0.177‡ | 96.5‡ | 94.6‡ | 93.9‡ | 81.1‡ | 86.5‡ | 90.2† | 75.7‡ | 86.9‡ | 56.4‡ | 56.2‡ | 33.1‡ | 53.2‡ |
| ckQE(32)xcMBR | 0.739‡ | 0.280‡ | 98.3‡ | 94.9‡ | 95.2‡ | 81.6‡ | 86.5‡ | 90.3 | 75.7‡ | 87.1‡ | 56.7‡ | 56.6‡ | 34.0‡ | 51.4‡ |
| mxQE(32)xcMBR | 0.684‡ | 0.213‡ | 98.1‡ | 94.8‡ | 94.7‡ | 81.3‡ | 86.4‡ | 90.1† | 75.6‡ | 87.0‡ | 56.1‡ | 55.9‡ | 33.1‡ | 52.1‡ |
| ckQE(32)mxMBR | 0.529‡ | 0.213‡ | 97.1‡ | 95.0‡ | 94.9‡ | 81.6‡ | 86.6‡ | 90.4 | 76.3† | 87.2‡ | 57.3‡ | 57.2‡ | 34.1‡ | 51.9‡ |

Table 19: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1$^{st}$ block), and ensembles (2$^{nd}$ block), on en-vi (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for $p<0.05$, † for $p<0.01$, ‡ for $p<0.001$. The green diagonal in the 1$^{st}$ block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 0.589 | 0.338 | 96.2 | 94.1 | 93.7 | 82.7 | 87.9 | 90.9 | 90.5 | 84.6 | 60.6 | 57.4 | 26.7 | 58.7 |
| MetricX | 0.117‡ | 0.0585‡ | 97.1‡ | 95.8‡ | 95.4‡ | 83.7‡ | 88.0 | 90.6† | 92.5‡ | 81.6‡ | 54.1‡ | 50.3‡ | 18.5‡ | 71.5‡ |
| MetricX-QE | 0.344‡ | 0.0166‡ | 96.0 | 94.2 | 93.9 | 82.6 | 87.6‡ | 89.9‡ | 90.6 | 81.1‡ | 52.5‡ | 48.6‡ | 17.2‡ | 71.9‡ |
| XCOMET-XXL | 0.367‡ | 0.162‡ | 99.2‡ | 95.6‡ | 96.0‡ | 83.5‡ | 87.9 | 90.4‡ | 91.4‡ | 81.8‡ | 54.1‡ | 50.4‡ | 18.7‡ | 69.8‡ |
| XCOMET-XL | 0.336‡ | 0.140‡ | 97.9‡ | 98.0‡ | 96.0‡ | 85.6‡ | 88.5‡ | 91.2‡ | 92.9‡ | 82.9‡ | 56.3‡ | 52.8‡ | 21.2‡ | 67.0‡ |
| CometKiwi23-XXL | 0.369‡ | 0.149‡ | 98.1‡ | 96.1‡ | 97.7‡ | 85.2‡ | 88.5‡ | 91.1* | 92.1‡ | 83.0‡ | 56.3‡ | 52.8‡ | 21.3‡ | 66.6‡ |
| CometKiwi23-XL | 0.375‡ | 0.155‡ | 97.5‡ | 96.9‡ | 96.3‡ | 87.5‡ | 88.6‡ | 91.1 | 92.2‡ | 82.8‡ | 56.1‡ | 52.5‡ | 20.8‡ | 66.6‡ |
| CometKiwi22 | 0.449‡ | 0.179‡ | 96.7* | 95.7‡ | 95.5‡ | 84.5‡ | 89.5‡ | 91.3‡ | 91.8‡ | 83.1‡ | 56.6‡ | 53.0‡ | 20.9‡ | 66.1‡ |
| COMET22 | 0.414‡ | 0.212‡ | 97.1‡ | 95.7‡ | 95.3‡ | 84.4‡ | 88.7‡ | 92.4‡ | 92.2‡ | 84.4* | 59.6‡ | 56.2‡ | 24.6‡ | 61.7‡ |
| BLEURT | 0.402‡ | 0.205‡ | 96.5 | 95.5‡ | 95.0‡ | 83.7‡ | 88.1† | 90.5‡ | 96.4‡ | 81.3‡ | 53.6‡ | 49.8‡ | 17.4‡ | 75.1‡ |
| YiSi | 0.561 | 0.330 | 96.0 | 94.4 | 94.0 | 83.1† | 88.1‡ | 91.1‡ | 90.9‡ | 85.3‡ | 61.2† | 58.0‡ | 27.4* | 57.8* |
| chrF | 0.506‡ | 0.317 | 96.0 | 94.4 | 94.1* | 83.2* | 88.0* | 91.2‡ | 91.0‡ | 85.1‡ | 62.4‡ | 59.0‡ | 27.3 | 59.0 |
| chrF++ | 0.516‡ | 0.317 | 96.0 | 94.4 | 94.2* | 83.1* | 88.1† | 91.2‡ | 91.1‡ | 85.1‡ | 62.1‡ | 58.8‡ | 27.6† | 59.1 |
| sentBLEU | 0.589 | 0.338 | 95.8* | 94.1 | 93.6 | 82.7 | 87.8 | 90.8 | 90.4 | 84.8* | 60.6 | 57.6 | 28.0‡ | 57.2‡ |
| TER | 0.621 | 0.351 | 95.7† | 93.8 | 93.0‡ | 82.2‡ | 87.6‡ | 90.6‡ | 89.7‡ | 84.6 | 59.7‡ | 56.7* | 27.4 | 55.4‡ |
| | | | | | | | | | | | | | | |
| rankAvg:all | 0.282‡ | 0.116‡ | 98.1‡ | 96.8‡ | 96.4‡ | 85.5‡ | 88.8‡ | 91.9‡ | 93.0‡ | 84.9† | 61.2* | 57.8 | 26.8 | 58.9 |
| rankAvg:qe | 0.274‡ | 0.0509‡ | 97.9‡ | 97.0‡ | 97.1‡ | 86.6‡ | 89.2‡ | 91.5‡ | 92.8‡ | 83.3‡ | 57.3‡ | 53.8‡ | 22.2‡ | 65.2‡ |
| rankAvg:top | 0.196‡ | 0.0476‡ | 98.7‡ | 97.4‡ | 97.1‡ | 86.4‡ | 88.7‡ | 91.4‡ | 93.0‡ | 83.1‡ | 57.1‡ | 53.5‡ | 21.6‡ | 66.4‡ |
| rankAvg:topQe | 0.280‡ | 0.0397‡ | 98.1‡ | 97.0‡ | 97.3‡ | 86.8‡ | 88.7‡ | 91.3‡ | 92.6‡ | 83.0‡ | 56.6‡ | 53.1‡ | 21.4‡ | 65.9‡ |
| rankAvg:mxmxqe | 0.127‡ | 0.0239‡ | 97.3‡ | 95.9‡ | 95.5‡ | 83.9‡ | 88.0 | 90.6† | 92.5‡ | 81.6‡ | 54.1‡ | 50.3‡ | 18.5‡ | 71.2‡ |
| rankAvg:noLex | 0.231‡ | 0.0761‡ | 98.5‡ | 97.2‡ | 96.8‡ | 86.0‡ | 89.0‡ | 92.0‡ | 93.7‡ | 84.5 | 60.0* | 56.5† | 25.0‡ | 61.9‡ |
| rankAvg:noNC | 0.296‡ | 0.117‡ | 97.3‡ | 95.9‡ | 95.6‡ | 84.5‡ | 88.5‡ | 91.8‡ | 92.8‡ | 85.1‡ | 61.5‡ | 58.3‡ | 27.5* | 57.8* |
| rankAvg:noNCnoLex | 0.226‡ | 0.0690‡ | 97.5‡ | 96.3‡ | 95.9‡ | 84.8‡ | 88.6‡ | 92.0‡ | 93.7‡ | 84.6 | 60.1* | 56.7* | 25.4‡ | 61.2‡ |
| allQE(32)allMBR | 0.288‡ | 0.119‡ | 98.1‡ | 96.8‡ | 96.4‡ | 85.6‡ | 88.8‡ | 91.9‡ | 93.0‡ | 84.8 | 60.8 | 57.5 | 26.7 | 59.2 |
| allQE(32)nolexMBR | 0.228‡ | 0.0929‡ | 98.5‡ | 97.2‡ | 96.7‡ | 85.8‡ | 88.9‡ | 91.9‡ | 93.7‡ | 84.3* | 59.8‡ | 56.4‡ | 24.6‡ | 61.6‡ |
| topQE(32)topMBR | 0.182‡ | 0.0686‡ | 98.7‡ | 97.4‡ | 96.7‡ | 85.7‡ | 88.6‡ | 91.3‡ | 92.9‡ | 82.9‡ | 56.5‡ | 52.9‡ | 21.0‡ | 66.9‡ |
| noncQE(32)noncMBR | 0.265‡ | 0.0709‡ | 97.4‡ | 96.1‡ | 95.6‡ | 84.5‡ | 88.4‡ | 91.6‡ | 92.7‡ | 84.4 | 59.8† | 56.4‡ | 25.3‡ | 60.3† |
| noncQE(32)noncnolexMBR | 0.206‡ | 0.0591‡ | 97.5‡ | 96.4‡ | 95.8‡ | 84.8‡ | 88.5‡ | 91.8‡ | 93.6‡ | 84.0‡ | 58.7‡ | 55.2‡ | 23.7‡ | 63.2‡ |
| mxQE(32)mxMBR | 0.126‡ | 0.0373‡ | 97.2‡ | 95.8‡ | 95.4‡ | 83.9‡ | 88.0 | 90.6† | 92.4‡ | 81.6‡ | 53.9‡ | 50.1‡ | 18.3‡ | 71.2‡ |
| ckQE(32)xcMBR | 0.346‡ | 0.139‡ | 99.1‡ | 95.9‡ | 96.9‡ | 84.5‡ | 88.3‡ | 90.8 | 92.0‡ | 82.4‡ | 55.2‡ | 51.5‡ | 19.9‡ | 68.1‡ |
| mxQE(32)xcMBR | 0.298‡ | 0.0600‡ | 98.9‡ | 95.8‡ | 96.0‡ | 83.8‡ | 88.0 | 90.5‡ | 91.4‡ | 81.9‡ | 54.1‡ | 50.4‡ | 18.7‡ | 69.2‡ |
| ckQE(32)mxMBR | 0.132‡ | 0.0554‡ | 98.0‡ | 96.3‡ | 96.8‡ | 84.8‡ | 88.3‡ | 91.0 | 92.8‡ | 82.3‡ | 55.5‡ | 51.8‡ | 19.8‡ | 68.8‡ |

Table 20: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-hu (FLORES200 test dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

## G.11 Results for English-German (en-de) on WMT2023 dataset

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.24 | 1.42 | 90.1 | 87.2 | 79.6 | 70.7 | 81.3 | 85.6 | 73.5 | 87.9 | 70.1 | 68.2 | 45.4 | 42.1 |
| MetricX | 0.571‡ | 0.794‡ | 92.0‡ | 87.8* | 79.7 | 70.4 | 80.1‡ | 83.9‡ | 73.2 | 82.2‡ | 58.9‡ | 55.9‡ | 27.4‡ | 63.4‡ |
| MetricX-QE | 0.630‡ | 0.494‡ | 91.8‡ | 87.5 | 79.9 | 70.3 | 80.3‡ | 83.7‡ | 72.8* | 82.2‡ | 58.2‡ | 55.1‡ | 26.3‡ | 64.8‡ |
| XCOMET-XXL | 0.915‡ | 1.03‡ | 94.5‡ | 88.3‡ | 81.6‡ | 71.6† | 80.6‡ | 84.0‡ | 72.9 | 83.1‡ | 59.7‡ | 57.0‡ | 29.8‡ | 60.3‡ |
| XCOMET-XL | 0.907‡ | 1.04‡ | 92.1‡ | 90.8‡ | 81.1‡ | 72.5‡ | 81.2 | 84.5‡ | 73.3 | 83.8‡ | 61.1‡ | 58.4‡ | 31.3‡ | 59.0‡ |
| CometKiwi23-XXL | 1.06† | 1.09‡ | 92.0‡ | 88.4‡ | 85.5‡ | 72.6‡ | 81.5* | 85.1‡ | 72.9 | 84.9‡ | 63.4‡ | 60.9‡ | 34.0‡ | 54.7‡ |
| CometKiwi23-XL | 1.05‡ | 1.15‡ | 91.5‡ | 89.1‡ | 82.6‡ | 75.7‡ | 81.9‡ | 85.1† | 73.5 | 84.9‡ | 63.8‡ | 61.3‡ | 34.5‡ | 55.3‡ |
| CometKiwi22 | 1.11‡ | 1.20‡ | 91.0‡ | 88.0‡ | 81.5‡ | 72.3‡ | 83.4‡ | 85.5 | 73.3 | 85.5‡ | 64.7‡ | 62.2‡ | 35.8‡ | 52.7‡ |
| COMET22 | 1.01‡ | 1.23† | 91.4‡ | 88.1‡ | 80.6‡ | 71.3† | 81.6* | 87.0‡ | 74.7‡ | 86.4‡ | 67.3‡ | 65.0‡ | 39.7‡ | 47.3‡ |
| BLEURT | 0.874‡ | 0.999‡ | 91.5‡ | 88.2‡ | 80.8‡ | 71.4† | 81.3 | 85.4 | 77.4‡ | 84.8‡ | 63.9‡ | 61.3‡ | 34.0‡ | 54.6‡ |
| YiSi | 1.27 | 1.43 | 90.1 | 87.2 | 79.2 | 70.5 | 81.1 | 85.7 | 73.8 | 88.1 | 70.0 | 67.9 | 44.3* | 42.4 |
| chrF | 1.22 | 1.42 | 90.1 | 87.2 | 79.7 | 71.0 | 81.4 | 85.8 | 73.8 | 87.7 | 70.4 | 68.2 | 43.3‡ | 44.7‡ |
| chrF++ | 1.23 | 1.42 | 90.2 | 87.3 | 79.7 | 70.9 | 81.3 | 85.7 | 73.8 | 87.7 | 70.3 | 68.3 | 44.0‡ | 43.7† |
| sentBLEU | 1.29 | 1.48 | 90.0 | 87.2 | 79.3 | 70.3* | 81.0† | 85.5 | 73.6 | 87.8 | 69.7 | 67.8 | 44.9 | 42.2 |
| TER | 1.40* | 1.55* | 90.0 | 87.2 | 78.6‡ | 69.8‡ | 80.7‡ | 85.0† | 73.3 | 87.2† | 68.5‡ | 66.5‡ | 44.0† | 41.5 |
| | | | | | | | | | | | | | | |
| rankAvg:all | 0.948‡ | 1.09‡ | 92.2‡ | 89.1‡ | 82.0‡ | 72.6‡ | 81.9‡ | 86.3‡ | 75.4‡ | 87.1† | 68.7‡ | 66.5‡ | 42.0‡ | 44.7‡ |
| rankAvg:qe | 0.868‡ | 0.800‡ | 92.3‡ | 88.9‡ | 83.9‡ | 74.1‡ | 82.6‡ | 85.6 | 74.4† | 85.2‡ | 64.3‡ | 61.8‡ | 35.3‡ | 53.7‡ |
| rankAvg:top | 0.762‡ | 0.822‡ | 93.3‡ | 89.7‡ | 83.3‡ | 73.7‡ | 81.6* | 85.2 | 74.4† | 84.3‡ | 63.1‡ | 60.5‡ | 33.4‡ | 56.0‡ |
| rankAvg:topQe | 0.798‡ | 0.738‡ | 92.5‡ | 88.9‡ | 84.1‡ | 74.3‡ | 81.8‡ | 85.2* | 74.0 | 84.6‡ | 62.8‡ | 60.2‡ | 33.3‡ | 56.3‡ |
| rankAvg:mxmxqe | 0.616‡ | 0.633‡ | 92.1‡ | 87.7* | 79.9 | 70.6 | 80.2‡ | 83.8‡ | 73.2 | 82.2‡ | 58.9‡ | 55.8‡ | 27.1‡ | 63.6‡ |
| rankAvg:noLex | 0.873‡ | 0.964‡ | 92.7‡ | 89.4‡ | 82.8‡ | 73.2‡ | 82.1‡ | 86.1† | 75.6‡ | 86.2‡ | 66.8‡ | 64.4‡ | 38.5‡ | 48.9‡ |
| rankAvg:noNC | 0.964‡ | 1.07‡ | 91.4‡ | 88.1‡ | 80.6‡ | 71.5‡ | 81.5* | 86.2‡ | 75.2‡ | 87.1† | 68.8‡ | 66.7‡ | 42.2‡ | 44.4‡ |
| rankAvg:noNCnoLex | 0.856‡ | 0.931‡ | 91.9‡ | 88.2‡ | 81.0‡ | 71.7‡ | 81.6† | 86.3‡ | 75.7‡ | 86.5‡ | 67.2‡ | 64.9‡ | 39.7‡ | 47.6‡ |
| allQE(32)allMBR | 0.945‡ | 1.09‡ | 92.2‡ | 89.1‡ | 82.1‡ | 72.7‡ | 82.0‡ | 86.2‡ | 75.1‡ | 86.8‡ | 68.0‡ | 65.7‡ | 40.9‡ | 46.1‡ |
| allQE(32)nolexMBR | 0.861‡ | 0.986‡ | 92.8‡ | 89.6‡ | 82.3‡ | 72.8‡ | 81.9‡ | 86.1† | 75.6‡ | 86.2‡ | 66.5‡ | 64.2‡ | 38.7‡ | 48.9‡ |
| topQE(32)topMBR | 0.739‡ | 0.828‡ | 93.6‡ | 89.9‡ | 82.4‡ | 72.8‡ | 81.4 | 85.0* | 74.3* | 84.2‡ | 62.7‡ | 60.0‡ | 33.3‡ | 55.6‡ |
| noncQE(32)noncMBR | 0.825‡ | 0.862‡ | 91.7‡ | 88.3‡ | 81.0‡ | 71.4‡ | 81.6† | 85.9 | 74.9‡ | 86.3‡ | 66.6‡ | 64.2‡ | 38.6‡ | 48.1‡ |
| noncQE(32)noncnolexMBR | 0.774‡ | 0.834‡ | 92.1‡ | 88.4‡ | 81.1‡ | 71.6‡ | 81.4 | 85.8 | 75.4‡ | 85.6‡ | 65.1‡ | 62.5‡ | 36.3‡ | 51.3‡ |
| mxQE(32)mxMBR | 0.552‡ | 0.666‡ | 92.0‡ | 87.9† | 80.0 | 70.4 | 80.2‡ | 83.9‡ | 73.3 | 82.5‡ | 59.2‡ | 56.1‡ | 27.5‡ | 63.5‡ |
| ckQE(32)xcMBR | 0.951‡ | 1.04‡ | 94.2‡ | 88.4‡ | 83.4‡ | 72.2‡ | 81.1 | 84.7‡ | 73.5 | 83.9‡ | 61.7‡ | 59.1‡ | 32.1‡ | 56.8‡ |
| mxQE(32)xcMBR | 0.810‡ | 0.826‡ | 94.3‡ | 88.4‡ | 81.6‡ | 71.4† | 80.9‡ | 84.4‡ | 73.6 | 83.3‡ | 60.2‡ | 57.4‡ | 30.1‡ | 59.6‡ |
| ckQE(32)mxMBR | 0.627‡ | 0.776‡ | 92.7‡ | 88.4‡ | 83.0‡ | 71.9‡ | 81.2 | 84.7‡ | 74.0 | 83.6‡ | 61.3‡ | 58.5‡ | 30.8‡ | 58.6‡ |

Table 21: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-de (WMT2023 dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

| MBR/QE Method | Evaluated Metric MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 2.00 | 1.82 | 87.3 | 89.0 | 76.5 | 68.4 | 79.3 | 85.4 | 74.7 | 88.8 | 68.1 | 66.5 | 46.0 | 39.5 |
| MetricX | 1.31‡ | 1.59 | 89.0‡ | 89.2 | 77.3 | 68.4 | 78.7† | 83.9‡ | 72.0‡ | 83.8‡ | 58.9‡ | 56.7‡ | 31.2‡ | 57.5‡ |
| MetricX-QE | 1.33‡ | 0.839‡ | 89.4‡ | 89.0 | 78.5‡ | 69.3* | 79.3 | 84.3‡ | 72.2‡ | 85.0‡ | 59.4‡ | 57.1‡ | 31.7‡ | 56.2‡ |
| XCOMET-XXL | 1.71* | 1.70 | 93.6‡ | 89.8‡ | 79.4‡ | 69.6‡ | 79.4 | 84.2‡ | 72.3‡ | 85.0‡ | 60.3‡ | 58.1‡ | 33.7‡ | 54.8‡ |
| XCOMET-XL | 1.65‡ | 1.64 | 90.4‡ | 92.4‡ | 79.3‡ | 71.0‡ | 80.0‡ | 85.1 | 74.0 | 86.4‡ | 62.6‡ | 60.7‡ | 36.9‡ | 50.2‡ |
| CometKiwi23-XXL | 1.72‡ | 1.53‡ | 90.2‡ | 90.3‡ | 83.2‡ | 70.8‡ | 80.3‡ | 84.9* | 73.3‡ | 86.4‡ | 62.6‡ | 60.6‡ | 36.8‡ | 50.1‡ |
| CometKiwi23-XL | 1.83 | 1.66* | 89.4‡ | 90.5‡ | 80.0‡ | 73.6‡ | 80.0‡ | 84.5‡ | 72.8‡ | 86.0‡ | 62.0‡ | 60.0‡ | 35.4‡ | 51.9‡ |
| CometKiwi22 | 1.88 | 1.65* | 88.8‡ | 89.9‡ | 78.8‡ | 70.3‡ | 81.8‡ | 85.0 | 73.7‡ | 86.3‡ | 63.3‡ | 61.2‡ | 37.1‡ | 49.2‡ |
| COMET22 | 1.84 | 1.75 | 89.3‡ | 89.8† | 77.8‡ | 69.1‡ | 79.7‡ | 86.2† | 75.1 | 87.2‡ | 65.8‡ | 64.0‡ | 42.1‡ | 42.9‡ |
| BLEURT | 1.66‡ | 1.56‡ | 89.2‡ | 89.9‡ | 78.0‡ | 69.2‡ | 79.8‡ | 85.8† | 76.6‡ | 88.0‡ | 66.2‡ | 64.3‡ | 42.2‡ | 42.5‡ |
| YiSi | 1.98 | 1.77 | 88.0* | 89.4 | 77.5‡ | 68.9† | 79.5 | 85.7* | 75.4‡ | 89.3† | 68.3 | 66.7 | 46.1 | 38.8 |
| chrF | 1.91 | 1.80 | 88.1* | 89.3 | 77.7‡ | 69.0‡ | 79.5* | 85.6 | 75.0 | 89.0 | 69.0 | 67.2 | 45.2 | 41.1 |
| chrF++ | 1.91 | 1.80 | 88.0* | 89.2 | 77.7‡ | 69.0† | 79.5 | 85.7* | 75.2 | 89.1 | 69.1* | 67.5* | 46.3 | 39.7 |
| sentBLEU | 2.00 | 1.79 | 87.4 | 89.2 | 76.9 | 68.7 | 79.3 | 85.4 | 75.1 | 88.9 | 68.0 | 66.5 | 46.5 | 38.5 |
| TER | 2.46† | 2.18† | 86.1‡ | 88.2† | 75.6* | 67.7† | 78.6‡ | 83.5‡ | 73.5‡ | 86.2‡ | 63.8‡ | 62.0‡ | 41.7‡ | 39.5 |
| | | | | | | | | | | | | | | |
| rankAvg:all | 1.73† | 1.53† | 90.2‡ | 90.6‡ | 79.3‡ | 70.3‡ | 80.1‡ | 85.9* | 75.9‡ | 88.4 | 67.8 | 66.1 | 44.9 | 40.0 |
| rankAvg:qe | 1.54‡ | 1.17‡ | 90.9‡ | 90.8‡ | 81.5‡ | 72.1‡ | 81.0‡ | 85.6 | 74.6 | 87.2‡ | 64.4‡ | 62.5‡ | 38.6‡ | 47.2‡ |
| rankAvg:top | 1.47‡ | 1.26‡ | 91.8‡ | 91.2‡ | 81.0‡ | 71.4‡ | 80.1‡ | 85.3 | 74.5 | 86.5‡ | 63.8‡ | 61.9‡ | 38.4‡ | 48.4‡ |
| rankAvg:topQe | 1.46‡ | 1.09‡ | 90.8‡ | 90.7‡ | 81.9‡ | 72.3‡ | 80.3‡ | 85.4 | 74.1 | 86.9‡ | 63.4‡ | 61.4‡ | 37.1‡ | 49.2‡ |
| rankAvg:mxmxqe | 1.36‡ | 1.14‡ | 89.1‡ | 89.0 | 78.0‡ | 69.1* | 79.1 | 84.2‡ | 72.4‡ | 84.4‡ | 60.1‡ | 58.0‡ | 33.0‡ | 54.9‡ |
| rankAvg:noLex | 1.60‡ | 1.41‡ | 91.2‡ | 90.9‡ | 80.2‡ | 70.9‡ | 80.4‡ | 86.0* | 76.0‡ | 88.0* | 66.7* | 64.9† | 42.6‡ | 42.2‡ |
| rankAvg:noNC | 1.78* | 1.56* | 89.2‡ | 89.8‡ | 78.1‡ | 69.3‡ | 79.7* | 85.7 | 75.7‡ | 88.1 | 67.6 | 66.0 | 45.2 | 39.9 |
| rankAvg:noNCnoLex | 1.63‡ | 1.43‡ | 89.7‡ | 89.9‡ | 78.4‡ | 69.5‡ | 79.7† | 85.8 | 75.8‡ | 87.9‡ | 66.6* | 64.9‡ | 43.5‡ | 41.6† |
| allQE(32)allMBR | 1.69‡ | 1.50‡ | 90.4‡ | 90.8‡ | 79.8‡ | 70.5‡ | 80.3‡ | 86.0* | 75.9‡ | 88.5 | 67.9 | 66.3 | 44.7 | 40.5 |
| allQE(32)nolexMBR | 1.56‡ | 1.42‡ | 91.1‡ | 91.0‡ | 80.0‡ | 70.7‡ | 80.4‡ | 86.1† | 76.0‡ | 88.2 | 67.1 | 65.3* | 43.4‡ | 41.6† |
| topQE(32)topMBR | 1.45‡ | 1.30‡ | 92.2‡ | 91.5‡ | 80.4‡ | 70.9‡ | 80.1‡ | 85.4 | 74.5 | 86.6‡ | 64.0‡ | 62.1‡ | 38.7‡ | 47.9‡ |
| noncQE(32)noncMBR | 1.51‡ | 1.24‡ | 89.8‡ | 90.1‡ | 78.6‡ | 69.6‡ | 79.8‡ | 85.8† | 75.4* | 88.2* | 66.4‡ | 64.6‡ | 42.8‡ | 42.0‡ |
| noncQE(32)noncnolexMBR | 1.40‡ | 1.20‡ | 90.2‡ | 90.1‡ | 78.5‡ | 69.6‡ | 79.9‡ | 86.0‡ | 75.6‡ | 88.1† | 65.8‡ | 64.0‡ | 41.9‡ | 43.2‡ |
| mxQE(32)mxMBR | 1.11‡ | 1.08‡ | 89.7‡ | 89.5* | 78.3‡ | 69.2* | 79.2 | 84.7† | 72.8‡ | 85.2‡ | 60.8‡ | 58.6‡ | 33.4‡ | 55.0‡ |
| ckQE(32)xcMBR | 1.67† | 1.61 | 93.2‡ | 90.5‡ | 81.2‡ | 70.3‡ | 80.0‡ | 84.7* | 73.3‡ | 85.8‡ | 62.2‡ | 60.1‡ | 36.0‡ | 50.9‡ |
| mxQE(32)xcMBR | 1.43‡ | 1.17‡ | 93.3‡ | 90.3‡ | 79.7‡ | 69.9‡ | 79.7† | 84.9 | 73.3‡ | 85.9‡ | 61.8‡ | 59.6‡ | 35.1‡ | 52.3‡ |
| ckQE(32)mxMBR | 1.25‡ | 1.31‡ | 90.7‡ | 90.1‡ | 80.6‡ | 70.0‡ | 79.8† | 85.0 | 73.8† | 86.0‡ | 62.5‡ | 60.5‡ | 36.0‡ | 51.8‡ |

Table 22: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1$^{st}$ block), and ensembles (2$^{nd}$ block), on de-en (WMT2023 dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1$^{st}$ block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

## G.13  Results for English-Chinese (en-zh) on WMT2023 dataset

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 1.36 | 1.24 | 89.4 | 85.9 | 75.5 | 70.1 | 80.2 | 87.0 | 73.3 | 88.0 | 46.9 | 41.0 | 11.6 | 97.1 |
| MetricX | 0.682‡ | 0.711‡ | 92.8‡ | 88.1‡ | 82.4‡ | 74.1‡ | 81.8‡ | 86.3‡ | 71.1‡ | 83.9‡ | 33.2‡ | 28.8‡ | 6.39‡ | 102.‡ |
| MetricX-QE | 0.827‡ | 0.553‡ | 92.2‡ | 87.5‡ | 82.3‡ | 74.0‡ | 81.7‡ | 85.6‡ | 70.0‡ | 83.5‡ | 32.2‡ | 27.8‡ | 6.04‡ | 103.‡ |
| XCOMET-XXL | 0.925‡ | 0.837‡ | 96.2‡ | 88.5‡ | 84.1‡ | 74.4‡ | 81.7‡ | 86.5‡ | 70.9‡ | 84.4‡ | 34.7‡ | 29.9‡ | 6.55‡ | 101.‡ |
| XCOMET-XL | 0.927‡ | 0.864‡ | 93.5‡ | 92.3‡ | 82.3‡ | 75.7‡ | 82.0‡ | 87.1 | 72.4‡ | 85.2‡ | 37.1‡ | 32.1‡ | 6.63‡ | 101.† |
| CometKiwi23-XXL | 0.996‡ | 0.842‡ | 93.5‡ | 87.8‡ | 88.4‡ | 75.3‡ | 82.1‡ | 86.2‡ | 70.6‡ | 84.5‡ | 34.3‡ | 29.8‡ | 6.50‡ | 102.‡ |
| CometKiwi23-XL | 1.02‡ | 0.876‡ | 92.4‡ | 89.1‡ | 83.6‡ | 78.9‡ | 82.3‡ | 86.4‡ | 70.9‡ | 84.6‡ | 34.4‡ | 29.8‡ | 6.31‡ | 101.‡ |
| CometKiwi22 | 0.995‡ | 0.869‡ | 92.2‡ | 88.2‡ | 82.5‡ | 75.0‡ | 84.2‡ | 87.1 | 71.8‡ | 85.2‡ | 35.9‡ | 31.1‡ | 7.23‡ | 103.‡ |
| COMET22 | 1.04‡ | 0.999‡ | 91.6‡ | 87.9‡ | 80.2‡ | 73.3‡ | 81.9‡ | 89.3‡ | 74.0‡ | 87.4‡ | 43.4‡ | 37.9‡ | 10.4‡ | 97.6 |
| BLEURT | 1.08‡ | 1.05‡ | 91.1‡ | 87.6‡ | 79.3‡ | 72.5‡ | 81.4‡ | 87.5‡ | 76.7‡ | 87.2‡ | 42.6‡ | 37.1‡ | 8.42‡ | 98.0 |
| YiSi | 1.29‡ | 1.21* | 89.9‡ | 86.3† | 77.0‡ | 71.0‡ | 80.7‡ | 87.7‡ | 74.2‡ | 89.0‡ | 48.4‡ | 42.3‡ | 11.5 | 96.6 |
| chrF | 1.28‡ | 1.20† | 90.0‡ | 86.3† | 77.1‡ | 71.0‡ | 80.7‡ | 87.8‡ | 74.2‡ | 88.8‡ | 49.6‡ | 43.4‡ | 12.4 | 97.6 |
| chrF++ | 1.28‡ | 1.19† | 89.9‡ | 86.3† | 77.1‡ | 71.0‡ | 80.7‡ | 87.8‡ | 74.2‡ | 88.8‡ | 49.4‡ | 43.6‡ | 12.7† | 97.5 |
| sentBLEU | 1.40 | 1.26 | 88.7‡ | 84.8‡ | 75.4 | 69.7* | 79.9† | 86.5‡ | 71.9‡ | 86.9‡ | 43.3‡ | 38.2‡ | 15.1‡ | 106.‡ |
| TER | 1.43† | 1.26 | 88.5‡ | 84.3‡ | 75.6 | 69.6† | 79.7‡ | 86.0‡ | 71.7‡ | 86.5‡ | 41.8‡ | 36.5‡ | 8.65‡ | 94.3* |
| | | | | | | | | | | | | | | |
| rankAvg:all | 0.931‡ | 0.869‡ | 93.4‡ | 89.7‡ | 83.0‡ | 74.9‡ | 82.4‡ | 88.6‡ | 75.3‡ | 88.2* | 46.7 | 40.9 | 11.6 | 94.5† |
| rankAvg:qe | 0.853‡ | 0.694‡ | 93.9‡ | 89.8‡ | 86.2‡ | 77.2‡ | 83.4‡ | 87.2 | 72.2‡ | 85.1‡ | 36.1‡ | 31.4‡ | 7.69‡ | 100.† |
| rankAvg:top | 0.792‡ | 0.698‡ | 94.9‡ | 90.9‡ | 85.8‡ | 76.9‡ | 82.7‡ | 87.3* | 72.6‡ | 85.2‡ | 37.0‡ | 32.2‡ | 7.70‡ | 99.4* |
| rankAvg:topQe | 0.854‡ | 0.673‡ | 93.8‡ | 89.6‡ | 86.7‡ | 77.5‡ | 82.6‡ | 86.7 | 71.6‡ | 84.7‡ | 35.2‡ | 30.6‡ | 7.23‡ | 100.† |
| rankAvg:mxmxqe | 0.712‡ | 0.608‡ | 93.1‡ | 88.3‡ | 82.9‡ | 74.4‡ | 81.9‡ | 86.3‡ | 71.1‡ | 83.9‡ | 33.4‡ | 29.0‡ | 6.49‡ | 101.‡ |
| rankAvg:noLex | 0.860‡ | 0.802‡ | 94.0‡ | 90.3‡ | 84.2‡ | 75.9‡ | 82.9‡ | 88.5‡ | 75.2‡ | 87.5‡ | 43.7‡ | 38.1‡ | 10.1‡ | 96.2 |
| rankAvg:noNC | 0.997‡ | 0.920‡ | 91.9‡ | 88.3‡ | 80.3‡ | 73.1‡ | 81.8‡ | 88.5‡ | 75.4‡ | 88.5‡ | 47.8† | 41.9† | 11.9 | 94.0† |
| rankAvg:noNCnoLex | 0.911‡ | 0.835‡ | 92.4‡ | 88.8‡ | 81.4‡ | 73.8‡ | 82.1‡ | 88.6‡ | 75.4‡ | 88.0 | 45.0‡ | 39.3‡ | 10.5‡ | 96.5 |
| allQE(32)allMBR | 0.901‡ | 0.841‡ | 93.6‡ | 89.8‡ | 83.4‡ | 75.2‡ | 82.6‡ | 88.5‡ | 75.0‡ | 87.7† | 44.9‡ | 39.3‡ | 10.8 | 95.2 |
| allQE(32)nolexMBR | 0.868‡ | 0.813‡ | 94.0‡ | 90.2‡ | 83.9‡ | 75.6‡ | 82.7‡ | 88.5‡ | 74.9‡ | 87.4‡ | 43.4‡ | 38.0‡ | 10.2‡ | 96.4 |
| topQE(32)topMBR | 0.778‡ | 0.730‡ | 95.1‡ | 91.0‡ | 84.8‡ | 76.2‡ | 82.4‡ | 87.3* | 72.6‡ | 85.3‡ | 37.2‡ | 32.4‡ | 8.18‡ | 99.0 |
| noncQE(32)noncMBR | 0.915‡ | 0.781‡ | 92.6‡ | 88.7‡ | 81.7‡ | 73.9‡ | 82.1‡ | 88.1‡ | 74.5‡ | 87.5‡ | 44.3‡ | 38.7‡ | 9.85‡ | 97.2 |
| noncQE(32)noncnolexMBR | 0.862‡ | 0.755‡ | 92.9‡ | 89.0‡ | 82.3‡ | 74.3‡ | 82.3‡ | 88.2‡ | 74.6‡ | 87.2‡ | 42.8‡ | 37.3‡ | 9.37‡ | 97.8 |
| mxQE(32)mxMBR | 0.695‡ | 0.664‡ | 93.1‡ | 88.3‡ | 82.7‡ | 74.2‡ | 81.8‡ | 86.3‡ | 71.2‡ | 84.0‡ | 33.6‡ | 29.1‡ | 6.46‡ | 102.‡ |
| ckQE(32)xcMBR | 0.912‡ | 0.823‡ | 96.0‡ | 88.7‡ | 85.8‡ | 75.0‡ | 82.0‡ | 86.6‡ | 71.2‡ | 84.5‡ | 35.1‡ | 30.4‡ | 7.20‡ | 100.* |
| mxQE(32)xcMBR | 0.860‡ | 0.724‡ | 95.9‡ | 88.8‡ | 84.4‡ | 74.8‡ | 82.1‡ | 86.5‡ | 71.2‡ | 84.4‡ | 34.8‡ | 30.1‡ | 6.93‡ | 101.† |
| ckQE(32)mxMBR | 0.723‡ | 0.720‡ | 93.5‡ | 88.6‡ | 85.2‡ | 74.9‡ | 82.1‡ | 86.7* | 71.7‡ | 84.5‡ | 34.8‡ | 30.2‡ | 7.37‡ | 100.† |

Table 23: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on en-zh (WMT2023 dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

| MBR/QE Method | MetricX | MetricX-QE | XCOMET-XXL | XCOMET-XL | CometKiwi23-XXL | CometKiwi23-XL | CometKiwi22 | COMET22 | BLEURT | YiSi | chrF | chrF++ | sentBLEU | TER |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Greedy | 2.19 | 2.00 | 90.9 | 88.1 | 78.4 | 70.3 | 79.6 | 82.7 | 71.4 | 82.5 | 55.0 | 52.4 | 26.2 | 65.2 |
| MetricX | 1.01‡ | 1.20‡ | 93.9‡ | 89.6‡ | 78.7 | 70.3 | 79.3* | 82.1‡ | 70.8† | 80.4‡ | 49.5‡ | 46.6‡ | 19.1‡ | 78.9‡ |
| MetricX-QE | 1.25‡ | 0.853‡ | 93.4‡ | 89.3‡ | 79.3‡ | 71.0‡ | 80.0‡ | 81.9‡ | 70.7‡ | 80.5‡ | 49.0‡ | 46.2‡ | 19.0‡ | 78.2‡ |
| XCOMET-XXL | 1.41‡ | 1.38‡ | 96.2‡ | 89.8‡ | 80.2‡ | 70.9‡ | 79.7 | 81.8‡ | 70.6‡ | 80.1‡ | 48.7‡ | 45.8‡ | 18.2‡ | 78.6‡ |
| XCOMET-XL | 1.45‡ | 1.41‡ | 93.9‡ | 92.4‡ | 80.3‡ | 72.5‡ | 80.7‡ | 83.1‡ | 72.5‡ | 81.9‡ | 52.4‡ | 49.7‡ | 22.3‡ | 72.8‡ |
| CometKiwi23-XXL | 1.69‡ | 1.52‡ | 93.6‡ | 89.9‡ | 85.1‡ | 72.7‡ | 80.8‡ | 82.8 | 71.6 | 81.6‡ | 51.8‡ | 49.0‡ | 21.4‡ | 72.7‡ |
| CometKiwi23-XL | 1.74‡ | 1.59‡ | 92.9‡ | 90.3‡ | 81.3‡ | 75.9‡ | 81.0‡ | 82.7 | 71.5 | 81.5‡ | 51.8‡ | 49.0‡ | 21.5‡ | 74.1‡ |
| CometKiwi22 | 1.74‡ | 1.51‡ | 92.8‡ | 89.7‡ | 80.7‡ | 72.4‡ | 82.8‡ | 83.1† | 72.2‡ | 81.9† | 52.1‡ | 49.4‡ | 22.2‡ | 72.6‡ |
| COMET22 | 1.73‡ | 1.63‡ | 92.5‡ | 89.6‡ | 79.8‡ | 71.6‡ | 80.6‡ | 84.6‡ | 72.6‡ | 82.8 | 54.8 | 52.0 | 25.1‡ | 65.7 |
| BLEURT | 1.68‡ | 1.60‡ | 92.8‡ | 89.9‡ | 79.9‡ | 71.5‡ | 80.5‡ | 83.4‡ | 74.4‡ | 82.5 | 54.1‡ | 51.3‡ | 24.5‡ | 67.8‡ |
| YiSi | 2.12* | 1.95*‡ | 91.2† | 88.4* | 78.9‡ | 70.5 | 79.8† | 83.0‡ | 71.8‡ | 83.7‡ | 55.9‡ | 53.3† | 26.6 | 64.8 |
| chrF | 2.16 | 2.01 | 91.1 | 88.1 | 78.8‡ | 70.7‡ | 79.8‡ | 82.9* | 71.8‡ | 83.2‡ | 56.9‡ | 54.1‡ | 26.5 | 67.4‡ |
| chrF++ | 2.14 | 1.99 | 91.1* | 88.2 | 78.9† | 70.8‡ | 79.9‡ | 83.0‡ | 72.0‡ | 83.3‡ | 56.9‡ | 54.3‡ | 26.9* | 66.6† |
| sentBLEU | 2.17 | 2.01 | 91.0 | 88.0 | 78.5 | 70.2 | 79.6 | 82.9* | 71.5 | 83.3‡ | 56.0‡ | 53.5‡ | 27.7‡ | 62.8‡ |
| TER | 2.28† | 2.11‡ | 90.5* | 87.3‡ | 77.2‡ | 69.2‡ | 78.8‡ | 82.1‡ | 70.5‡ | 82.4 | 53.1‡ | 50.5‡ | 25.3* | 60.4‡ |
| | | | | | | | | | | | | | | |
| rankAvg:all | 1.57‡ | 1.44‡ | 93.6‡ | 90.5‡ | 81.3‡ | 72.6‡ | 81.0‡ | 84.0‡ | 73.4‡ | 83.5‡ | 56.4‡ | 53.7‡ | 27.0† | 63.7‡ |
| rankAvg:qe | 1.44‡ | 1.16‡ | 94.0‡ | 90.9‡ | 83.3‡ | 74.4‡ | 81.9‡ | 83.6‡ | 73.0‡ | 82.2 | 53.2‡ | 50.5‡ | 23.2‡ | 71.1‡ |
| rankAvg:top | 1.26‡ | 1.17‡ | 95.0‡ | 91.6‡ | 82.7‡ | 73.9‡ | 81.1‡ | 83.5‡ | 72.9‡ | 82.1* | 53.1‡ | 50.4‡ | 22.7‡ | 71.8‡ |
| rankAvg:topQe | 1.39‡ | 1.11‡ | 94.0‡ | 90.8‡ | 83.4‡ | 74.5‡ | 81.1‡ | 83.3‡ | 72.6‡ | 82.0† | 52.7‡ | 49.9‡ | 22.4‡ | 71.9‡ |
| rankAvg:mxmxqe | 1.07‡ | 0.972‡ | 93.9‡ | 89.9‡ | 79.6‡ | 71.1‡ | 80.0‡ | 82.3* | 71.4 | 80.8‡ | 50.5‡ | 47.6‡ | 20.0‡ | 76.9‡ |
| rankAvg:noLex | 1.43‡ | 1.30‡ | 94.2‡ | 91.1‡ | 82.1‡ | 73.3‡ | 81.4‡ | 84.1‡ | 73.7‡ | 83.2‡ | 55.4 | 52.7 | 25.7 | 66.8† |
| rankAvg:noNC | 1.64‡ | 1.47‡ | 92.8‡ | 89.7‡ | 80.0‡ | 71.4‡ | 80.4‡ | 83.7‡ | 73.0‡ | 83.4‡ | 56.2‡ | 53.5‡ | 26.9* | 63.2‡ |
| rankAvg:noNCnoLex | 1.45‡ | 1.30‡ | 93.4‡ | 90.2‡ | 80.4‡ | 71.8‡ | 80.7‡ | 84.0‡ | 73.4‡ | 83.2‡ | 55.4 | 52.7 | 25.8 | 65.6 |
| allQE(32)allMBR | 1.52‡ | 1.41‡ | 93.8‡ | 90.8‡ | 81.6‡ | 72.8‡ | 81.2‡ | 84.1‡ | 73.6‡ | 83.5‡ | 56.3‡ | 53.7‡ | 27.1† | 64.4 |
| allQE(32)nolexMBR | 1.40‡ | 1.33‡ | 94.3‡ | 91.2‡ | 81.7‡ | 73.0‡ | 81.3‡ | 84.2‡ | 73.8‡ | 83.3‡ | 55.7† | 53.1* | 26.3 | 65.8 |
| topQE(32)topMBR | 1.22‡ | 1.21‡ | 95.1‡ | 91.8‡ | 81.8‡ | 73.1‡ | 80.9‡ | 83.4‡ | 72.9‡ | 82.1* | 53.1‡ | 50.4‡ | 23.1‡ | 71.5‡ |
| noncQE(32)noncMBR | 1.47‡ | 1.23‡ | 93.3‡ | 90.1‡ | 80.4‡ | 71.8‡ | 80.7‡ | 83.7‡ | 73.1‡ | 82.9* | 54.9 | 52.2 | 25.3† | 66.0 |
| noncQE(32)noncnolexMBR | 1.36‡ | 1.18‡ | 93.6‡ | 90.4‡ | 80.6‡ | 71.9‡ | 80.8‡ | 83.9‡ | 73.4‡ | 82.7 | 54.2† | 51.5‡ | 24.4‡ | 67.6‡ |
| mxQE(32)mxMBR | 1.04‡ | 1.08‡ | 93.9‡ | 89.7‡ | 79.3‡ | 70.8† | 79.7 | 82.3† | 71.1* | 80.5‡ | 49.8‡ | 46.9‡ | 19.6‡ | 77.8‡ |
| ckQE(32)xcMBR | 1.42‡ | 1.36‡ | 95.8‡ | 90.3‡ | 82.8‡ | 72.0‡ | 80.5‡ | 82.5 | 71.6 | 81.0‡ | 50.6‡ | 47.7‡ | 20.4‡ | 75.4‡ |
| mxQE(32)xcMBR | 1.32‡ | 1.16‡ | 95.8‡ | 90.0‡ | 80.5‡ | 71.4‡ | 80.0‡ | 82.2† | 71.3 | 80.5‡ | 49.7‡ | 46.7‡ | 19.4‡ | 77.3‡ |
| ckQE(32)mxMBR | 1.12‡ | 1.21‡ | 94.4‡ | 90.5‡ | 82.4‡ | 72.1‡ | 80.5‡ | 83.0† | 72.2‡ | 81.5‡ | 51.8‡ | 49.0‡ | 21.5‡ | 73.7‡ |

Table 24: Reference-based and QE evaluation scores for greedy and MBR/QE decoding (1st block), and ensembles (2nd block), on zh-en (WMT2023 dataset). Higher scores are better, except MetricX, MetricX-QE, and TER, where lower is better. Green is better than greedy, red is worse. Ensembles are defined in Table 2. Significant differences from greedy (pairwise t-test) indicated by * for p<0.05, † for p<0.01, ‡ for p<0.001. The green diagonal in the 1st block shows metrics prefer outputs from MBR/QE decoding using the same utility metric.

# Beyond Human-Only: Evaluating Human-Machine Collaboration for Collecting High-Quality Translation Data

**Zhongtao Liu, Parker Riley, Daniel Deutsch, Alison Lui,**
**Mengmeng Niu, Apu Shah and Markus Freitag**
Google
{zhongtao,prkriley,dandeutsch,alisonlui,mniu,apurva,freitag}@google.com

## Abstract

Collecting high-quality translations is crucial for the development and evaluation of machine translation systems. However, traditional human-only approaches are costly and slow. This study presents a comprehensive investigation of 11 approaches for acquiring translation data, including human-only, machine-only, and hybrid approaches. Our findings demonstrate that human-machine collaboration can match or even exceed the quality of human-only translations, while being more cost-efficient. Error analysis reveals the complementary strengths between human and machine contributions, highlighting the effectiveness of collaborative methods. Cost analysis further demonstrates the economic benefits of human-machine collaboration methods, with some approaches achieving top-tier quality at around 60% of the cost of traditional methods. We release a publicly available dataset[1] containing nearly 18,000 segments of varying translation quality with corresponding human ratings to facilitate future research.

## 1 Introduction

Collecting high-quality translations efficiently presents significant challenges. Traditional approaches rely heavily on different tiers of human translators, ranging from professional linguists to junior bilingual speakers (Zouhar and Bojar, 2024). While these approaches can produce high-quality translations, they are often expensive, time-consuming, and challenging to scale for large datasets.

Recent advancements in machine translation with large language models (OpenAI, 2024; Gemini, 2024) have demonstrated models' impressive abilities to generate human-like translations. However, recent research (Yan et al., 2024) tends to



Figure 1: Our 11 translation systems, organized by initial translation type (human or machine) and post-editing type (none, human, or machine). Detailed system descriptions are provided in Section 2.2.

position human translators and machine translation systems as competitors rather than potential collaborators, which could result in efficient alternatives for addressing the limitations of traditional translation data collection methods.

In this paper, we aim to fill the gap by comprehensively investigating the potential of human-machine collaboration to efficiently collect high-quality translation data. We hypothesize that combining the strengths of humans and machines could lead to higher quality, cost-efficient translation collection methods. To verify the hypothesis, we explore 11 different methods for acquiring translation data, including human-only, machine-only, and various hybrid methods.

Our research seeks to answer the following key questions:

- Can human-machine collaborative approaches produce translations of comparable or higher quality than traditional human-only or machine-only methods?

- How do different collaborative methods impact translation quality, and where do the improvements primarily originate?

- What are the cost implications of these various approaches, and can human-machine collabo-

---

[1]The dataset can be found at https://github.com/google-research/google-research/tree/master/collaborative-tr-collection.

ration offer a more cost-efficient solution for high-quality translation collection?

Our findings demonstrate that human-machine collaboration can match or even exceed human-only translation quality while being more cost-efficient. We present detailed error analyses to reveal the complementary strengths of the collaborative methods and conduct a thorough cost analysis to illustrate the economic benefits of collaborative approaches.

To support future research, we also release a publicly available dataset containing nearly 18,000 segments of varying translation quality with corresponding human ratings.

## 2 Collecting Translations

Translating text from one language to another can either be done by bilingual annotators or machine translation systems. However, both cases are prone to producing errors in their translations, including well-trained expert translators (Freitag et al., 2023). As such, translations can be post-edited, a process of correcting a translation, either manually or with a model, that often yields higher-quality translations.

Both steps of this process — the initial translation collection and the post-editing — can either be done with humans or with models, each with their own advantages and disadvantages in terms of speed, quality, cost, and scalability. In this work, we explore how combinations of human and machines for both steps of this pipeline can combine to produce high-quality translations.

### 2.1 Data Sources

We use the test sets provided by the WMT23 General MT Shared Task (Kocmi et al., 2023) and collect new translations using several methods. These data sets comprise 460 English-German (EnDe) paragraph-level segments and 1175 Chinese-English (ZhEn) sentence-level segments with human rating annotations.

### 2.2 Data Collection Systems

Figure 1 illustrates the combinations from the two dimensions: initial translation and post-editing methods from either human annotators or machines. Machines may be either large language models (LLMs) or machine translation (MT) systems. This results in the 11 systems in the figure, named ac-

cording to the source of the initial translation with a suffix representing the post-editing approach.

In this work, we use several different sources for the initial translation:

- **OrigHumanRef** and **HumanRef** are human translations collected by professional translators. We refer to the original reference provided by the WMT23 General MT Shared Task (Kocmi et al., 2023) as ORIGHUMAN-REF. We collected a new from-scratch professional translation HUMANREF following the standard annotation steps.

- **BestWMT** is the top-ranked MT system picked from the official results of WMT23 General Translation Task: GPT4-5shot (OpenAI, 2024) for EnDe and Lan-BridgeMT (Wu and Hu, 2023) for ZhEn, representing the state-of-the-art MT capability we can access.

- **MidWMT** is a middle-ranked MT system from the official results of WMT23 General Translation Task: ONLINE-G for both EnDe and ZhEn, representing the conventional MT quality we can use.

We additionally explore the following different methods for post-editing translations:

- **HumanPE** refers to the post-edit service provided by a separate batch of linguists. **HumanPEx2** means the translation going through two independent rounds of post-edits from professional translators. There is no translator overlap between the two batches.

- **LLMRefine** (Xu et al., 2024) is one of the state-of-the-art post-edit approaches leveraging error feedback for pin-pointing corrections. Here we reproduced its error-feedback process and leverage Gemini-1.0 Ultra (Gemini, 2024) with the reported prompts to generate post-edited text.

### 2.3 Evaluation

In this paper, we use Multidimensional Quality Metrics (MQM; Lommel et al., 2014; Freitag et al., 2021) to evaluate translation quality. MQM is the state-of-the-art human evaluation framework for MT. In MQM, expert raters identify error spans within translations, which are automatically converted to numeric scores. Lower scores indicate fewer errors and thus higher quality.

Figure 2: Cross-BLEU scores for different EnDe translation collection approaches.

## 3 Data Quality Analysis

In this section, we seek to understand how our collected translations differ from each other (§3.1), how those differences correspond to changes in quality (§3.2), and what those results indicate about the value of human-machine collaboration in terms of quality (§3.3).

### 3.1 Lexical Overlap and Similarity

Figure 2 presents a cross-BLEU (Freitag et al., 2022) similarity matrix for English-German translations, which measures lexical similarity between pairs of translations. See Figure 8 in Appendix A.1 for Chinese-English results. Higher scores indicate greater similarity.

One prominent pattern in these results is that systems based on the same initial translation retain high similarity even after post-editing. This indicates that **post-editing still preserves some characteristics of the original translation**. Also, translation systems based on MT (BESTWMT, MIDWMT, and their post-edited versions) are more similar to each other than to translations based on an initial human translation.

Table 1 presents a subset of the information in Figures 2 and 8, to emphasize the interaction between using human- vs. model-based approaches for the initial translation and post-edit. This illustrates the trend that **humans and machines tend to make more changes to translations from the other group**.

| Source | + HUMANPE | + LLMREFINE |
|---|---|---|
| EN-DE | | |
| HUMANREF | 95 | **72** |
| BESTWMT | **81** | 83 |
| MIDWMT | **70** | 77 |
| ZH-EN | | |
| HUMANREF | 88 | **79** |
| BESTWMT | **84** | 89 |
| MIDWMT | **68** | 71 |

Table 1: Cross-BLEU score comparison between different post-edited versions of the same translation. Lower numbers indicate less similarity and more changes from the initial translation.

| Reference | MQM per segment |
|---|---|
| HUMANREF+LLMREFINE | 2.76 |
| ORIGHUMANREF | 2.80 |
| BESTWMT+HUMANPE | 2.98 |
| BESTWMT | 3.30 |
| BESTWMT+LLMREFINE | 3.36 |
| HUMANREF+HUMANPE | 3.53 |
| HUMANREF+HUMANPEx2 | 3.70 |
| HUMANREF | 3.79 |
| MIDWMT+HUMANPE | 3.83 |
| MIDWMT+LLMREFINE | 4.02 |
| MIDWMT | 6.45 |

Table 2: English-German MQM human evaluation results. Lower scores represent higher translation quality.

### 3.2 MQM Quality Evaluation

Tables 2 and 3 present the MQM human evaluation results. The solid lines denote *significance clusters*, where every system in a cluster is statistically significantly better than every system below that cluster, based on random permutation tests with 10,000 trials, where a $p$-value of less than $\alpha = 0.05$ is considered significant.

The results reveal that **HumanRef+LLMRefine** and **BestWMT+HumanPE** are the overall winners, with each appearing in the best significance cluster in both language pairs.

Tables 2 and 3 show that post-edits, both HUMANPE and LLMREFINE, **demonstrate a positive impact on initial translations**. These methods consistently either elevate the translation quality to a higher level of significance or preserve the existing quality.

| Reference | MQM per segment |
|---|---|
| HUMANREF+LLMREFINE | 1.82 |
| HUMANREF+HUMANPEx2 | 1.82 |
| BESTWMT+HUMANPE | 1.87 |
| HUMANREF+HUMANPE | 1.91 |
| BESTWMT+LLMREFINE | 1.94 |
| HUMANREF | 2.05 |
| BESTWMT | 2.22 |
| MIDWMT+HUMANPE | 2.23 |
| MIDWMT+LLMREFINE | 2.45 |
| MIDWMT | 3.98 |
| ORIGHUMANREF | 5.63 |

Table 3: Chinese-English MQM human evaluation results. Lower scores represent higher translation quality.

### 3.3 Human-Machine Collaboration

Table 4 and Figure 3 present a detailed analysis that highlights the quality benefits of human-machine collaboration.

Figure 3 shows the gains in quality that can be provided by our post-edit approaches. The gains are more pronounced when starting with a lower-quality translation (MIDMT), but even high-quality translations (HUMANREF, BESTWMT) can be improved. The quality differences between initial translations are greatly reduced after post-editing, but not eliminated.

Table 4 provides a finer-grain analysis of the effect of each post-edit approach on different initial translations. In both language pairs, an LLM-based method provides the most benefit when starting with human translation, while human post-editing provides the most benefit for machine translation. Recall that in Table 1 we showed that model-based methods and humans make more changes to translations from the other group; here we see that these changes are also net-positive. This indicates that **human-machine collaboration is an effective way to achieve high-quality translations**.

## 4 Error Analysis

Here we present more detailed error analysis of both initial translation and post-edit stages to investigate where the quality improvements originate. Using English-German as an example, we first present analysis of the initial translations in Section 4.1 to understand the initial error distributions. Then, we further investigate the error-correction dynamics during post-editing in Section 4.2 to un-



Figure 3: MQM Scores for different translation systems across two language pairs: Chinese-English and English-German. Bars represents the average MQM scores for each translation system. The systems are grouped and colored by initial translation and further categorized by post-editing method with different fill patterns. Lower MQM scores indicate better quality.

derstand why human-machine collaboration stands out from other approaches.

### 4.1 Error distribution from Initial Translation

We present the error type and severity distribution of the English-German initial translations in Table 5 and that of Chinese-English in Appendix A.2. It shows that accuracy-related errors are the primary source of major errors and that the distribution of minor errors is more evenly spread across different categories. Importantly, these error distributions are similar for both human and machine-based initial translations. This consistency provides a solid foundation for comparing different post-editing techniques in our downstream analysis.

### 4.2 Error Correction from Post-Editing

To understand the error-correction dynamics of different post-editing methods and find the origin of the improvements of human-machine collaboration, we explored three key questions:

- Do different post-editing methods agree on which segments to modify?

- How do different post-editing methods affect the total number of major and minor errors?

- How do different post-editing methods affect

| Source | Init. Translation | + HumanPE | | + HumanPEx2 | | + LLMRefine | |
|---|---|---|---|---|---|---|---|
| | | Score ↓ | Δ | Score ↓ | Δ | Score ↓ | Δ |
| EN-DE | | | | | | | |
| HUMANREF | 3.79 | 3.53 | -0.26 | 3.70 | -0.09 | **2.76** | **-1.03** |
| BESTWMT | 3.30 | **2.98** | **-0.32** | - | - | 3.36 | +0.06 |
| MIDWMT | 6.45 | **3.83** | **-2.62** | - | - | 4.02 | -2.43 |
| ZH-EN | | | | | | | |
| HUMANREF | 2.05 | 1.91 | -0.14 | **1.82** | **-0.23** | 1.82 | -0.23 |
| BESTWMT | 2.22 | **1.87** | **-0.35** | - | - | 1.94 | -0.28 |
| MIDWMT | 3.98 | **2.23** | **-1.75** | - | - | 2.45 | -1.53 |

Table 4: MQM human evaluation comparison of each post-edit approach on different initial translations. Lower MQM scores indicates better quality.

| Error Type | OrigHumanRef | HumanRef | BestWMT | MidWMT |
|---|---|---|---|---|
| **No-error** | 200 | 175 | 144 | 116 |
| **Major** | 186 | 263 | 211 | 470 |
| Fluency | 27 (15%) | 29 (11%) | 25 (12%) | 88 (19%) |
| Accuracy | 114 (61%) | 149 (57%) | 108 (51%) | 275 (59%) |
| Style | 22 (12%) | 54 (21%) | 45 (21%) | 60 (13%) |
| **Minor** | 659 | 745 | 891 | 1084 |
| Fluency | 257 (39%) | 242 (32%) | 439 (49%) | 538 (50%) |
| Accuracy | 181 (27%) | 211 (28%) | 186 (21%) | 226 (21%) |
| Style | 141 (21%) | 190 (26%) | 178 (20%) | 213 (20%) |

Table 5: Error type and severity distributions of English-German MQM human evaluation results.



Figure 4: Agreement between HumanPE and LLM-Refine in identifying segments requiring post-edit on English-German data. Each pie chart[2] represents a different initial translation source.

the total number of high- and low-quality segments?

We first examine how often human post-editors (HUMANPE) and machine post-editing methods (LLMREFINE) agree on which segments need correction. Figure 4 shows that both methods identify more segments for editing in lower-quality initial translations as evidenced by the shrinking "No Change" (yellow) section. Notably, agreement between HUMANPE and LLMREFINE increased from 23.9% for high-quality HUMANREF translations to 67.4% for lower-quality MidWMT translations as observed by the expanded "HumanPE & LLMRefine" (purple) section. This suggests more consensus on obvious errors in lower-quality texts while greater divergence for higher-quality translations in editing approaches. A detailed numerical breakdown is shown in Table 8 in Appendix A.3.

Another interesting observation from Figure 4 is that HUMANPE (purple + red) identifies a larger proportion of segments needing correction in BEST-WMT (48.9% + 19.3% = 68.2%) compared to HUMANREF (23.9% + 11.5% = 35.4%), despite

the superior quality of BESTWMT over HUMAN-REF as shown in Table 2. This suggests that human post-editors might overlook certain errors in human translations due to their familiar patterns. Conversely, the unfamiliar patterns in machine-generated text may make errors more salient to human editors. This interpretation is consistent with the pattern depicted in Figure 2, where human post-editors make fewer changes to human translations than to machine translations. This observation provides one plausible explanation for the necessity of human-machine collaboration in achieving high-quality translations. Detailed statistics for English-German are shown in Table 8 and similar trends are also observed in Chinese-English in Table 9 in Appendix A.3.

We wish to investigate how many errors are corrected during post-editing, but because it is difficult to automatically determine whether an individual error was corrected, we instead examine how the total number of errors changes after post-editing, also considering severity. Figure 5 shows that both human and machine post-editing reduce overall error counts across different initial translation qual-

Figure 5: Error changes percentages by different post-editing approaches on English-German data. The percentages present the changes in error counts for each post-editing method compared to its initial translation. A negative indicates a decrease in errors, while positive value indicates an increase in the error type.



(a) From HUMANREF to HUMANREF+HUMANPE



(b) From HUMANREF to HUMANREF+LLMREFINE

Figure 6: Segment-level quality shift through HU-MANPE and LLMREFINE from English-German HU-MANREF. Each segment is categorized into one of three groups based on its MQM score: 1) high-scoring segments with MQM >= 5; 2) low-scoring segments with 0 < MQM < 5; 3) error-free segments with MQM=0. Higher MQM scores indicate more numerous/severe errors and accordingly lower translation quality.

ities. LLMREFINE outperforms HUMANPE on HUMANREF initial translation in reducing major errors (-27.8% vs. -2.7%), while HUMANPE is superior for BESTWMT, decreasing major-error segments compared to the increase for LLMREFINE (-4.3% vs. +8.5%). On MIDWMT, both methods show substantial improvements, with HUMANPE moderately ahead of LLMREFINE (-41.3% vs. -37.0% decrease in major-error segments). These findings highlight the complementary strengths of human and machine post-editing methods, indicating that a hybrid method is likely the most effective strategy for reducing errors, regardless of the initial translation's origin. Similar trends are also observed in Chinese-English in Figure 9 in Appendix A.3.

To understand the error correction dynamics for each segment, we analyzed how MQM scores change before and after post-editing. Ideally, post-editing would fix existing errors while minimizing the introduction of new ones. However, as Figure 6 illustrates, post-editing is not guaranteed to improve every segment: while some segments are improved, others are worsened.

Figure 6 compares HUMANPE and LLM-REFINE on HUMANREF initial translations for English-German data. Both methods reduce the number of high-scoring (low-quality) segments (MQM >= 5). Notably, LLMREFINE outperforms HUMANPE by showing fewer quality-degrading corrections and more quality-improving ones. LLMREFINE minimizes low-to-high-scoring

transitions with a narrower flow from low-scoring segments (0 < MQM < 5) to high-scoring segments compared to HUMANPE. Moreover, LLM-REFINE achieves a significant reduction in high-scoring segments by 6.5% (from 29.1% to 22.6%) compared to HUMANPE's 1.3% (from 29.1% to 27.8%), suggesting that it is more effective at achieving post-editing gains while preserving originally good translations. A similar trend is observed for Chinese-English with HUMANREF initial translation in Figure 10 in Appendix A.3.

To demonstrate that the quality improvement is not solely due to the capabilities of LLMRE-FINE, we conducted further experiments as shown in Figure 7. This figure compares HUMANPE and LLMREFINE with BESTWMT initial translations. Interestingly, the results are reversed: HUMANPE

(a) From BESTWMT to BESTWMT+HUMANPE



(b) From BESTWMT to BESTWMT+LLMREFINE

Figure 7: Segment-level quality shift through HU-MANPE and LLMREFINE from English-German BESTWMT.

outperforms LLMREFINE in this scenario, showing fewer quality-degrading corrections and more quality-improving ones. HUMANPE demonstrates a significantly wider flow from high-scoring segments to low-scoring ones. It achieves a notable reduction in high-scoring segments by 3.2% (from 26.5% to 23.3%), while LLMREFINE sees an increase of 0.7% (from 26.5% to 27.2%). Furthermore, HUMANPE significantly increases the No-Error segments (MQM=0) by 12.2% (from 33.5% to 45.7%) compared to LLMREFINE's 2.2% increase.

The distinct performance differences of LLM-REFINE and HUMANPE in these two sets of experiments highlight that **the quality improvement stems primarily from the complementary strengths of human and machine collaboration**, rather than the superior capability of either LLM-REFINE or HUMANPE alone. This underscores the importance of leveraging both human and machine strengths in achieving optimal translation quality.

| Systems | Quality Rank | | Costs |
| --- | --- | --- | --- |
| | EnDe | ZhEn | |
| HUMANREF | 3 | 2 | 1X |
| HUMANREF+HUMANPE | 2 | 1 | 1.6X |
| HUMANREF+HUMANPEx2 | 2 | 1 | 2.2X |
| HUMANREF+LLMREFINE | **1** | **1** | 1X |
| BESTWMT+HUMANPE | **1** | **1** | **0.6X** |

Table 6: Quality rank and costs comparison of different data collection systems. 1st rank indicates the translation quality belongs to the highest quality significance cluster in Table 2 and 3.

## 5 Costs Analysis

So far we have focused on comparing quality between various translation data collection approaches. However, practical considerations make it important to consider the trade-off between quality and costs. Table 6 analyzes relative human annotation costs between various approaches, along with the rank of the significance cluster that each method appeared in. The exact costs for the human annotation conducted in this study are confidential (although all annotators were paid fair market wages), so we instead use relative costs, based on the industry standard that post-editing text of a given length takes less time (and accordingly costs less) than producing a translation of that length. We specifically assume that human post-editing costs around 60% of what human translation does. According to existing literature (Plitt and Masselot, 2010; Zouhar et al., 2021; Green et al., 2013) and internal statistics, we believe it's a fair assumption, although the exact costs can vary upon different vendors, languages, task size, etc.

With this framework, the best combination of quality and cost appears to be human post-editing of high-quality MT (BESTWMT+HUMANPE), attaining quality in the top significance cluster in both language pairs with only 60% of the human annotation cost of collecting an initial human translation. Meanwhile, we see that one or two rounds of human post-editing of an initial human translation increases costs without a meaningful gain in quality, while just applying an LLM post-editor (HUMANREF+LLMREFINE) brings quality to the top significant cluster with no additional human annotation cost, making it a viable option when human translations are already collected. It's worth noting that LLM inference costs are negligible (on the order of dollars per million tokens) compared to human annotation costs, further enhancing the

cost-effectiveness of LLM-based approaches. This indicates that **human-machine collaboration can be a faster, more cost-efficient alternative** to traditional collection of translations from humans, optimizing both quality and resource allocation by leveraging the strengths of both humans and machines.

# 6    Related Work

There have been a few studies investigating methods of acquiring high-quality translations. Recently, Zouhar et al. (2024) proposed collecting high-quality translations by building consensus between multiple translators. Zouhar and Bojar (2024) proposed collecting multiple translations from different tiers of human translators with careful budget calculations to optimize cost-efficiency.

**Human Post-Edits**    Computer-aided translation tools are now widely used by professional translators for interactive translation and post-editing (Alabau et al., 2014; Federico et al., 2014; Green et al., 2014; Denkowski, 2015; Sin-wai, 2014; Kenny, 2012). Carl et al. (2011) have shown that human translators work faster and make fewer mistakes when editing machine translations than when translating from scratch. Toral et al. (2018) supports this, demonstrating even greater improvements with neural machine translation compared to phrase-based systems. Zouhar et al. (2021) investigates the relationship between machine translation quality and post-editing efforts and found no straightfoward relationship. On the other hand, Popovic et al. (2016) suggested that post-edits should be used carefully for MT evaluation due to the bias of each postedit towards its MT system. Further, Toral (2019) showed that human post-edits are simpler and more normalised in language than human translations from scratch.

**Automatic Refinement**    Lin et al. (2022) showed how the errors that humans make differ from those made by MT systems. They constructed a Translation Error Correction (TEC) corpus with professional translators and showed that models trained on it outperform Automatic Post-Editing (APE) models (Knight and Chander, 1994) that are trained to correct MT output. Since the emergence of LLMs, new refinement approaches based on detailed MQM annotations have appeared (Xu et al., 2023; Fernandes et al., 2023). Xu et al. (2024) showed that these refinement method can be used

to improve the quality of human translations.

Meanwhile, machines have been extensively evaluated and utilized as an alternative to human annotators for data collection (Zouhar et al., 2021; Yan et al., 2024).

In contrast to the above methods, we investigate the interaction between humans and machines in the initial translation and post-editing stages, including detailed analysis of the resulting changes in quality while also considering cost-efficiency.

# 7    Conclusion

We investigate various approaches for gathering translation data, including human-only, machineonly, and hybrid approaches. Our results demonstrate that human-machine collaboration can consistently generate high-quality translations at a lower cost than human-only methods. Through detailed error analysis, we uncovered the nuances of error correction dynamics and highlighted the advantages of human-machine collaborative methods. Our cost analysis also demonstrates the costefficiency of human-machine collaboration methods. Finally, we release to the public a dataset of roughly 18,000 translation segments of varying quality from different collection methods along with human ratings, to facilitate further research in this area.

## Limitations

This study focuses on two language pairs, English-German and Chinese-English. They are chosen due to the extensive study in the WMT23 metrics shared task (Freitag et al., 2023) and the availability of data from various translation systems from the WMT23 general shared task (Kocmi et al., 2023). While our analysis provides support for the findings presented in this work and we offer a plausible explanation for the observed results, it is important to acknowledge certain variables are not accounted for in this work, including using translators or post-editors with varying quality levels, different systems for translation and post-editing, utilizing sentence or paragraph datasets from other domains, and higher or lower resource language pairs beyond the two investigated here. Therefore, we cannot guarantee the observed trends will generalize to different datasets.

We want to especially highlight the need for further exploration of the quality variance observed among human translators, such as ORIGHUMAN-REF and HUMANREF in the English-German translation task. The current study's limited annotation budget and timeline restricted the depth of this investigation. Future experiments aimed at examining the impact of post-editing on annotator agreement would be particularly interesting and valuable.

## Ethical Statement

The source data used for translation and post-edits is accessible to the public. We're certain that the data annotated by human labors is free from risk or toxic content. We used an internal, proprietary tool to collect human translation, post-edits, and evaluation data. The annotators were compensated fairly and were not required to disclose any personal details during the annotation process. All the test data used in this study are publicly available and annotators were allowed to label sensitive information if necessary. The annotators are fully informed that the data they collected will be used for research purposes.

## References

Vicent Alabau, Christian Buck, Michael Carl, Francisco Casacuberta, Mercedes García-Martínez, Ulrich Germann, Jesús González-Rubio, Robin L Hill, Philipp Koehn, Luis A Leiva, et al. 2014. Casmacat: A computer-assisted translation workbench. In *Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 25–28.

Michael Carl, Barbara Dragsted, Jakob Elming, Daniel Hardt, and Arnt Lykke Jakobsen. 2011. The process of post-editing: A pilot study. *Copenhagen Studies in Language*, 41(1):131–142.

Michael Denkowski. 2015. Machine translation for human translators.

Marcello Federico, Nicola Bertoldi, Mauro Cettolo, Matteo Negri, Marco Turchi, Marco Trombetti, Alessandro Cattelan, Antonio Farina, Domenico Lupinetti, Andrea Martines, et al. 2014. The matecat tool. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: System Demonstrations*, pages 129–132.

Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André Martins, Graham Neubig, Ankush Garg, Jonathan Clark, Markus Freitag, and Orhan Firat. 2023. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1066–1083, Singapore. Association for Computational Linguistics.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Team Gemini. 2024. Gemini: A family of highly capable multimodal models.

Spence Green, Jeffrey Heer, and Christopher D. Manning. 2013. The efficacy of human post-editing for language translation. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '13, page 439–448, New York, NY, USA. Association for Computing Machinery.

Spence Green, Sida I. Wang, Jason Chuang, Jeffrey Heer, Sebastian Schuster, and Christopher D. Manning. 2014. Human effort and machine learnability in computer aided translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1225–1236, Doha, Qatar. Association for Computational Linguistics.

Dorothy Kenny. 2012. *Electronic Tools and Resources for Translators*.

Kevin Knight and Ishwar Chander. 1994. Automated postediting of documents. In *AAAI*, volume 94, pages 779–784.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Jessy Lin, Geza Kovacs, Aditya Shastry, Joern Wuebker, and John DeNero. 2022. Automatic correction of human translations. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 494–507, Seattle, United States. Association for Computational Linguistics.

Arle Lommel, Hans Uszkoreit, and Aljoscha Burchardt. 2014. Mqm: Un marc per declarar i descriure mètriques de qualitat de la traducció. *Tradumàtica: traducció i tecnologies de la informació i la comunicació*, (12):455–463.

Team OpenAI. 2024. Gpt-4 technical report.

Mirko Plitt and François Masselot. 2010. A productivity test of statistical machine translation post-editing in a typical localisation context. In *Prague Bulletin of Mathematical Linguistics*.

Maja Popovic, Mihael Arčan, and Arle Lommel. 2016. Potential and limits of using post-edits as reference translations for MT evaluation. In *Proceedings of the 19th Annual Conference of the European Association for Machine Translation*, pages 218–229.

Chan Sin-wai. 2014. The development of translation technology 1967–2013. In *Routledge Encyclopedia of Translation Technology*, pages 2–31. Routledge.

Antonio Toral. 2019. Post-editese: an exacerbated translationese. In *Proceedings of Machine Translation Summit XVII: Research Track*, pages 273–281, Dublin, Ireland. European Association for Machine Translation.

Antonio Toral, Martijn Wieling, and Andy Way. 2018. Post-editing effort of a novel with statistical and neural machine translation. *Frontiers in Digital Humanities*, 5:9.

Yangjian Wu and Gang Hu. 2023. Exploring prompt engineering with GPT language models for document-level machine translation: Insights and findings. In *Proceedings of the Eighth Conference on Machine Translation*, pages 166–169, Singapore. Association for Computational Linguistics.

Wenda Xu, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, Biao Zhang, Zhongtao Liu, William Yang Wang, Lei Li, and Markus Freitag. 2024. Llmrefine: Pinpointing and refining large language models via fine-grained actionable feedback.

Wenda Xu, Danqing Wang, Liangming Pan, Zhenqiao Song, Markus Freitag, William Wang, and Lei Li. 2023. INSTRUCTSCORE: Towards explainable text generation evaluation with automatic feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5967–5994, Singapore. Association for Computational Linguistics.

Jianhao Yan, Pingchuan Yan, Yulong Chen, Judy Li, Xianchao Zhu, and Yue Zhang. 2024. Gpt-4 vs. human translators: A comprehensive evaluation of translation quality across languages, domains, and expertise levels.

Vilém Zouhar, Martin Popel, Ondřej Bojar, and Aleš Tamchyna. 2021. Neural machine translation quality and post-editing performance. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10204–10214, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Vilém Zouhar and Ondřej Bojar. 2024. Quality and quantity of machine translation references for automatic metrics.

Vilém Zouhar, Věra Kloudová, Martin Popel, and Ondřej Bojar. 2024. Evaluating optimal reference translations.

## A Appendix

### A.1 Cross-BLUE scores

Figure 8 presents the cross-BLEU similarity matrix for Chinese-English translation systems.

### A.2 Error Distribution of Initial Translation

Table 7 presents the error type and severity distributions of Chinese-English MQM human evaluation results.

Figure 8: Cross-BLEU scores for different Chinese-English translation collection approaches

| Error Type | OrigHumanRef | HumanRef | BestWMT | MidWMT |
|---|---|---|---|---|
| **No-error** | 225 | 453 | 490 | 304 |
| **Major** | | | | |
| Total | 955 | 284 | 260 | 648 |
| Fluency | 18 (2%) | 15 (5%) | 5 (2%) | 23 (4%) |
| Accuracy | 851 (89%) | 229 (81%) | 221 (85%) | 568 (88%) |
| Style | 31 (3%) | 12 (4%) | 5 (2%) | 18 (3%) |
| **Minor** | | | | |
| Total | 1902 | 1303 | 1238 | 1704 |
| Fluency | 625 (33%) | 431 (33%) | 364 (29%) | 522 (31%) |
| Accuracy | 599 (31%) | 412 (32%) | 388 (31%) | 447 (26%) |
| Style | 629 (33%) | 402 (31%) | 436 (35%) | 677 (40%) |

Table 7: Error type and severity distributions of Chinese-English MQM human evaluation results.

## A.3 Error Correction from Post-Editing

Figures 9, 10, and 11 present Chinese-English results comparable to the English-German results presented in Section 4.2. Table 8 presents the same data as in Figure 4 for English-German, and Table 9 presents the same for Chinese-English.



Figure 9: Error changes percentages by different post-editing approaches on Chinese-English data. The percentages present the changes in error counts for each post-editing method compared to its initial translation. A negative indicates a decrease in errors, while positive value indicates an increase in the error type.



(a) From HUMANREF to HUMANREF+HUMANPE



(b) From HUMANREF to HUMANREF+LLMREFINE

Figure 10: Segment-level quality shift through HUMANPE and LLMREFINE from Chinese-English HUMANREF.

| Initial Translation | Total Seg | HumanPE | LLMRefine | HumanPE & LLMRefine | Human Only | LLMRefine Only |
|---|---|---|---|---|---|---|
| HumanRef | 460 | 163 (35.4%) | 235 (51.1%) | 110 (23.9%) | 53 (11.5%) | 125 (27.2%) |
| BestWMT | 460 | 314(68.3%) | 281 (61.1%) | 225 (48.9%) | 89 (19.3%) | 56 (12.2%) |
| MidWMT | 460 | 400 (87.0%) | 324 (70.4%) | 310 (67.4%) | 90 (19.6%) | 14 (3%) |

Table 8: Numerical breakdown of the agreement between HUMANPE and LLMREFINE in identifying segments requiring post-editing in English-German

| Initial Translation | Total Seg | HumanPE | LLMRefine | HumanPE & LLMRefine | Human Only | LLMRefine Only |
|---|---|---|---|---|---|---|
| HumanRef | 1175 | 558 (47.5%) | 225 (19.1%) | 161 (13.7%) | 397 (33.8%) | 64 (5.4%) |
| BestWMT | 1175 | 830 (70.6%) | 133 (11.3%) | 123 (10.5%) | 707 (60.2%) | 10 (0.9%) |
| MidWMT | 1175 | 1006 (85.6%) | 408 (34.7%) | 399 (34.0%) | 607 (51.7%) | 9 (0.8%) |

Table 9: Numerical breakdown of the agreement between HumanPE and LLMRefine in identifying segments requiring post-editing in Chinese-English



(a) From BESTWMT to BESTWMT+HUMANPE



(b) From BESTWMT to BESTWMT+LLMREFINE

Figure 11: Segment-level quality shift through HU-MANPE and LLMREFINE from Chinese-English BESTWMT.

# How Effective are State Space Models for Machine Translation?

**Hugo Pitorro**[*,1,3], **Pavlo Vasylenko**[*,2,3], **Marcos Treviso**[3], **André F. T. Martins**[2,3,4,5]

[1]TU Munich, [2]Instituto Superior Técnico, Universidade de Lisboa
[3]Instituto de Telecomunicações, [4]Unbabel, [5]ELLIS Unit Lisbon

## Abstract

Transformers are the current architecture of choice for NLP, but their attention layers do not scale well to long contexts. Recent works propose to replace attention with linear recurrent layers—this is the case for state space models, which enjoy efficient training and inference. However, it remains unclear whether these models are competitive with transformers in machine translation (MT). In this paper, we provide a rigorous and comprehensive experimental comparison between transformers and linear recurrent models for MT. Concretely, we experiment with RetNet, Mamba, and hybrid versions of Mamba which incorporate attention mechanisms. Our findings demonstrate that Mamba is highly competitive with transformers on sentence and paragraph-level datasets, where in the latter both models benefit from shifting the training distribution towards longer sequences. Further analysis show that integrating attention into Mamba improves translation quality, robustness to sequence length extrapolation, and the ability to recall named entities.

## 1 Introduction

The inherent design of attention—the underlying mechanism of transformers—leads to quadratic computational costs and challenges in length generalization (Varis and Bojar, 2021). As an alternative, recent works propose to replace attention with linear recurrent approaches, which enjoy efficient training and inference, and obtain competitive results in language modeling tasks (Katharopoulos et al., 2020; Gu et al., 2022; Peng et al., 2023; Sun et al., 2023a; Gu and Dao, 2023).

In machine translation (MT), there is an increasing demand for supporting longer context lengths, such as paragraphs or entire documents (Fernandes et al., 2021; Wang et al., 2023; Kocmi et al., 2023). Given this trend, it has become increasingly important to design models capable of efficiently handling longer sequences. Previous research indicates that models like state space models (SSMs), exemplified by S4 (Gu et al., 2022), still lag behind transformers in MT (Vardasbi et al., 2023). However, it remains unclear whether these findings hold true for recent, more expressive variations of linear recurrent models, such as RetNet (Sun et al., 2023a) and Mamba (Gu and Dao, 2023), especially on settings that involve the use of pretrained models and long context datasets.

In this paper, we provide a rigorous and comprehensive experimental comparison between transformers, RetNet, Mamba, as well as hybrid versions of Mamba that incorporate attention mechanisms (§4). We also compare with pretrained Mamba and Pythia (Biderman et al., 2023) at two parameter scales, ∼400M and 1.4B. Building on existing literature that explores the capabilities of linear recurrent models in language modeling (Arora et al., 2024a; Jelassi et al., 2024), we further investigate the performance of models trained from scratch in recalling context tokens during the translation process (§4.2). Moreover, we extend our analysis by investigating the models' ability to handle long contexts, on paragraph-level datasets (§5), along with measuring their sensitivity to different sequence lengths (§5.2) and inference cost (§5.4). Overall, our main findings are:[1]

- For sentence-level experiments, we show that Mamba exhibits competitive performance compared to transformers, for both trained-from-scratch and pretrained models.

- At the paragraph level, we find that Mamba is sensitive to the training distribution's sequence length and struggles with longer inputs. However, shifting the distribution towards longer sequence lengths helps to close the gap with transformers.

- We observe that integrating attention and state

---

[*]Equal contribution.

[1]https://github.com/deep-spin/ssm-mt

space models creates a strong model in terms of translation quality, robustness to sequence length extrapolation, and ability to recall named entities.

## 2 Background

In this section, we present an overview of transformers, and the foundation of the linear recurrent models covered in this paper: linear attention (RetNet) and state space models (Mamba).

### 2.1 Transformers

The key component in the transformer architecture is the attention mechanism, which is responsible for contextualizing information within and across input sequences. Concretely, given query $Q \in \mathbb{R}^{n \times d}$, key $K \in \mathbb{R}^{n \times d}$, and value $V \in \mathbb{R}^{n \times d}$ matrices as input, where $n$ is the sequence length and $d$ the hidden size, the single head *self-attention mechanism* is defined as follows (Vaswani et al., 2017):

$$Y = \mathsf{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) V \in \mathbb{R}^{n \times d}. \quad (1)$$

For decoder-only models, a causal mask is used to ignore future tokens. Notably, the $QK^\top$ operation leads to a $\mathcal{O}\left(n^2\right)$ cost during training, and $\mathcal{O}\left(n\right)$ during inference with caching and causal masking.

### 2.2 Linear Attention

Denote by $q_i, k_i, v_i, y_i \in \mathbb{R}^d$ respectively the (column) vectors corresponding to the $i^{\text{th}}$ rows of the matrices $Q, K, V, Y$ defined above. Katharopoulos et al. (2020) reformulate the attention mechanism by casting the role of the softmax as a similarity function $\mathsf{sim}\left(q, k\right) = \exp\left(q^\top k / \sqrt{d}\right)$:

$$y_i = \frac{\sum_{j=1}^n \mathsf{sim}(q_i, k_j) v_j}{\sum_{j=1}^n \mathsf{sim}(q_i, k_j)}. \quad (2)$$

However, any kernel $k(x, y) : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ is a suitable candidate for the similarity function (Smola and Schölkopf, 1998; Tsai et al., 2019). In particular, a kernel $k(x, y) = \phi(x)^\top \phi(y)$, where $\phi : \mathbb{R}^d \to \mathbb{R}^r$ is a feature map, leads to:

$$
\begin{aligned}
y_i &= \frac{\sum_{j=1}^n \phi(q_i)^\top \phi(k_j) v_j}{\sum_{j=1}^n \phi(q_i)^\top \phi(k_j)} \\
&= \frac{\sum_{j=1}^n v_j \phi(k_j)^\top \phi(q_i)}{\sum_{j=1}^n \phi(k_j)^\top \phi(q_i)} \\
&= \frac{S^\top \phi(q_i)}{z^\top \phi(q_i)}, \quad (3)
\end{aligned}
$$

where $S = \sum_{j=1}^n \phi(k_j) v_j^\top \in \mathbb{R}^{r \times d}$ and $z = \sum_{j=1}^n \phi(k_j) \in \mathbb{R}^r$. Notably, if initial states are initialized as $S_0 = \mathbf{0}_{r \times d}$ and $z_0 = \mathbf{0}_r$, intermediate states can be computed in a recurrent fashion:

$$
\begin{aligned}
S_i &= S_{i-1} + \phi(k_i) v_i^\top, \\
z_i &= z_{i-1} + \phi(k_i). \quad (4)
\end{aligned}
$$

Since we can reuse the same $S_i$ and $z_i$ for all queries, this recurrent variant offers a $\mathcal{O}\left(n\right)$ complexity during training and enjoys a $\mathcal{O}\left(1\right)$ complexity for inference.[2]

**Retentive Networks (RetNet).** Sun et al. (2023a) set $\phi$ as the identity function, i.e., $k(q, k) = q^\top k$, ignore the normalizer in Equation 2, and introduce an exponential decay mask $\gamma$, leading to:

$$
\begin{aligned}
S_i &= \gamma S_{i-1} + k_i v_i^\top, \\
y_i &= S_i^\top q_i. \quad (5)
\end{aligned}
$$

This formulation effectively biases the attention mechanism to focus on closer token interactions. RetNet also uses XPos (Sun et al., 2023b), a relative positional encoding method, to improve its context extrapolation abilities.

### 2.3 State Space Models (SSMs)

SSMs (Gu et al., 2020) provide an alternative sequence mixing layer by processing sequences $x_1, ..., x_n$, where each $x_i \in \mathbb{R}^d$, through a linear recurrence. Letting $H_i \in \mathbb{R}^{r \times d}$ denote the "state" at the $i^{\text{th}}$ time step, a discrete SSM is defined as follows:[3]

$$
\begin{aligned}
H_i &= A H_{i-1} + b x_i^\top, \\
y_i &= H_i^\top c, \quad (6)
\end{aligned}
$$

where $A \in \mathbb{R}^{r \times r}$, $b \in \mathbb{R}^r$, and $c \in \mathbb{R}^r$ are (discrete) parameters.[4] Since the same parameters are used for both relevant and irrelevant inputs, this model is deemed *input-independent*, which, in turn,

---

[2]In practice, however, this recurrent view is not parallelizable, leading to chunkwise-recurrent variations for training (Hua et al., 2022; Sun et al., 2023a; Yang et al., 2024).

[3]A discretization step is needed in order to obtain discrete parameters. For example, a possible method for this step is the zero-order hold rule, used by Mamba (Gu and Dao, 2023).

[4]The SSM equations are commonly written independently for each input dimension $j \in [d]$ as

$$h_i^{(j)} = A h_{i-1}^{(j)} + b x_i^{(j)}, \quad y_i^{(j)} = c^\top h_i^{(j)},$$

with $A$, $b$, and $c$ shared across input dimensions. This is equivalent to (6), where the $j^{\text{th}}$-column of $H_i$ equals $h_i^{(j)}$.

makes the model unable to reset or overwrite its hidden states. S4 (Gu et al., 2022) is an instance of this model, which enjoys a $\mathcal{O}(n \log n)$ time complexity during training, and $\mathcal{O}(1)$ during inference. Vardasbi et al. (2023) shows that S4 still underperforms transformers for MT. Finally, note the similarity between Eq. 5 and Eq. 6: RetNets can be seen as state space models with $\boldsymbol{A} = \gamma \boldsymbol{I}$ and data-dependent $\boldsymbol{b}$ and $\boldsymbol{c}$.

**Mamba.** To make the SSM parameters *data-dependent*, Mamba (Gu and Dao, 2023) introduces a selection mechanism that uses learnable linear projections over $\boldsymbol{x}$ prior to the discretization step, effectively making all parameters dependent on the $i^{\text{th}}$ input. This leads to:

$$\boldsymbol{H}_i = \boldsymbol{A}_i \odot \boldsymbol{H}_{i-1} + \boldsymbol{B}_i \odot \boldsymbol{X}_i,$$
$$\boldsymbol{y}_i = \boldsymbol{H}_i^{\top} \boldsymbol{c}_i, \tag{7}$$

where $\boldsymbol{X}_i = \boldsymbol{1}_r \boldsymbol{x}_i^{\top} \in \mathbb{R}^{r \times d}$ is an $r$-sized stack of the input, $\boldsymbol{A}_i \in \mathbb{R}^{r \times d}$ represents $d$ diagonal matrices of size $r \times r$, $\boldsymbol{B}_i \in \mathbb{R}^{r \times d}$, $\boldsymbol{c}_i \in \mathbb{R}^r$, and $\odot$ is the Hadamard product. Note that, unlike S4, where the same $\boldsymbol{A}$ and $\boldsymbol{B}$ parameters are shared across all hidden dimensions $1 \leq h \leq d$, Mamba defines $\boldsymbol{A}_i$ and $\boldsymbol{B}_i$ with a shape of $(\ldots, d)$, allowing for unique parameters in each hidden dimension. While this formulation makes Mamba more expressive, it disrupts the convolutional approach used for training in S4. To address this, Gu and Dao (2023) propose an efficient IO-aware and parallelizable associative scan algorithm for training (Smith et al., 2023). Nonetheless, the recurrent view can still be used for inference with a $\mathcal{O}(1)$ time complexity.

# 3 Experimental Setup

We conduct experiments with transformers, RetNet, and Mamba for MT in §4 and §5. In this section, we detail the sentence and paragraph-level datasets used in our experiments, along with the settings for our models, which are trained in two distinct regimes: from scratch, or finetuned from a pretrained checkpoint.

## 3.1 Datasets

For sentence-level experiments, we focus on WMT14 DE↔EN and WMT16 RO↔EN for consistency with previous works (Vardasbi et al., 2023), but also include WMT16 FI↔EN using the standard training, validation and test splits. For paragraph level, we use the more recent WMT23

| DATASET | # SAMPLES | # TOKENS |
|---|---|---|
| IWSLT17 (DE↔EN) | 200K | 45.2 ± 29.5 |
| WMT16 (RO↔EN) | 610K | 58.9 ± 31.1 |
| WMT16 (FI↔EN) | 2.08M | 52.8 ± 33.1 |
| WMT14 (DE↔EN) | 4.5M | 62.1 ± 45.6 |
| WMT23-6M (DE↔EN) | 6M | 58.4 ± 32.9 |
| WMT23-CAT-5 (DE↔EN) | 2M | 171.3 ± 134.9 |
| WMT23-CAT-10 (DE↔EN) | 1M | 312.4 ± 282.3 |
| WMT23 Test (DE→EN) | 549 | 135.1 ± 147.7 |
| WMT23 Test (EN→DE) | 557 | 185.2 ± 188.2 |
| Ted Talks Val. (DE↔EN) | 995 | 268.5 ± 189.6 |
| Ted Talks Test (DE↔EN) | 2247 | 939.2 ± 594.0 |

Table 1: Sentence and paragraph-level datasets statistics.

dataset (Kocmi et al., 2023), which contains ∼300M training samples and ∼1K test samples incorporating multi-sentence passages. In order to obtain a small high-quality subset for training, we exclude ParaCrawl and CommonCrawl samples from the original dataset and clean the remaining data. Our cleaning process includes three steps. First, we identify and remove samples in incorrect languages via `langdetect`[5]. Second, we eliminate duplicates. Third, we rank the samples using COMETKIWI-22 (Rei et al., 2022b) a state-of-the-art translation quality estimator, and keep only the top 6M samples. We call the refined dataset WMT23-6M. Datasets statistics are shown in Table 1.

## 3.2 Models

We make a broad selection of models spanning both trained-from-scratch and finetuned versions. In the first setting, we compare standard transformers, linear recurrent models, and also hybrid approaches that integrate attention into Mamba. For finetuned models, we experiment with released Pythia and Mamba checkpoints. We describe each model next.

### 3.2.1 Standard Models

**Transformers.** We select two variants of the transformer model as baselines: a base encoder-decoder formulation and a modern decoder-only version. The **Transformer Enc-Dec.** model, as described in the original paper (Vaswani et al., 2017), has 77M parameters, and uses sinusoidal positional embeddings and standard ReLU activations. The second variant, **Transformer++**, is a decoder-only formulation incorporating recent advancements, such as rotary positional embeddings (Su et al., 2024) and the SwiGLU layer (Shazeer,

---

[5] https://github.com/Mimino666/langdetect

2020). Specifically, we use the LLaMA architecture (Touvron et al., 2023), adjusting the embedding dimension to match the parameter count of the base transformer (79M), consistent with the version employed in (Gu and Dao, 2023).

**Linear recurrent models.** We select two representative recurrent models, **RetNet** (Sun et al., 2023a) and **Mamba** (Gu and Dao, 2023). Both models are tested with 77M parameters to approximately match the number of parameters in the transformer models.

### 3.2.2 Hybrid Models

Previous work has shown that incorporating attention into linear recurrent models leads to strong performance in language modeling (Fu et al., 2023; Arora et al., 2024b; De et al., 2024). Therefore, we aim to examine if this is also the case for MT by exploring three hybrid variants, detailed next.

**Mamba-MHA.** The simplest hybrid formulation involves replacing some of the Mamba layers with attention. Some natural questions then arise: how many attention layers are needed, and where to place them? After careful ablations, detailed in Appendix B, we use two attention layers placed at the middle and at the output of the network, resembling the hybrid version of H3 (Fu et al., 2023).

**Mamba-Local.** While aiming to achieve robust performance, the introduction of full attention to Mamba disrupts its efficiency gains. Thus, we consider local attention variants such as sliding window attention (Beltagy et al., 2020; Child et al., 2019), employed in recent hybrid models (Arora et al., 2024b; De et al., 2024). We use a window size of 64 based on the average sequence length shown in Table 1 and ablations in Appendix B.

**Mamba Enc-Dec.** Lastly, inspired by the S4-based encoder-decoder model from Vardasbi et al. (2023), we replace the self-attention mechanism in transformers with a Mamba block and keep the cross-attention component intact. In terms of complexity, since this variant computes attention over the source sentence, it incurs an additional $\mathcal{O}\left(n^2\right)$ cost for training and $\mathcal{O}\left(n\right)$ for inference.

### 3.2.3 Pretrained Models

In order to fairly evaluate the relative performance between pretrained models, we need to ensure consistency between their pretraining data. Taking this into account, we consider two strong models pretrained on The Pile (Gao et al., 2020): Pythia (Biderman et al., 2023), a modern transformer, and Mamba, a modern SSM. Note, however, that Pythia was pretrained on more tokens than Mamba (see Table 6), hence the comparison might be slightly unfavorable to Mamba. We experiment with two model scales, *small* (S) and *medium* (M). Concretely, we experiment with Pythia 410M and 1.4B, and with Mamba 370M and 1.4B.

### 3.3 Training and Evaluation

For models trained from scratch, we follow the settings proposed in (Vardasbi et al., 2023), whereas for pretrained models, we follow the finetuning settings used by Mamba (Gu and Dao, 2023). For decoder-only models, we pass a concatenation of the source and target sequences separated by a special token as input. We evaluate all models with BLEU (Post, 2018)[6] and COMET (Rei et al., 2022a).[7] We base our analysis on the latter, given its strong correlation with human judgments on sentence and paragraph-level data (Freitag et al., 2022, 2023). More training details can be found in §A.

## 4 Sentence-level Translation

We start by evaluating our standard, hybrid, and finetuned models on the sentence-level WMT16 RO↔EN, FI↔EN and WMT14 DE↔EN datasets. Results can be found in Table 2 in terms of BLEU and COMET. Next, we discuss the key findings.

### 4.1 Discussion

**Mamba is competitive when trained from scratch.** Mamba, a decoder-only model, not only outperforms a decoder-only transformer (Transformer++) across the board, but also an encoder-decoder transformer (Transf. Enc-Dec) in the larger WMT14 for both DE↔EN language pairs. This creates a contrast with the S4 results obtained by Vardasbi et al. (2023). We hypothesize that Mamba's good results are due to its data-dependent state updates (Eq. 7), which allows for more precise information retention in its hidden state. On the other hand, RetNet's performance is generally subpar compared to other models, likely due to its strong locality bias (induced by $\gamma$ in Eq. 5), which may hinder performance in MT, a task where the source

---

[6]SacreBLEU signature: `|1|mixed|no|13a|exp|`
[7]huggingface.co/Unbabel/wmt22-comet-da

| | | WMT16 | | | | | | | | WMT14 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | RO→EN | | EN→RO | | FI→EN | | EN→FI | | DE→EN | | EN→DE | |
| MODEL | SIZE | BLEU | COMET | BLEU | COMET | BLEU | COMET | BLEU | COMET | BLEU | COMET | BLEU | COMET |
| *Trained from scratch* | | | | | | | | | | | | | |
| Transf. Enc-Dec | 77M | **29.2** | <u>74.8</u> | <u>22.0</u> | **78.6** | 15.3 | 70.5 | **14.8** | **78.2** | 27.4 | 78.6 | 22.3 | 77.1 |
| Transformer++ | 79M | 26.4 | 72.6 | 21.7 | 72.7 | 14.9 | 69.3 | 14.2 | 75.5 | 26.9 | 79.0 | 22.8 | 77.9 |
| RetNet | 77M | 26.4 | 72.4 | 19.9 | 76.0 | 14.5 | 70.2 | 11.0 | 70.2 | 23.4 | 74.7 | 19.6 | 71.7 |
| Mamba | 77M | 27.0 | 73.8 | 21.4 | 77.9 | 16.0 | 72.7 | 13.0 | 75.4 | **27.5** | 80.2 | 22.4 | 77.8 |
| Mamba-MHA | 78M | <u>28.5</u> | **75.1** | 21.7 | 78.3 | **17.5** | **73.8** | <u>14.3</u> | 76.4 | <u>27.4</u> | **80.6** | **23.2** | **79.9** |
| Mamba-Local | 78M | 25.9 | 73.9 | 20.9 | 76.9 | 16.3 | 73.1 | 13.2 | 75.4 | 27.2 | <u>80.1</u> | <u>23.2</u> | <u>79.5</u> |
| Mamba Enc-Dec | 82M | 28.5 | 74.4 | **22.7** | 77.9 | <u>17.0</u> | <u>73.6</u> | <u>14.3</u> | <u>77.0</u> | 27.2 | 80.0 | 21.6 | 78.8 |
| *Finetuned* | | | | | | | | | | | | | |
| Pythia-S | 410M | 33.4 | 82.0 | 24.1 | 85.8 | 19.8 | 80.1 | 16.5 | 84.6 | 30.9 | 83.6 | 25.2 | 84.0 |
| Mamba-S | 370M | **34.1** | **83.2** | 24.2 | <u>86.4</u> | **21.4** | 81.5 | 16.5 | 85.5 | 29.8 | 83.3 | 25.0 | 83.2 |
| Pythia-M | 1.4B | <u>33.9</u> | **83.2** | **24.9** | **87.1** | 20.9 | <u>81.7</u> | **17.8** | **87.1** | **32.2** | **84.5** | **26.7** | **84.9** |
| Mamba-M | 1.4B | 33.8 | 83.1 | <u>24.5</u> | 86.2 | <u>21.3</u> | **82.1** | 18.4 | <u>86.8</u> | <u>31.9</u> | **84.5** | 26.5 | <u>84.2</u> |

Table 2: Sentence-level results in terms of BLEU and COMET for models trained from scratch (top) and models finetuned from a pretrained checkpoint (bottom). **Bold** represents top results; <u>underline</u> represents second-best.

input servers as a prefix to the translation, and it requires "focused attention" during decoding.

**Attention benefits Mamba.** By including attention layers in Mamba's architecture, we find that Mamba-MHA, which employs only two attention layers, is able to outperform both transformers and Mamba for almost all language pairs. While Mamba-Local retains constant inference complexity via windowed attention, it is not as strong as the full attention variant. Finally, Mamba Enc-Dec's performance is also competitive, falling just short of Mamba-MHA and echoing the S4 encoder-decoder findings of Vardasbi et al. (2023).

**Pretraining improves all models.** We note a large COMET gap, roughly 4-8 COMET points, between the finetuned models and those trained from scratch for all language pairs. This is expected, since not only are these models bigger, but they also have strong data-driven priors, which are beneficial in downstream tasks (Amos et al., 2024).

**Larger models achieve top results.** For small models, Mamba outperforms Pythia for RO↔EN and FI↔EN in terms of COMET and BLEU. However, Pythia is superior on the larger DE↔EN datasets. Moving to larger models, we note that Mamba improves COMET scores by ∼1 point on EN↔DE and 0.6-1.3 point on EN↔FI while dropping only 0.1-0.2 on EN↔RO datasets. On the other hand, Pythia improves results consistently for all language pairs with a larger model size, outperforming or matching the results of other models. On average, we find that both their gaps decrease

when moving from smaller to medium-sized models but Pythia benefits more in terms of COMET. It is worth noting that Mamba is pretrained on fewer samples than Pythia (see Table 6) and that the impact of pretraining data quality can also play a role in downstream task performance.

## 4.2 Recall of Named Entities

Following our discussion of sentence-level translation, we now focus on how well these models recall context tokens during translation. Inspired by prior studies investigating the recall of context tokens in language modeling with state space models (Arora et al., 2024a; Jelassi et al., 2024), we conduct a similar experiment for MT. Unlike language modeling, where token patterns often recur within a near context, MT presents a challenge due to the distinct spelling of words across languages. Therefore, we focus on the recall of named entities (NEs) that appear verbatim in both source and target sentences, using NLTK for NE recognition (Bird, 2006).

We assess the models' ability to recall NEs from the WMT16 RO→EN dataset according to their frequency in the training set, as illustrated in Figure 1. The results reveal a clear correlation between NE frequency and their chance to be recalled in the translation process, as more frequent NEs are recalled more often. Notably, we note a disparity in performance with unseen entities, which provides a more illustrative view of recall ability. In this respect, transformers and Mamba perform on par, while RetNet shows inferior results. As before, the hybrid models are promising, with Mamba-

Figure 1: Recall in recovering named entities on the WMT16 RO→EN dataset by their training set frequency: *unseen* entities do not appear in the training data, while *regular* and *frequent* entities appear $[1, 16)$ and $16+$ times.

MHA outperforming all models, followed closely by Mamba Enc-Dec. We include additional analyses for other datasets in the Appendix §C.

## 5 Paragraph-level translation

While Mamba shows competitive performance with transformers on sentence-level datasets (see Table 2), it was originally designed to handle long sequences. Thus, we now turn our attention to paragraph-level datasets. This allows us to study the models' sensitivity to long sequence lengths along with their robustness in handling sequences that are longer than the ones seen during training.[8]

To this end we focus on the WMT23-6M dataset (§3.1), from which the training and test sets are composed of sentence and paragraph-level passages, respectively. In order to see the impact of training on long sequences, we propose to artificially construct multi-sentence datasets, which we call WMT23-CAT-$N$. Our procedure is as follows:

1. We first retain the original training samples from WMT23-6M with a probability of $50\%$.

2. Next, for the remaining part, we concatenate $N$ random training samples.

This approach ensures that we consistently maintain a $50\%$ ratio between single-sentence and multi-sentence samples. For validation, we sample 1-to-10-sentence passages from the TED Talks dataset (Cettolo et al., 2012). Statistics for CAT-$N$ datasets can be found in Table 1. COMET scores on the WMT23 EN↔DE test sets are shown in Table 3. We provide additional BLEU scores in Table 9 in Appendix E. Next, we discuss our key findings.

### 5.1 Discussion

**Concatenation helps.** Our strategy of concatenating sentences proves beneficial for almost all

models, as COMET scores typically improve with the CAT-5 and CAT-10 datasets, whether models are trained from scratch or finetuned. Among models trained from scratch, Transformer Enc-Dec, Mamba-MHA, and Mamba Enc-Dec show substantial improvements, with Mamba Enc-Dec achieving the best overall results. For finetuned models, concatenation benefits larger models; Mamba-M outperforms Pythia-M in DE→EN but underperforms in EN→DE. Interestingly, for both training regimes, the concatenation strategy can lead to COMET gains of up to 5 points.

**Finetuning outperforms training from scratch.** Finetuned models consistently achieve higher COMET scores, with larger models attaining the top results overall. Similar to sentence-level experiments, Pythia outperforms Mamba when trained on the original, WMT23-6M dataset. However, both Pythia and Mamba benefit from our concatenation strategy. While these results indicate that our concatenation strategy helps in translating long inputs, it remains unclear whether performance on short inputs is compromised or if the models can handle longer inputs than those seen during training. We investigate these issues next.

### 5.2 Sensitivity to Input Length

Based on the previous observations, we notice that performance between models varies considerably after being exposed to different sequence lengths during training. In this subsection, we investigate how robust each model is to length distribution shifts between training and test. Results are shown in Figure 2 for both training regimes on the WMT23 DE→EN dataset. Results are consistent for EN→DE, shown in Figure 6, Appendix D.

**Discussion.** When training on WMT23-6M, we observe a decline in performance for all models on long sequences, affecting both trained-from-scratch

---

[8]We dropped RetNet and Mamba-Local as they already achieve poor results on long *sentence-level* inputs (see Fig. 5).

| MODEL | SIZE | DE→EN | | | EN→DE | | |
|---|---|---|---|---|---|---|---|
| | | ORIG. | CAT5 | CAT10 | ORIG. | CAT5 | CAT10 |
| *Trained from scratch* | | | | | | | |
| Transf. Enc-Dec | 77M | 72.4 | 74.6 | 69.6 | 65.2 | <u>70.3</u> | <u>70.3</u> |
| Transformer++ | 79M | 70.7 | 73.6 | 72.8 | 64.8 | 69.1 | 68.8 |
| Mamba | 77M | 70.0 | 73.3 | 72.3 | 63.3 | 67.5 | 67.8 |
| Mamba-MHA | 78M | 72.7 | 74.2 | <u>74.5</u> | 67.0 | 68.6 | 69.7 |
| Mamba Enc-Dec | 82M | 70.7 | 73.8 | **75.6** | 65.3 | **71.0** | 70.1 |
| *Finetuned* | | | | | | | |
| Pythia-S | 410M | 77.4 | 78.4 | 79.0 | 76.7 | <u>77.8</u> | 77.1 |
| Mamba-S | 370M | 77.2 | 78.2 | 78.3 | 72.4 | 74.2 | 73.1 |
| Pythia-M | 1.4B | 76.2 | 78.6 | 79.4 | 75.8 | 77.4 | **79.0** |
| Mamba-M | 1.4B | 74.6 | **79.6** | <u>79.5</u> | 73.4 | 77.5 | 77.3 |

Table 3: Paragraph-level results in terms of COMET for models trained from scratch (top) and models finetuned from a pretrained checkpoint (bottom) on WMT23 EN↔DE test set, according to the training dataset: original WMT23-6M, WMT23-CAT-5 and WMT23-CAT-10. **Bold** represents top results; <u>underline</u> represents second-best.



Figure 2: Sensitivity to input length, measured by the number of sources tokens, on the WMT23 DE→EN datset, for models trained from scratch (top) and finetuned from a pretrained checkpoint (bottom).

and finetuned variants. Interestingly, this degradation is evident in Mamba, as expected due to its finite hidden state capacity. However, it is also challenging for transformers despite their relative positional embeddings. Moreover, both finetuned and hybrid models exhibit more consistent performance across different sequence lengths, even on the original sentence-level dataset, suggesting a more robust capability for dealing with long-context inputs.

Overall, our concatenation approach has largely mitigated the performance issues with long inputs present in models trained on WMT23-6M, as

models trained on CAT datasets produce higher-quality translations for longer sequences. This improvement is uniform across all models, with CAT-10 yielding consistently better translations in the longest bin (257+ tokens). However, the CAT-10 dataset seems to reduce translation quality for shorter samples in some models, though this effect is minimal or absent in hybrid and finetuned models. Next, we further examine the ability to extrapolate to even longer sentences than those seen during training.

Figure 3: Sensitivity to input length, measured by the number of sources tokens, on the Ted Talks DE→EN dataset, for models trained from scratch (top) and finetuned from a pretrained checkpoint (bottom). The dashed vertical line represents the bin containing the longest sentence in the training set.

## 5.3 Extrapolation to Longer Sequences

Following the previous discussion, to further explore the impact of sequence length on our models, we create a new test set sampled from TED Talks DE→EN passages that is larger (2200 samples) and contains even longer sequences (up to 2048 tokens) than WMT23. Details on this dataset can be found in Table 1. The source length distribution can be seen in Figure 7. After evaluating our models in this dataset, we plot COMET scores per sentence length in Figure 3, where we include a dashed vertical line representing the bin containing the longest sentences the model has been exposed to during training.

**Discussion.** When training from scratch, we highlight the sharp decline in translation quality decline in the Transformer++ model when considering samples larger than those it has been exposed to during training, this finding is consistent with the generalization task in (Jelassi et al., 2024). In contrast, Transformer Encoder-Decoder and Mamba exhibit a steady decline overall with the first being robust to generalization problems when trained with larger-context datasets. Additionally, the hybrid models prove to excel at generalization, providing good translation quality even when trained with the WMT23-6M dataset. With the finetuned models, we also see decreasing translation quality over longer sequences which is consistent with previous experiments. Nonetheless, Mamba models show

a more robust trend when generalizing to unseen lengths. In particular, the larger Mamba-M, when trained on the WMT23-CAT-10 dataset, exhibits a much lower performance degradation on longer samples in comparison to other finetuned models.

## 5.4 Inference Cost

In §2 we covered the theoretical time complexity of our models in training and inference time. Here, we examine the wallclock time and memory usage of Pythia and Mamba in a realistic setting where inputs are WMT23 DE→EN test samples, and outputs continue to be generated until they reach exactly $L \in \{512, 1024\}$ tokens. Table 4 shows that Mamba's memory usage is significantly lower than Pythia's, with gaps of $\sim$ 3-5x overall. The wallclock time difference is not as notable but still substantial, especially for larger models, where Mamba-M is 2x faster than Pythia-M for $L = 1024$. In other words, Mamba-M throughputs $\sim$806 tokens/s while Pythia-M outputs $\sim$405 tokens/s, aligning with (Gu and Dao, 2023).[9]

## 6 Related Works

**Linear recurrent models for MT.** Our work is closely related to (Vardasbi et al., 2023), which compares SSMs and transformers. Furthermore, they also experiment with hybrid architectures composed of S4 and attention layers, an approach that has since become common (Arora et al., 2024b; De

[9]Computed as batch-size $\times$ $L$/wallclock-time.

| Model | 512 | | 1024 | |
|---|---|---|---|---|
| | T (s) | M (GB) | T (s) | M (GB) |
| Pythia-S | 11.52 | 2.472 | 25.80 | 3.934 |
| Mamba-S | 10.38 | 0.839 | 20.59 | 1.607 |
| Pythia-M | 14.88 | 4.789 | 40.41 | 7.841 |
| Mamba-M | 10.29 | 0.913 | 20.31 | 1.668 |

Table 4: Average time (T) and maximum allocated memory (M) of 30 inference runs with batch size 16 on WMT23 DE→EN.

et al., 2024; Glorioso et al., 2024). In this work, we experiment with more recent linear recurrent models and their respective hybrid versions while also including larger and pretrained variants. Our analysis further includes investigating each model's ability to recall named entities, along with measuring translation performance across different sequence lengths on paragraph-level datasets. In contrast to Vardasbi et al. (2023)'s results showing that S4 lags behind transformer baselines in MT tasks, we observe that Mamba, a modern SSM, is competitive with transformers on sentence and paragraph-level datasets, whether trained from scratch or fine-tuned from a pretrained checkpoint, especially in the first setting when equipped with attention mechanisms.

**Linear recurrent models' limitations.** Recent works show that Mamba struggles in tasks that involve recalling context tokens (Arora et al., 2024a; Jelassi et al., 2024), such as the synthetic Multi-Query Associative Recall task. In MT, however, context tokens (source and translation prefix) are not often replicated in the output (translation). In this work, we study this phenomenon with named entities and analyze the recall ability of transformers and linear recurrent models in §4.2.

**Sentence concatenation** Kondo et al. (2022); Varis and Bojar (2021) analyze transformers' generalization towards sequence length. They show that transformers are susceptible to the training distribution of context length and that concatenating multiple sentences can improve the translation of longer sentences. Specifically, Kondo et al. (2022) augment the original data with samples containing concatenations of two random sentences, while Varis and Bojar (2021) concatenate up to six sentences. While these studies focused on sentence-level translation with sequence lengths up to 120 tokens, in this work, we extend the analysis to much longer sequences and test on paragraph-level data from the WMT2023 dataset.

## 7 Conclusion

We set out to evaluate recent linear recurrent models, particularly RetNet and Mamba, in MT tasks while thoroughly comparing them to transformer baselines and hybrid models, which combine Mamba and attention. We find that Mamba models are competitive with transformers, both when they are trained from scratch and when they are finetuned from a pretrained checkpoint; however, the performance delta is smaller in the latter regime. Our paragraph-level experiments reveal that models are hindered by the mismatch in the training and test length distributions; however, a simple concatenation approach helps to mitigate the issue. We find that hybrid models are only slightly affected by this issue while also being competitive or outperforming transformers. Finally, we note that Mamba models also exhibit a faster runtime, consume less memory, and extrapolate better to longer inputs than decoder-only transformers.

## Limitations

We point out some limitations of the presented study. First, one limitation is that we refrain from pretraining the hybrid models due to the high associated compute costs. To this effect, while our trained-from-scratch results are promising, validating them with a larger scale and strong language priors would strengthen our claim of their good performance. Secondly, our experiments (§5.3) appear to indicate larger models are more robust to sequence length issues. Nonetheless, we limited our study to models with parameter scales between 370M and 1.4B since, in preliminary sentence-level experiments, translation quality gains plateaued at the latter scale.

In another direction, we mainly rely on automated metrics for evaluating translation quality,

which might not fully capture the accuracy of the translation. We alleviate this fault by considering the recollection of NEs in translations (§4.2). Furthermore, our experiments in §5.2 do not have a notion of translation difficulty, which might help explain the differences between models and associated datasets in different length buckets (albeit sentence length and difficulty may be connected).

## Potential Risks

Translation biases and error modes inherent in transformed-based LLMs could also be manifested in the linear recurrent models studied in this paper. Careful evaluation and mitigation strategies, such as detecting and overcoming hallucinations (Guerreiro et al., 2023; Dale et al., 2023), can alleviate these risks and ensure models' responsible use. It should also be noted that although SSMs are potentially more energy efficient than transformer-based models, they still pose energy consumption concerns, particularly due to the large size of the models.

## References

Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. 2024. In-context language learning: Architectures and algorithms. In *Forty-first International Conference on Machine Learning*.

Ido Amos, Jonathan Berant, and Ankit Gupta. 2024. Never train from scratch: Fair comparison of long-sequence models requires data-driven priors. In *The Twelfth International Conference on Learning Representations*.

Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Re. 2024a. Zoology: Measuring and improving recall in efficient language models. In *The Twelfth International Conference on Learning Representations*.

Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, James Zou, Atri Rudra, and Christopher Re. 2024b. Simple linear attention language models balance the recall-throughput tradeoff. In *Forty-first International Conference on Machine Learning*.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer.

Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar Van Der Wal. 2023. Pythia: a suite for analyzing large language models across training and scaling. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Steven Bird. 2006. NLTK: The Natural Language Toolkit. In *Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions*, pages 69–72, Sydney, Australia. Association for Computational Linguistics.

Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuhito Sudoh, Koichiro Yoshino, and Christian Federmann. 2017. Overview of the IWSLT 2017 evaluation campaign. In *Proceedings of the 14th International Conference on Spoken Language Translation*, pages 2–14, Tokyo, Japan. International Workshop on Spoken Language Translation.

Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. WIT3: Web inventory of transcribed and translated talks. In *Proceedings of the 16th Annual Conference of the European Association for Machine Translation*, pages 261–268, Trento, Italy. European Association for Machine Translation.

Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019. Generating long sequences with sparse transformers.

David Dale, Elena Voita, Loic Barrault, and Marta R. Costa-jussà. 2023. Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity Even better. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 36–50, Toronto, Canada. Association for Computational Linguistics.

Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Albert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Desjardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas, and Caglar Gulcehre. 2024. Griffin: Mixing gated linear recurrences with local attention for efficient language models.

Patrick Fernandes, Kayo Yin, Graham Neubig, and André F. T. Martins. 2021. Measuring and increasing context usage in context-aware machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6467–6478, Online. Association for Computational Linguistics.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references

are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, and Christopher Re. 2023. Hungry hungry hippos: Towards language modeling with state space models. In *The Eleventh International Conference on Learning Representations*.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The pile: An 800gb dataset of diverse text for language modeling.

Paolo Glorioso, Quentin Anthony, Yury Tokpanov, James Whittington, Jonathan Pilault, Adam Ibrahim, and Beren Millidge. 2024. Zamba: A compact 7b ssm hybrid model.

Albert Gu and Tri Dao. 2023. Mamba: Linear-time sequence modeling with selective state spaces.

Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. 2020. Hippo: Recurrent memory with optimal polynomial projections. In *Advances in Neural Information Processing Systems*, volume 33, pages 1474–1487. Curran Associates, Inc.

Albert Gu, Karan Goel, and Christopher Re. 2022. Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations*.

Nuno M. Guerreiro, Elena Voita, and André Martins. 2023. Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1059–1075, Dubrovnik, Croatia. Association for Computational Linguistics.

Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc Le. 2022. Transformer quality in linear time. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 9099–9117. PMLR.

Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and eran malach. 2024. Repeat after me: Transformers are better than state space models at copying. In *Forty-first International Conference on Machine Learning*.

Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. 2020. Transformers are rnns: fast autoregressive transformers with linear attention. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Seiichiro Kondo, Naoya Ueda, Teruaki Oka, Masakazu Sugiyama, Asahi Hentona, and Mamoru Komachi. 2022. Japanese named entity recognition from automatic speech recognition using pre-trained models. In *Proceedings of the 36th Pacific Asia Conference on Language, Information and Computation*, pages 102–108, Manila, Philippines. Association for Computational Linguistics.

Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, Xingjian Du, Matteo Grella, Kranthi Gv, Xuzheng He, Haowen Hou, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong, Bartłomiej Koptyra, Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song, Xiangru Tang, Johan Wind, Stanisław Woźniak, Zhenyuan Zhang, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. 2023. RWKV: Reinventing RNNs for the transformer era. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14048–14077, Singapore. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T.

Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Noam Shazeer. 2020. Glu variants improve transformer.

Jimmy T.H. Smith, Andrew Warrington, and Scott Linderman. 2023. Simplified state space layers for sequence modeling. In *The Eleventh International Conference on Learning Representations*.

Alex J Smola and Bernhard Schölkopf. 1998. *Learning with kernels*, volume 4. Citeseer.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958.

Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomput.*, 568(C).

Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. 2023a. Retentive network: A successor to transformer for large language models.

Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. 2023b. A length-extrapolatable transformer. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14590–14604, Toronto, Canada. Association for Computational Linguistics.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.

Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2019. Transformer dissection: An unified understanding for transformer's attention via the lens of kernel. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4344–4353, Hong Kong, China. Association for Computational Linguistics.

Ali Vardasbi, Telmo Pessoa Pires, Robin Schmidt, and Stephan Peitz. 2023. State spaces aren't enough: Machine translation needs attention. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 205–216, Tampere, Finland. European Association for Machine Translation.

Dusan Varis and Ondřej Bojar. 2021. Sequence length is a domain: Length-based overfitting in transformer models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.

Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. 2024. Gated linear attention transformers with hardware-efficient training. In *Forty-first International Conference on Machine Learning*.

## A  Implementation and Training Details

All experiments were carried on Nvidia RTX A6000 GPUS with 48GB VRAM, and the training framework is constructed around PyTorch Lightning.[10] To train and generate translations in batches, we use a left-padding strategy. However, for Mamba, additional functionality is required to avoid processing padding tokens. To address this, we zero out inputs before and after convolution at the positions of the padding tokens and sacrifice some efficiency by using the slow path in Mamba.[11] Notably, during inference, the slow path affects only the initial processing of the prompt and does not impact the actual generation. Moreover, we added Dropout (Srivastava et al., 2014) to Mamba blocks, which was missing in the original implementation. Specifically, dropout is applied after the last linear projection of the Mamba block. Additionally, following the findings in (Vardasbi et al., 2023), we calculate cross-entropy loss only for target tokens. During training, we use greedy decoding and select the top model using BLEU as the validation metric, as it is faster to compute in comparison to COMET. For inference, we use beam search with a beam size of 5. Due to the time-consuming nature of our experiments, we report the results of a single run for all experiments. The overall model structure and hyperparameters across both training regimes, from-scratch (§A.1) and finetuning (§A.2), are shown in Table 5. Furthermore, all models were trained with `bfloat16` precision.

### A.1  Training from Scratch

Regarding tokenization, we leverage the Hugging-Face *tokenizers* library[12] and construct a separate BPE tokenizer (Sennrich et al., 2016) per dataset. The total vocabulary size is 32000 tokens. We carried out a hyperparameter search to select appropriate dropout values, learning rates and architectural decisions, with the latter two detailed in Table 5. We employ a dropout of $0.3$ for WMT16 EN↔RO, $0.1$ for WMT14 EN↔DE, WMT16 EN↔FI and the different variations of WMT23. Other hyperparameters were kept intact. Concretely, we use the Inverse Square Root learning rate scheduler (Vaswani et al., 2017) with 4000 warmup steps and a weight

---

[10] https://lightning.ai/docs/pytorch/
[11] https://github.com/state-spaces/mamba/issues/216
[12] https://github.com/huggingface/tokenizers

| MODEL | SIZE | LR | L | H | D | FFN |
|---|---|---|---|---|---|---|
| *Trained from scratch* | | | | | | |
| Transf. Enc-Dec | 77M | 4e-4 | 6-6 | 8 | 512 | 2048 |
| Transf.++ | 79M | 4e-4 | 12 | 8 | 496 | 1984 |
| RetNet | 77M | 1e-3 | 12 | 4 | 512 | 1024 |
| Mamba | 77M | 1e-3 | 24 | - | 610 | - |
| Mamba-MHA | 78M | 7e-4 | 24 | 8 | 624 | - |
| Mamba-Local | 78M | 7e-4 | 24 | 8 | 624 | - |
| Mamba Enc-Dec | 82M | 7e-4 | 8-6 | 8 | 512 | 2048 |
| *Finetuned* | | | | | | |
| Pythia-S | 410M | 1e-5 | 24 | 16 | 1024 | 4096 |
| Mamba-S | 370M | 3e-4 | 24 | - | 1024 | - |
| Pythia-M | 1.4B | 1e-5 | 24 | 16 | 2048 | 8192 |
| Mamba-M | 1.4B | 3e-4 | 24 | - | 2048 | - |

Table 5: Detailing the full set of hyperparameters for the different models. Encoder-Decoder models have their number of layers separated by each module. LR represents the Learning Rate; L represents the number of layers; H is the number of Attention Heads; D is the model dimension; FFN is the size of the inner feed-forward network.

| MODEL | SIZE | TRAINING TOKENS | CONTEXT TOKENS |
|---|---|---|---|
| Pythia-S | 410M | 300B | 2048 |
| Pythia-M | 1.4B | 300B | 2048 |
| Mamba-S | 370M | 7B | 2048 |
| Mamba-M | 1.4B | 26B | 2048 |

Table 6: Pre-training details. All models were pretrained on The Pile (Gao et al., 2020).

decay of $0.001$.

### A.2  Finetuning Pretrained Checkpoints

We employ pretrained models and corresponding tokenizers from the Huggingface library. Table 6 shows the number of tokens and the size of the context window used during pretraining. For finetuning, in all experiments, we use a dropout of 0.1 with the exception of WMT16 EN↔RO and Pythia-S + EN↔FI, where dropout varies from 0.1 to 0.3 for the former and 0 for the latter. Moreover, we use weight decay only in Mamba-M, with a value of $2 \cdot 10^{-4}$. Additionally, learning rates and models' attributes are shown in Table 5.

### A.3  Inference Cost

For the inference cost experiments, we measure overall wallclock time using cuda events and cuda synchronization from `torch.cuda` module. The overall reported time measures the entire generation pipeline, including the use of beam search. Moreover, we use

| MODEL | 512 | | 1024 | |
|---|---|---|---|---|
| | T (S) | M (GB) | T (S) | M (GB) |
| Transformer++ | 12.33 | 5.862 | 34.00 | 10.711 |
| Mamba | 11.71 | 0.562 | 29.37 | 0.554 |
| Mamba-MHA | 12.77 | 1.250 | 25.28 | 1.536 |
| Mamba Enc-Dec | 7.46 | 0.394 | 14.36 | 0.394 |

Table 7: Average time (T) and maximum allocated memory (M) of 30 inference runs with batch size 16 on WMT23 DE→EN.

| | DE→EN | | EN→DE | |
|---|---|---|---|---|
| | BLEU | COMET | BLEU | COMET |
| *Mamba-MHA* | | | | |
| Interleaved | 30.81 | 77.98 | 24.40 | 72.48 |
| L1,11 | 30.52 | 78.10 | **24.99** | 73.76 |
| L11,23 | **30.81** | **78.30** | 24.40 | **73.94** |
| *Mamba-Local* | | | | |
| Interleaved - w64 | 28.85 | 76.76 | 23.61 | 72.10 |
| L11,23 - w16 | 29.37 | 77.19 | 24.12 | 72.88 |
| L11,23 - w32 | 28.24 | 76.44 | 23.20 | 72.22 |
| L11,23 - w64 | 29.40 | 77.56 | 24.41 | 72.98 |
| L11,23 - w128 | 30.49 | 77.98 | 24.85 | 73.58 |

Table 8: Hybrid models ablations with BLEU and COMET scores on the IWSLT17 dataset. Different window sizes are denoted as w$\{16, 32, 64, 128\}$. *Interleaved* refers to alternating Mamba and attention layers. L*1,11* and L*11,23* refer to placing attention in layers 2 - $N/2$ and $N/2$ - $N$, respectively.

`torch.cuda.max_memory_allocated` to measure memory usage.

We additionally include the profiling measurements for the trained-from-scratch models in Table 7. Crucially, we advise that these metrics are rough estimates since the models are not optimized to perform at their best capacity. To this end, we do not include the Transformer Encoder-Decoder as the implementation used is not efficient.

## B  Hybrid Models Ablation

Building on the shortcomings of linear models (Akyürek et al., 2024; Arora et al., 2024a; Jelassi et al., 2024), we designed hybrid models to complement SSMs with attention mechanisms. In this section, we ablate the design choices leading to the construction of our hybrid models. These experiments were conducted using the IWSLT17 DE↔EN dataset (Cettolo et al., 2017). Results are shown in Table 8.

Since our Mamba-MHA model replaces a set of Mamba layers with attention modules, we ablated various configurations to determine the optimal

number and placement of attention layers. Our analysis of COMET scores indicated that incorporating two attention layers significantly boosted performance, aligning with findings in previous studies (Fu et al., 2023). The placement of these layers had a minimal effect, leading us to select the configuration with layers at positions $N/2$ and $N$ for further experiments due to its consistently higher COMET scores.

In the case of Mamba-Local, which uses a sliding window attention, we explored various window sizes. Our experiments revealed that performance generally improved with window size in a linear way. Ultimately, a 128-token window nearly matched full attention performance, and two layers of 64-token windowed attention provided a good balance between performance and efficiency for our experiments.

## C  Named Entity Recall Experiments

Following up on the discussion from §4.2, we extend our evaluation of NE recall accuracy to the WMT14 DE↔EN dataset and two paragraph-level datasets, WMT23-6M and WMT23-CAT-5, both in the DE↔EN translation direction. The results, detailed in Figure 4, offer further insights into the models' recall accuracy performance across other datasets and context length settings.

**Sentence-Level (WMT14 DE↔EN).** The NE recall results on the WMT14 DE↔EN dataset align closely with those obtained in WMT16 RO→EN, shown in Figure 1; we still observe Mamba's recall accuracy to be closer to that of the transformer models, while the hybrid models continue to (slightly) outperform their unmodified counterparts. Note, however, that overall, the gap between models is narrower, as also reflected in their close results in terms of BLEU.

**Paragraph-Level Datasets.** When assessing the WMT23-6M and WMT23-CAT-5 DE↔EN datasets, contrary to the WMT16 RO↔EN experiments, the Transformer Encoder-Decoder model outperforms all other models in recalling unseen entities. Additionally, while the hybrid models remain comparable to the Transformer++ model, Mamba's performance declines. This presents a striking contrast to the sentence-level experiments, suggesting that transformers may have an advantage in NE recall when shifting to longer contexts. Nonetheless, the transition from the 6M dataset to the CAT-5 dataset

Figure 4: Recall in recovering named entities on the WMT14 (top), WMT23-6M (middle) and WMT23-CAT-5 (bottom) DE→EN datasets, by their training set frequency: *unseen* entities do not appear in the training data, while *regular* and *frequent* entities appear $[1, 16)$ and $16+$ times.



Figure 5: COMET scores per sequence length on WMT14 DE→EN for trained-from-scratch models.

leads to recall improvements across all models, particularly for unseen entities. This indicates that the additional context provided during training in the CAT-5 dataset aids the recall of named entities.

## D Exploring Length-related Issues

### D.1 Preliminary Sentence-level Experiments

Before experimenting with paragraph-level data, we analyze how our trained-from-scratch models perform on different sequence lengths. To this

end, we study their sensitivity to input length when trained and tested on WMT14 DE→EN. The results are shown in Figure 5. While all models show a deterioration in performance as sequence length increases, this effect is more pronounced for Transformer++, RetNet, and Mamba-Local, with a significant drop in performance for samples longer than 64 tokens.

### D.2 Sensitivity to Input Length

Following the discussion in §5.2, we further investigate the sensitivity of our models to input length using the WMT23 EN→DE test set, with results shown in Figure 6. Notably, our takeaways remain broadly the same: concatenating samples in the training data is indeed helpful when handling longer sequences, and models trained on the WMT23-CAT-10 dataset are much better in the longer bin (257+) with minimal translation quality degradation in shorter samples. However, when considering each of the training datasets' histograms in Figure 7, we can observe that models have been exposed to the longest samples during training, even if in low quantities. This implies that the previous experiments do not represent an extrap-

Figure 6: Sensitivity to input length, measured by the number of sources tokens, on the WMT23 EN→DE datset, for models trained from scratch (top) and finetuned from a pretrained checkpoint (bottom).

olation setting, where inference is done on longer sequence lengths than those seen during training. We cover extrapolation to longer sequences next.

## E    Full Paragraph-Level Results

For completeness, we report paragraph-level results in terms of BLEU and COMET for all language pairs and models in Table 9.

## F    AI assistants

We have used Github Copilot[13] during code development, and ChatGPT[14] during paper writing for paraphrasing or polishing original contents.

Figure 7: Distribution of source length in 1) the training datasets: WMT23 DE→EN (top left), WMT23 EN→DE (top right), and 2) the test datasets: WMT23 DE↔EN (bottom left), our custom TED Talks DE→EN (bottom right).

| MODEL | TRAINING DATA | DE→EN | | EN→DE | |
|---|---|---|---|---|---|
| | | BLEU | COMET | BLEU | COMET |
| *Trained from scratch* | | | | | |
| Transformer Enc-Dec | | 25.4 | 72.4 | 22.4 | 65.2 |
| Transformer++ | | 21.6 | 70.7 | 20.2 | 64.8 |
| Mamba | WMT23-6M | 19.0 | 70.0 | 15.8 | 63.3 |
| Mamba-MHA | | **23.9** | **72.7** | **23.2** | **67.0** |
| Mamba Enc-Dec | | 22.7 | 70.7 | 21.5 | 65.3 |
| Transformer Enc-Dec | | **30.8** | **74.6** | **29.9** | 70.3 |
| Transformer++ | | 28.9 | 73.6 | 28.1 | 69.1 |
| Mamba | WMT23-CAT-5 | 26.1 | 73.3 | 23.8 | 67.5 |
| Mamba-MHA | | 29.5 | 74.2 | 23.5 | 68.6 |
| Mamba Enc-Dec | | 27.3 | 73.8 | 29.1 | **71.0** |
| Transformer Enc-Dec | | 28.3 | 69.6 | 29.3 | **70.3** |
| Transformer++ | | 29.8 | 72.8 | 29.1 | 68.8 |
| Mamba | WMT23-CAT-10 | 25.9 | 72.3 | 25.5 | 67.8 |
| Mamba-MHA | | 27.8 | 74.5 | 25.9 | 69.7 |
| Mamba Enc-Dec | | **31.4** | **75.6** | **30.1** | 70.1 |
| *Finetuned* | | | | | |
| Mamba-S | | 21.8 | 77.2 | 21.4 | 72.4 |
| Pythia-S | WMT23-6M | 23.9 | **77.4** | **25.9** | **76.7** |
| Mamba-M | | 20.7 | 74.6 | 22.5 | 73.4 |
| Pythia-M | | **26.0** | 76.2 | 25.2 | 75.8 |
| Mamba-S | | 24.3 | 78.2 | 23.3 | 74.2 |
| Pythia-S | WMT23-CAT-5 | **27.0** | 78.4 | **28.6** | **77.8** |
| Mamba-M | | 26.4 | **79.6** | 27.5 | 77.5 |
| Pythia-M | | 25.8 | 78.6 | 27.5 | 77.4 |
| Mamba-S | | 25.6 | 78.3 | 22.5 | 73.1 |
| Pythia-S | WMT23-CAT-10 | 26.8 | 79.0 | 29.3 | 77.1 |
| Mamba-M | | 32.5 | **79.5** | 27.5 | 77.3 |
| Pythia-M | | **33.4** | 79.4 | **33.9** | **79.0** |

Table 9: Paragraph-level results in terms of BLEU and COMET on the WMT23 EN↔DE test set.

# Evaluation and large-scale training for contextual machine translation

**Matt Post** and **Marcin Junczys-Dowmunt**
Microsoft
{mattpost,marcinjd}@microsoft.com

## Abstract

Despite the fact that context is known to be vital for resolving a range of translation ambiguities, most traditional machine translation systems continue to be trained and to operate at the sentence level. A common explanation is the lack of document-level annotations for existing training data. This work investigates whether having such annotations would be helpful for training traditional MT systems at scale. We build large-scale, state-of-the-art contextual MT systems into German, French, and Russian, fixing the datasets while comparing the effect of sourcing contextual training samples from both parallel and back-translated data. We then evaluate these contextual models across a range of contextual test sets from the literature, where we find that (a) document annotations from both mined parallel and back-translated monolingual data are helpful, but that the best contextual MT systems do not draw *contextual* training samples from the parallel data. We also make two points related to evaluation: (b) contrastive score-based metrics on challenge sets are not discriminative; instead, models must be tested directly on their ability to generate correct outputs, and (c) standard corpus-level metrics such as COMET work best in settings that are dense in contextual phenomena.

## 1 Introduction

By nature of its sentence-based design, traditional machine translation (MT) is unable to correctly translate any sentence with extra-sentential dependencies, such as pronouns in languages with grammatic gender, except by chance (Table 1). Despite significant prior work on the topic (§ 2), and general acknowledgment of the need to move on (Sennrich, 2018), contextual translation has never managed to overtake MT research, and sentence-level systems continue to dominate. This "sentence-level ceiling" leaves a gap between them and their increasingly powerful LLM counterparts, and raises

| English | German |
|---|---|
| I lost my hat. *Have you seen it?* | Ich verlor meinen Hut. *Hast du **es** gesehen?* |

Table 1: The sentence-level translation ceiling. Selecting the correct pronoun (*ihn*, masc.) requires context.

the question of whether this gap can be narrowed or closed, if traditional MT systems could be trained properly with context.

A common explanation for the lack of context in MT has to do with the relative dearth of document-level annotations that are available for mined parallel and even monolingual data. At the same time, it has long been understood (Venugopal et al., 2011) and recently corroborated (Thompson et al., 2024) that crawled bitext is rife with machine translation output, which—though high quality at the sentence level—may attenuate the contextual signal. We explore this central problem by building the first large-scale, state-of-the-art translation systems trained on data with complete document annotations. We are able to do this because instead of public data, we use a private, in-house dataset (§ 3) that we have crawled ourselves. This crucially allows us to explore the effects of document annotations sourced from both parallel and monolingual (backtranslated data), together and in isolation, in order to quantify their effects. We find that:

- **It is best to source contextual training examples from backtranslated data only**. We find gains in contextual metrics from systems trained with contextual signals from both parallel and backtranslated data. However, the best systems source these samples from backtranslated data only.

- **Generative evaluation is crucial**. Contrastive metrics, where the task is to discriminate good and bad translations using model

scores, are often used to evaluate contextual MT. We show that contextual systems that are trained on mined parallel documents do well on this task, but perform poorly when asked to generate correct translations. Only generative evaluation, which looks at whether correct words were produced, distinguishes good from bad contextual systems.

- **Standard metrics are most useful on discourse-dense datasets**. Standard sentence-level metrics like COMET are much more discriminative between sentence- and contextual systems when applied to datasets that are dense in discourse phenomena.

Together, these results raise important considerations for the construction and evaluation of contextual translation systems.

## 2 Background and Related Work

The transition to neural architectures was a paradigm enabler for document translation, since it eliminated the Markov limitations of statistical MT (Maruf et al., 2019). Much work has focused on special architectures and input encodings. This includes cache models (Tu et al., 2018; Kuang et al., 2018), hierarchical attention (Miculicich et al., 2018), separately encoding context (Voita et al., 2018; Zhang et al., 2018), allowing attention across a batch of pseudo-documents (Wu et al., 2023), encoding sentence position (Bao et al., 2021; Lupo et al., 2023), and sparse attention mechanisms (Guo et al., 2019). A number of approaches work on base systems outputs, such as post-editing with contextual language models (Voita et al., 2019a) and using contextual language models to rerank sentence-level system output Yu et al. (2020). Junczys-Dowmunt (2019) built one of the earliest contextual systems to perform well at WMT. Sun et al. (2022) also proposed to use standard transformer models, testing small architectures with no backtranslated data, and using a "multi-resolutional" training approach that creates overlapping documents. **We focus instead on standard architectures, judging them to be sufficient at large enough sizes.**

The lack of document-annotated parallel data is a longstanding problem. Datasets with document annotations are relatively small and specialized: they include OpenSubtitles (Lison and Tiedemann, 2016), WIT[3] (Cettolo et al., 2012), News Commentary, and Europarl (Koehn, 2005). Liu and

Zhang (2020) provide a nice survey, and release a small amount of government-crawled new data for Chinese–Portuguese. Very recently, document annotations on Paracrawl data have become available Pal et al. (2024); Wicks et al. (2024). These annotations are available for only a relatively small subset of the data, however; even so, their results corroborate what we find here. (2024, Table 2) see *drops* in performance from systems trained with their parallel data annotations, unless the gold target context is provided; (2024) see small but consistent gains when the parallel data has been sufficiently filtered. The Conference on Machine Translation (WMT) began releasing limited document-level data for DE-EN and CS-EN in 2019 (Barrault et al., 2019). This limitation has forced researchers to get creative. Voita et al. (2019b) built a monolingual post-editing system that took the output of a baseline system and used it for document-level "repair". Sugiyama and Yoshinaga (2019) also used target-side data for backtranslation, evaluating in small-data settings with BLEU and contrastive metrics. **Our work is unique in that we have complete document annotations on very large web-crawled datasets**, and shows that these annotations on parallel data, as a whole, are not as useful.

Contextual metrics work has been important. PROTEST (Guillou and Hardmeier, 2016) used hand-designed pronoun test cases and also evaluated generatively. Many special test sets have been developed isolating important contextual phenomena and largely evaluating discriminatively (more in § 4). Läubli et al. (2018) provided early evidence that document-level metrics would be helpful. Several document-level metrics have been proposed, including BlonDe (Jiang et al., 2022), which compares automatically-identified phenomena in the output to those in a reference, and Doc-COMET (Vernikos et al., 2022), which incorporates contextual sentence representations. Both metrics are interesting but await deeper evaluation and we did not explore them in this paper. Vamvas and Sennrich (2021) have also noted the problem with the disconnect between contrastive evaluation and generative ability for machine translation. Both Fernandes et al. (2023) and Wicks and Post (2023) developed rules to identify contextually-dependent sentences. In this work, we show that **datasets dense in contextual phenomena are important for evaluating contextual ability**, and that **discriminative contextual evaluation is of limited use**.

## 3 The data challenge

Large publicly-available parallel datasets do not have document annotations. While the Conference on Machine Translation (WMT) has made overtures in this direction,[1] including ensuring that test data is source-language-natural and contains document information, parallel and monolingual data is limited to a small subset of all data[2] for which such information is easily retained. This is also true of recent work extracting document annotations from Paracrawl (Pal et al., 2024; Wicks et al., 2024).

We wish to experiment with and compare annotations sourced from both parallel and back-translated monolingual datasets. We therefore turn instead to a state-of-the-art, large collection of in-house data. We work with three language pairs: English→German, English→French, and English→Russian, which were selected because of the availability of good contextual evaluation data in each of them (§ 4). Our data comprises the following sources (Table 2):

- Monolingual data, crawled from expected-native sites: news (10%), data linked from the Open Directory Project[3] (40%), filtered webcrawl (40%), and Wikipedia and its outlinks (10%).

- Crawled parallel web data (similar to ParaCrawl)

- CCMatrix parallel data (Schwenk et al., 2021b), which has no document information.

Datasets have been filtered using bicleaner (Ramírez-Sánchez et al., 2020), with additional boilerplate and document deduplication.

Although the dataset is private, there is nothing in it that would surprise any researcher; the data was crawled from the web using standard techniques. The parallel data sources include a rough equivalent of ParaCrawl (Bañón et al., 2020) and also CCMatrix (Schwenk et al., 2021b). The monolingual data sources focus on sites where we expect data to have been written natively.

We emphasize that experiments at the scale presented in this paper are only possible with our

private dataset, since document annotations are only available for small-data training settings like the TED talks data (Cettolo et al., 2012) used by IWSLT.[4] In a nod to the importance of repeatable work, we include results on the subset of our experiments that are possible on English–German public data and show that they corroborate corresponding results on private data (Section 7.6).

## 4 Contextual evaluation

A basic hurdle in the path to contextual translation is the difficulty of evaluation. We expect that contextual systems will produce improved translations of discourse-level phenomena, however, the frequency of these phenomena in standard corpora is not known, and we expect them to be relatively rare. This paper includes three types of evaluation.

### 4.1 Corpus-level metrics

The conventional way to test system performance is with standard metrics such as chrF (Popović, 2015) or COMET (Rei et al., 2020), which accumulate sentence-level scores to compute a single score for a test set. If the test set is organized into documents (as many are, including those from WMT), its sentences can be translated contextually and then split back out to sentences for evaluation. The expectation is that contextual translation will produce gains. However, a key consideration is whether the dataset is dense enough with contextual phenomena. Attempts to automatically identify sentences requiring context have shown the task to be difficult (Bawden et al., 2018) though possible with hand-created rules (Fernandes et al., 2023; Wicks and Post, 2023), but are often rare. Consequently, improvements may be invisible without the right test set.

We compare the performance of contextual systems using a standard corpus-level metric, COMET[5], on the following two test sets:

- WMT15. We use newstest2015 (Bojar et al., 2015) for EN→FR, and newstest2022 for EN→DE and EN→RU (Kocmi et al., 2022).

- OpenSubtitles (Lison and Tiedemann, 2016). We use the CTXPro (Wicks and Post, 2023) gender dataset, which is large and focuses on pronouns and anaphora.

---

[1] statmt.org

[2] Parallel: europarl, news-commentary, CzEng, and Rapid; Monolingual: news-crawl (en, de and cs), europarl, and news-commentary. Source: http://www2.statmt.org/wmt23/translation-task.html

[3] https://odp.org

[4] iwslt.org

[5] wmt20-comet-da

| | English–French | | | English–German | | | English–Russian | | |
|---|---|---|---|---|---|---|---|---|---|
| source | lines | docs | mean | lines | docs | mean | lines | docs | mean |
| mono | 166.4 | 5.5 | 29.7 | 205.4 | 7.0 | 29.1 | 202.7 | 6.5 | 31.1 |
| parallel | | | | | | | | | |
| → crawled | 123.1 | 3.7 | 33.0 | 116.7 | 4.7 | 16.6 | 72.4 | 4.7 | 13.2 |
| → ccmatrix | 65.1 | 0 | - | 45.4 | 0 | - | 2.4 | 0 | - |

Table 2: Statistics of the training data used in our experiments (lines and docs in millions). The *mean* column is the mean document length in sentences of documents with $\geq 2$ sentences.

Because the CTXPro dataset was constructed to select them, we expect it to be much denser in discourse phenomena. Data sizes are listed above the results in Table 3.

## 4.2 Contrastive test sets

The dominant paradigm for evaluation of long-tail document phenomena has been so-called *contrastive evaluation* (Sennrich, 2017), in which a system is tested on its ability to discriminate (via assigned model score) between correct and incorrect translation pairs. The correct examples are usually taken from found text; the incorrect ones are created by inserting an error of some sort. We look at three such test sets, examples of which can be found in Appendix A.

**ContraPro (EN-DE)** Müller et al. (2018) focus on the German pronouns *es*, *er*, and *sie*. They pair sentences containing naturally-found instances of pronouns drawn from OpenSubtitles with two variants where the incorrect pronoun has been used.

**ContraPro (EN-FR)** Lopes et al. (2020) extended ContraPro for EN-FR; the main difference is that there is only one incorrect example, since French has only two grammatical genders.

**GTWiC (EN-RU) (Voita et al., 2019b)** *Good Translation, Wrong in Context* (GTWiC) tests verb selection (500 instances) and morphology (500) in the presence of source-side ellipsis.

## 4.3 Testing generative ability

The challenge sets above test whether a model can discriminate between good and bad examples with using model score. However, this is at best a proxy for the true test of a machine translation system, which is to determine whether it generates the correct word or phrase. As we will show, many document models perform extremely well on these tasks,

but when asked to actually translate the source sentence, produce the wrong word (Table 5). The contrastive nature of these test sets is at odds with the actual task: what is needed are metrics that directly evaluate a model's *generative*, rather than its *discriminative*, ability.

Fortunately, because these test sets were distributed with rich annotation information, we can transform them into generative test sets, where we test for the correct word in the output. A test set $\mathcal{T}$ comprises a set of test examples in the form of tuples $(S, R, w)$, where $S$ is the source sentence, $R$ the reference, and $w \in R$ the target word or phrase that is expected to be found in the translation output. Let $\{T_i\}$ be the set of translations of the source sentences $\{S_i\}$. We compute accuracy[6] as

$$\text{acc}(T, \mathcal{T}) = \frac{1}{|T|} \sum_{i=1}^{|T|} \delta(w_i \in T_i)$$

This is not a perfect metric, since a correct translation may have paraphrased around the pronoun, but we do not expect that to systematically favor any particular system.

We have further opportunity to test this kind of accuracy with **CTXPro** (Wicks and Post, 2023), which expands ContraPro's coverage to many other languages and linguistic phenomena (auxiliaries, formality, gender, and inflection). CTXpro is evaluated only generatively, and has been been tested only on a single system, DeepL,[7] which is known to make use of context.

## 5 Experimental setup

We train and compare five models on the exact same data from two sources: parallel ($\mathcal{P}$) and back-translated monolingual ($\mathcal{B}$) data; the only difference among the models is whether document sam-

---

[6] Here accuracy is the same as both precision and recall.
[7] deepl.com

ples are drawn from neither, one, or both of the datasets.

**Training** All models are transformers trained with Marian (Junczys-Dowmunt et al., 2018a,b). We create two classes of models: first, those for backtranslation, and second, a set of models that constitute our primary comparative evaluation. For each language pair, we build a single joint unigram subword model (Kudo, 2018) with a vocabulary size of 32k that is used for both sets of models. Models are trained on random permutations over the training data for a predetermined number of updates. We use a batch size of 500k target-side tokens and a maximum sample length (whether sentences or pseudo-documents) of $L = 256$ tokens.

**Backtranslated data** The monolingual data is backtranslated (Sennrich et al., 2016) using sentence-level transformer systems (Vaswani et al., 2017) with 12 encoder and 6 decoder layers, an embedding size of 1024, and a feed-forward dimension of 8192. These models are trained for 20 virtual epochs.

This backtranslated data will be used to train contextual systems, but we note that this is not a problem, for two reasons. The major reason is that the target-side contextual signal is unaffected by backtranslation; since the original document boundaries are retained, any mistakes introduced by sentence-level backtranslation will appear just as normal source-side noise that the model must learn to overcome. Losses will be computed against the original, intact, target-side context. Second, even if this were not the case, our backtranslation models are into English, which is morphologically simpler than the evaluated translation direction.

**Models** For our contextual models, we also train transformers with a 12-layer encoder, a 6-layer decoder, and an embedding dimension of 1,024, but increase the feed-forward network size to 16,384. These models are trained for 40 virtual epochs to reflect the larger amounts of training data.

All of our models are trained on the complete parallel ($\mathcal{P}$) and backtranslated ($\mathcal{B}$) data. They vary only in whether the training procedure is permitted to construct multiple-sentence samples (also called *pseudo-documents* or *chunks*) from both, neither, or exactly one of these two pools of data. We compare the following systems, using the syntax NAME(pool$_1$, pool$_2$) to denote the pools of data each draws from; the presence of a box around

the data source notes that pseudo-documents were drawn from it.

- SENT($\mathcal{P}$,$\mathcal{B}$). A sentence-level baseline.
- RAND($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$). A contextual system, but trained with completely random contexts.
- DOC($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$). A contextual system, with documents from parallel and back-translated data.
- DOC($\ddot{\mathcal{P}}$,$\mathcal{B}$). A contextual system, with documents drawn from parallel data only.
- DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$). A contextual system, with documents drawn from backtranslated data only.

**Creating samples** We create our training data on the fly using SOTASTREAM (Post et al., 2023), which iterates over $\mathcal{P}$ and $\mathcal{B}$. At each iteration, each data source is permuted randomly at the document level. To generate each sample, SOTASTREAM first chooses randomly between the two data pools. If documents are disabled on the pool, it simply returns the next sentence pair. If documents are enabled, it then chooses a maximum token length, and concatenates sentences on both sides until this length is reached on the source side, or the document's end is reached. Concatenated sentences are joined with a special ⟨SEP⟩ token, which facilitates sentence alignment at inference time for evaluation. Contextual samples are *chunked*, our term for the 1:1 concatenative construction described in Tiedemann and Scherrer (2017).[8] The training toolkit is then responsible for buffering as many samples as are needed to sort and form batches for training.

**Inference** For inference, we use the *last sentence* approach as defined in Herold and Ney (2023): each input sentence (the *payload*) is prepended with left sentence context, up to a maximum token length, $L$, which includes the payload. The translation system translates this as a single unit. The ⟨SEP⟩ token is then used to extract the payload's translation. This is repeated for all sentences in a test set, allowing standard sentence-level metrics to be applied to the results.

## 6 Results

**Sentence-level metrics** We begin by establishing baseline scores with a standard corpus-level metric, COMET, in Table 3. We include a commercial

---

[8]This can be contrasted with the "multi-resolution" approach of Sun et al. (2022), which creates training samples of different lengths from many overlapping sub-sequences of each input document

| | | EN→DE | | EN→FR | | EN→RU | |
|---|---|---|---|---|---|---|---|
| | | WMT | CTXPro | WMT | CTXPro | WMT | CTXPro |
| | #lines | 1,500 | 31,640 | 2,307 | 43,375 | 2,307 | 32,948 |
| | Microsoft | 62.0 | 27.7 | 67.6 | 36.4 | 67.3 | 39.1 |
| sent-level | SENT($\mathcal{P},\mathcal{B}$) | 61.1 | 24.4 | 67.4 | 34.5 | 70.0 | 38.5 |
| | RAND($\ddot{\mathcal{P}},\ddot{\mathcal{B}}$) | 59.2 | 22.7 | 67.6 | 33.6 | 68.9 | 36.6 |
| | DOC($\ddot{\mathcal{P}},\ddot{\mathcal{B}}$) | 60.2 | 23.4 | 67.0 | 33.5 | 70.5 | 38.8 |
| | DOC($\ddot{\mathcal{P}},\mathcal{B}$) | 59.7 | 22.6 | 68.8 | 34.1 | 70.0 | 37.8 |
| | DOC($\mathcal{P},\ddot{\mathcal{B}}$) | 60.9 | 24.5 | 67.8 | 34.7 | 70.3 | 38.2 |
| context | RAND($\ddot{\mathcal{P}},\ddot{\mathcal{B}}$) | 58.8 | 20.4 | 66.8 | 32.1 | 68.7 | 35.4 |
| | DOC($\ddot{\mathcal{P}},\ddot{\mathcal{B}}$) | 60.7 | 26.9 | 67.2 | 37.8 | 69.2 | 43.2 |
| | DOC($\ddot{\mathcal{P}},\mathcal{B}$) | 60.2 | 25.4 | 67.9 | 37.6 | 68.5 | 40.3 |
| | DOC($\mathcal{P},\ddot{\mathcal{B}}$) | 60.8 | 31.6 | 68.7 | 42.2 | 70.6 | 45.8 |

Table 3: COMET20 scores on WMT (22/15) and OpenSubtitles (CTXPro/gender) test sets translating alone (top block) and with context (bottom block). Numbers within a column are comparable. The gains from DOC($\mathcal{P},\ddot{\mathcal{B}}$) (with context) over SENT($\mathcal{P},\mathcal{B}$) (without it) are much larger for the discourse-dense OpenSubtitles data.

baseline (Microsoft, accessed via API). As another baseline, we present sentence-level results for the sentence-level system trained on all of our data. We then present results for all our models translating the test corpora (WMT and OpenSubtitles, using the CTXPro/gender dataset) in two modes: at the sentence level (top block), and with context (bottom block). In this way, we can look at the effect of context at both training and inference time.

**Accuracy-based generative evaluation** Next, we look at the broader CTXPro datasets and evaluate them using word accuracy on their relevant phenomena. Table 4 contains results for all three language pairs for all CTXPro datasets.

**Contrastive suites** Finally, we turn to the document-level contrastive and generative metrics described in § 4.2–4.3. Table 5 contains results for all three language pairs.

## 7 Discussion

### 7.1 Standard sentence-level metrics show gains if the dataset is dense enough

Table 3 shows state-of-the-art performance for all models when translating at the sentence level (without context), compared to the commercial system. This confirms the large-scale, state-of-the-art nature of our experiments. On the WMT datasets, we see a fairly a regular small drop on sentence-level translation with SENT($\mathcal{P},\mathcal{B}$) (first row top sent-level section), that is slowly regained as we

move down to DOC($\mathcal{P},\ddot{\mathcal{B}}$). We note that we do not expect the contextual translation systems to perform *better* at sentence-level translation, but hope they retain performance there.

Next, Table 3 allows comparison of sentence-level translation to contextual translation (top versus bottom section). On the WMT datasets, the effects gains are fairly small (-0.3 for EN→DE, +0.6 for EN→RU). Looking at the CTXPro columns, however, we observe fairly large, consistent gains when translating contextually with nearly all the (non-randomized) DOC systems, but especially for the DOC($\mathcal{P},\ddot{\mathcal{B}}$) system across all three languages (+7.2 for EN→DE, +7.7 for EN→FR, and +7.3 for EN→RU). The CTXPro dataset is the OpenSubtitles gender-identified portion, so it is extremely dense in phenomena that require context to resolve compared to the WMT datasets, and is better able to discriminate systems with contextual abilities.

### 7.2 Domain and context both play a role

The DOC($\mathcal{P},\ddot{\mathcal{B}}$) system showed large gains in Table 3 when translating CTXPro contextually. One explanation is that CTXPro is, by construction, "discourse dense". But it also represents a domain shift, from news to conversational domains. We would like to have an idea of how much of the gain is due to each.

We therefore conduct a followup experiment in EN→DE that compares two datasets in the Open-Subtitles domain: the CTXPro/gender "dense" test

|  | EN→DE | | | EN→FR | | EN→RU | | | |
|---|---|---|---|---|---|---|---|---|---|
|  | AUX | FORm | GEN | FORm | GEN | AUX | FORm | GEN | INFl |
| #lines | 3,180 | 45,000 | 31,640 | 30,000 | 43,375 | 8,667 | 40,075 | 32,948 | 30,000 |
| SENT($\mathcal{P}$,$\mathcal{B}$) | 4.7 | 42.1 | 44.4 | 38.2 | 38.9 | 5.3 | 51.2 | 37.5 | |
| RAND($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | 4.7 | 39.6 | 42.4 | 36.9 | 38.2 | 5.5 | 51.4 | 36.7 | 32.2 |
| DOC($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | 4.9 | 41.7 | 50.7 | 38.7 | 47.6 | 20.9 | 58.6 | 45.5 | 39.8 |
| DOC($\ddot{\mathcal{P}}$,$\mathcal{B}$) | 4.2 | 41.4 | 47.2 | 42.7 | 45.2 | 16.7 | 56.8 | 39.5 | 37.4 |
| DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) | 7.5 | 45.0 | 66.0 | 43.8 | 54.8 | 25.2 | 58.7 | 53.5 | 42.6 |

Table 4: Generative accuracy on CTXPro datasets, where the task is to translate a source sentence and then determine whether an exact form of the required target word is in the output. The contextual systems trained on documents from mined parallel data perform notably worse than the DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) system.

| model | EN→DE gender | | EN→FR gender | | EN→RU NP ellipsis | | VP ellipsis | |
|---|---|---|---|---|---|---|---|---|
|  | contr. | gen. | contr. | gen. | contr. | gen. | contr. | gen. |
| RAND($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | 43.3 | 35.5 | 71.2 | 40.1 | 18.0 | 24.8 | 52.6 | 4.8 |
| DOC($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | 77.0 | 40.9 | 91.2 | 56.2 | 20.9 | 58.6 | 45.5 | 39.8 |
| DOC($\ddot{\mathcal{P}}$,$\mathcal{B}$) | 75.1 | 37.0 | 92.5 | 52.5 | 16.7 | 56.8 | 39.5 | 37.4 |
| DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) | **80.8** | **66.8** | **93.4** | **68.5** | **25.2** | **58.7** | **53.5** | **42.6** |

Table 5: Document contrastive test scores (contr.) and their generative (gen.) variants. All accuracies are over items with extra-sentential antecedents only. DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) consistently performs best on generative metrics by wide margins, while for contrastive metrics, other contextual systems are often similar or exhibit no consistent pattern.

| context | Dense true | Sparse true | Dense rand |
|---|---|---|---|
| SENT($\mathcal{P}$,$\mathcal{B}$) | 24.4 | 30.5 | 24.4 |
| DOC($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | 26.9 | 31.4 | 24.8 |
| DOC($\ddot{\mathcal{P}}$,$\mathcal{B}$) | 25.4 | 32.4 | 25.4 |
| DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) | 31.6 | 31.7 | 21.8 |

Table 6: EN→DE COMET scores on a dense dataset (OpenSubtitles CTXpro/gender) with true and random contexts; next, a sparse dataset (random sample of OpenSubtitles) with true contexts. DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) gains most over the sentence baseline on dense with true contexts and is harmed most on dense with random contexts. The doc systems are similar on the sparse dataset.

set, and another test set, which contains a random sample of 500 ten-sentence documents from Open-Subtitles 2016, yielding a corpus size of 4,973 sentences. We label this second one "sparse": since it was selected randomly, it is likely to be much less dense in contextual phenomena. For contextual systems, we translate each of these as a single chunk, and then split them out for evaluation with

COMET. The results are in Table 6.

The differences between the first two columns shows that the DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) gains over the sentence system are much larger on the "dense" dataset (+7.2 vs. +1.2). Performance among the contextual systems is closer, as we saw with WMT datasets. This suggests that the flat performance with WMT data was likely due to it, too, being sparse with contextual phenomena. **For standard, sentence-based metrics like COMET to separate these systems, dense test sets are needed.**

Table 6 (column 3) contains the results of another experiment, where we replace the context of each sentence in the "dense" dataset with a random context. This hurts performance, and the effect is most pronounced on the DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) system, suggesting that this model is most dependent on a reliable contextual clue.

### 7.3 Generative word-based accuracy corroborates these differences

Table 4 presents the results of word-based accuracy on the CTXPro datasets, across a range of linguistic phenomena. With word-based accuracy, we

are testing whether a word is present in the output. This leaves open the possibility of metric mistakes. For example, if the pronoun *er* is expected in the output, a system could be penalized for translating the sentence correctly with no pronouns, or it could be rewarded for generating a semantically unrelated instance of *er*. We do not expect this to systematically favor any one system.

Here, we see a similar gap between (a) contextual systems versus a random context and (b) especially, a gap between $\text{Doc}(\mathcal{P},\ddot{\mathcal{B}})$ and the other contextual systems. For EN→DE and EN→FR, the gender categories are similar to the ContraPro test sets for those languages, but much larger. This is most true for the GENder category (with gains of +23.6, +16.6, and +10.4), but also for other categories, including auxiliaries (+19.7 for EN→FR) and EN→RU inflection (+10.4).

## 7.4 The general trend favors BT-only contextual data

Figure 1 visualizes the metric score gains from Tables 3 and 5 for all four contextual models over the sentence-level baselines. The $x$-axis is arranged by the percentage of the contextual examples that are drawn from parallel data. This makes clearer the observations from the discussion so far: contextual annotations from parallel data are better than nothing, but they are inferior to those from the backtranslation monolingual data, and removing them is preferable.

## 7.5 Contrastive test sets are less discriminative

Table 5 contains results that pair contrastive accuracies (§ 4.2) with their generative counterparts. Across all three language pairs, there is an interesting pattern: in the contrastive metrics, the document systems improve over the sentence baseline, as a block. However, *the generative metrics see their best results with* $\text{Doc}(\mathcal{P},\ddot{\mathcal{B}})$, *often by a large margin*. Together with the observations in the previous section, we believe this calls into question the reliability of contrastive metrics. What we really care about in an MT system is its ability to *generate* the correct results at inference time. Discriminative ability is at best a proxy for this ability; if its results do not correlate with such metrics, it calls into question its reliability.



Figure 1: Contextual metric gains over the sentence baseline for COMET and accuracy metrics for the four systems, arranged by the percentage of contextual samples sourced from parallel data.

## 7.6 Experiments with public data provide some corroboration

Since complete document annotations for publicly available large-scale parallel data do not exist, we were unable to build $\text{Doc}(\ddot{\mathcal{P}},\ddot{\mathcal{B}})$ and $\text{Doc}(\ddot{\mathcal{P}},\mathcal{B})$ on open data. However, we can build the $\text{Sent}(\mathcal{P},\mathcal{B})$ and $\text{Doc}(\mathcal{P},\ddot{\mathcal{B}})$ systems with a subset of the WMT22 EN→DE data with monolingual document annotations, and see whether they exhibit the same pattern.

We use all available parallel data provided for WMT22 (Kocmi et al., 2022):[9] Europarl v10 (Koehn, 2005), Paracrawl v9 (Bañón et al., 2020), Common Crawl,[10] News Commentary, Wiki Ti-

---

[9] statmt.org/wmt22/translation-task.html
[10] https://commoncrawl.org/

| | |
|---|---|
| context | But let's not give in just yet.⟨SEP⟩ Right now, this is our one chance to be different.⟨SEP⟩ We could do something great with it.⟨SEP⟩ Like save the science museum.⟨SEP⟩ We grew up going to that place our whole lives.⟨SEP⟩ It's gave us so much.⟨SEP⟩ This is an opportunity to give something back.⟨SEP⟩ Besides, aren't you curious?⟨SEP⟩ So, three wishes are granted to whoever discovers <u>the box</u>. |
| source | But we all found **it**. And touched it at the same time. |
| SENT($\mathcal{P}$,$\mathcal{B}$) | Aber wir haben **es** alle gefunden. Und **es** gleichzeitig berührt. |
| RAND($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | Aber wir haben **es** alle gefunden und gleichzeitig berührt. |
| DOC($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | Aber wir haben **es** alle gefunden. Und haben **es** gleichzeitig berührt. |
| DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) | Aber wir haben **sie** alle gefunden und gleichzeitig angefasst. |
| ref | Aber wir haben **sie** alle gleichzeitig entdeckt und berührt. |
| context | Mark it.⟨SEP⟩ If Mr. Wick isn't dead already, he soon will be.⟨SEP⟩ Will you mark it, sir?⟨SEP⟩ You have no idea, what's coming do you?⟨SEP⟩ I have everyone in New York looking for him.⟨SEP⟩ I doubt we will see him again.⟨SEP⟩ Do you now?⟨SEP⟩ You stabbed the devil in the back, and forced him back into the life that he had just left.⟨SEP⟩ You incinerated the priest's <u>temple</u>. |
| source | Burned **it** to the ground. |
| SENT($\mathcal{P}$,$\mathcal{B}$) | Verbrannte **es** bis auf die Grundmauern. |
| RAND($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | Verbrannte **es** zu Boden. |
| DOC($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | Hast **es** zu Boden gebrannt. |
| DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) | Sie haben **ihn** niedergebrannt. |
| ref | Und **ihn** niedergebrannt. |

Table 7: Translation examples from the CTXPro gender dataset demonstrating DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$)'s superior performance. Pronouns are **in bold** with antecedents <u>underlined</u>. For all but SENT($\mathcal{P}$,$\mathcal{B}$), the source is translated together with the context, and then the context is discarded.

.

tles v3, Tilde MODEL Corpus (Rozis and Skadiņš, 2017), and Wikimatrix (Schwenk et al., 2021a). A few of these resources have document-level information, but we do not use any of it. For monolingual data, the only data available with document metadata is News Crawl.[11] We used all even years from 2008–2020, backtranslating it from German to English with an internal system. No filtering is applied. From this data, we train the only two of our systems supported by this setup: SENT($\mathcal{P}$,$\mathcal{B}$) and DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$). These are trained for 40 virtual epochs each using the same settings described in Section 6.[12]

Results can be found in Table 8. They are encouraging: we see the same pattern of improvement between SENT($\mathcal{P}$,$\mathcal{B}$) and DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$), although the absolute numbers are lower. Compared to our in-house data, the document metrics are even better for SENT($\mathcal{P}$,$\mathcal{B}$).

| | | gender | |
|---|---|---|---|
| system | COMET | contr. | gen. |
| SENT($\mathcal{P}$,$\mathcal{B}$) | 60.6 | 56.7 | 23.9 |
| DOC($\ddot{\mathcal{P}}$,$\ddot{\mathcal{B}}$) | x | x | x |
| DOC($\ddot{\mathcal{P}}$,$\mathcal{B}$) | x | x | x |
| DOC($\mathcal{P}$,$\ddot{\mathcal{B}}$) | 59.4 | 83.4 | 64.3 |

Table 8: Metrics on the only two models we are able to build on public data. Similar patterns are observable to those seen in Tables 3 and 5.

### 7.7 MT output in crawled parallel data

We do not undertake an exploration of the causes for the results and analysis discussed in Figure 1 and throughout this section, but there is an obvious explanation: we suspect that parallel web-crawled data is full of machine-translated output. Widespread use of translation across the web, especially since the release of Google Translate in 2006, is a commercial success story that has unfortunately produced a kind of "poisoning of the

---

[11] https://data.statmt.org/news-crawl/de-doc/
[12] Mono data: 311.2m lines, 14.1m docs, with a mean sentence length of 21.9 sentences. Parallel data: 297.6m lines.

| English | German |
|---------|--------|
| Unique Moorish style **villa** set in a tropical oasis with pool, guest accommodation and amazing views. ⟨SEP⟩ Property Reference 1846 ⟨SEP⟩ **It** was built by the current owner... | Einzigartige maurische **Villa** in einer tropischen Oase mit Pool, Gästeunterkunft und herrlicher Aussicht. ⟨SEP⟩ Referenznummer 1846 ⟨SEP⟩ <span style="color:red">**Es**</span> wurde vom jetzigen Besitzer gebaut... |

Table 9: An example of bad data drawn from the parallel data pool. While the sentence-level translations are fine, the incorrect pronoun *Es* in the third sentence suggests sentence-level machine or low-quality human translations.

well", where machine translation outputs are later collected as training data for new systems (Venugopal et al., 2011). Recent work has corroborated how extensive this is in multi-way parallel data (Thompson et al., 2024).

Quantifying this awaits further work, but it is easy to source examples from our parallel data (Table 9). While we don't know if this was generated by machine or a human, we do know that even large NMT systems are sensitive to small amounts of poor data.[13]. This data may still be of high quality at the sentence level; it is only *inter-sentence contextual* information that is affected. If true, this suggests that **contextual translation introduces a new quality dimension that is invisible** in the standard sentence-level training paradigm, and the problem may in fact be quite large, since all machine translation content in the wild will have been generated at the sentence level.

We suspect that our monolingual data—which by design was sourced from known target-native sites, such as newspapers—is largely immune from these problems. Training on sentence-level translations is primarily a problem for data translated in the *forward* direction. Backtranslation introduces noise into the *source language text*, while preserving the target-language contextual signal.

We leave to future work an investigation into detecting and removing machine translation output from parallel data at high enough precision.

## 8 Conclusions

Machine translation research and production systems continue to be dominated by sentence-level approaches. A common explanation for this shortcoming is the lack of document-annotated parallel data. We have compared the effectiveness of constructing contextual translation models for three translation directions in large-data settings. Our results suggests that while mined parallel data is of high-enough quality for building sentence systems and contains some contextual signal, **it is best to construct contextual training samples from back-translated data only**. Although we have not investigated the reasons for this, we consider it a strong possibility that our parallel data, which is mostly crawled from the web and has had only sentence-level filtering applied, contains large amounts of data that was machine-translated at the sentence level, a finding that is very likely to hold for publicly available data, as well. This suspicion makes sense a priori, and is bolstered in other recent work (Thompson et al., 2024; Wicks et al., 2024; Pal et al., 2024).

We have also shown the importance of **evaluating contextual machine translation output in its generative capacity**, rather than in its ability to discriminate good outputs from bad ones. This can be done by using provided challenge sets like CTX-Pro or converting existing contrastive metrics like ContraPro and its variants, or by using standard corpus-level metrics like COMET on test sets that are sufficiently dense with contextual phenomena.

A fruitful avenue for followup work is to automatically identify sentences that require context to translate correctly, which could be used to filter training data and also in the construction of new test sets. Though we have focused on "traditionally"-trained MT, it will also be useful to learn how LLMs perform on these tasks.

---

[13]A classic example is source-copy data (Ott et al., 2018)

## Limitations

With respect to reproducibility, the deepest limitation of our paper is our use of private data. There is therefore a risk that our findings might not be reproducible by other teams working with (necessarily) different datasets. Finally, although we suspect our results will hold for language pairs beyond the three we investigated, it is possible they will not generalize.

## References

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

Guangsheng Bao, Yue Zhang, Zhiyang Teng, Boxing Chen, and Weihua Luo. 2021. G-transformer for document-level machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3442–3455, Online. Association for Computational Linguistics.

Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. Findings of the 2019 conference on machine translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy. Association for Computational Linguistics.

Rachel Bawden, Thomas Lavergne, and Sophie Rosset. 2018. Detecting context-dependent sentences in parallel corpora. In *Actes de la Conférence TALN. Volume 1 - Articles longs, articles courts de TALN*, pages 393–400, Rennes, France. ATALA.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Barry Haddow, Matthias Huck, Chris Hokamp, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Carolina Scarton, Lucia Specia, and Marco Turchi. 2015. Findings of the 2015 workshop on statistical machine translation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 1–46, Lisbon, Portugal. Association for Computational Linguistics.

Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. WIT3: Web inventory of transcribed and translated talks. In *Proceedings of the 16th Annual Conference of the European Association for Machine Translation*, pages 261–268, Trento, Italy. European Association for Machine Translation.

Patrick Fernandes, Kayo Yin, Emmy Liu, André Martins, and Graham Neubig. 2023. When does translation require context? a data-driven, multilingual exploration. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 606–626, Toronto, Canada. Association for Computational Linguistics.

Liane Guillou and Christian Hardmeier. 2016. PROTEST: A test suite for evaluating pronouns in machine translation. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 636–643, Portorož, Slovenia. European Language Resources Association (ELRA).

Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. 2019. Star-transformer. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1315–1325, Minneapolis, Minnesota. Association for Computational Linguistics.

Christian Herold and Hermann Ney. 2023. On search strategies for document-level neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12827–12836, Toronto, Canada. Association for Computational Linguistics.

Yuchen Jiang, Tianyu Liu, Shuming Ma, Dongdong Zhang, Jian Yang, Haoyang Huang, Rico Sennrich, Ryan Cotterell, Mrinmaya Sachan, and Ming Zhou. 2022. BlonDe: An automatic evaluation metric for document-level machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1550–1565, Seattle, United States. Association for Computational Linguistics.

Marcin Junczys-Dowmunt. 2019. Microsoft translator at WMT 2019: Towards large-scale document-level neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 225–233, Florence, Italy. Association for Computational Linguistics.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018a. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Marcin Junczys-Dowmunt, Kenneth Heafield, Hieu Hoang, Roman Grundkiewicz, and Anthony Aue. 2018b. Marian: Cost-effective high-quality neural machine translation in C++. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 129–135, Melbourne, Australia. Association for Computational Linguistics.

Jungo Kasai, Nikolaos Pappas, Hao Peng, James Cross, and Noah A. Smith. 2020. Deep encoder, shallow decoder: Reevaluating the speed-quality tradeoff in machine translation. *CoRR*, abs/2006.10369.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand.

Shaohui Kuang, Deyi Xiong, Weihua Luo, and Guodong Zhou. 2018. Modeling coherence for neural machine translation with dynamic and topic caches. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 596–606, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics.

Samuel Läubli, Rico Sennrich, and Martin Volk. 2018. Has machine translation achieved human parity? a case for document-level evaluation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4791–4796, Brussels, Belgium. Association for Computational Linguistics.

Pierre Lison and Jörg Tiedemann. 2016. OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 923–929, Portorož, Slovenia. European Language Resources Association (ELRA).

Siyou Liu and Xiaojun Zhang. 2020. Corpora for document-level neural machine translation. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3775–3781, Marseille, France. European Language Resources Association.

António Lopes, M. Amin Farajian, Rachel Bawden, Michael Zhang, and André F. T. Martins. 2020. Document-level neural MT: A systematic comparison. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 225–234, Lisboa, Portugal. European Association for Machine Translation.

Lorenzo Lupo, Marco Dinarelli, and Laurent Besacier. 2023. Encoding sentence position in context-aware neural machine translation with concatenation. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 33–44, Dubrovnik, Croatia. Association for Computational Linguistics.

Sameen Maruf, Fahimeh Saleh, and Gholamreza Haffari. 2019. A survey on document-level machine translation: Methods and evaluation. *CoRR*, abs/1912.08494.

Lesly Miculicich, Dhananjay Ram, Nikolaos Pappas, and James Henderson. 2018. Document-level neural machine translation with hierarchical attention networks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2947–2954, Brussels, Belgium. Association for Computational Linguistics.

Mathias Müller, Annette Rios, Elena Voita, and Rico Sennrich. 2018. A large-scale test set for the evaluation of context-aware pronoun translation in neural machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 61–72, Brussels, Belgium. Association for Computational Linguistics.

Myle Ott, Michael Auli, David Grangier, and Marc'Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. *CoRR*, abs/1803.00047.

Proyag Pal, Alexandra Birch, and Kenneth Heafield. 2024. Document-level machine translation with large-scale public parallel corpora. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13185–13197, Bangkok, Thailand. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post, Thamme Gowda, Roman Grundkiewicz, Huda Khayrallah, Rohit Jain, and Marcin Junczys-Dowmunt. 2023. SOTASTREAM: A streaming approach to machine translation training. In *Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023)*, pages 110–119, Singapore. Association for Computational Linguistics.

Gema Ramírez-Sánchez, Jaume Zaragoza-Bernabeu, Marta Bañón, and Sergio Ortiz Rojas. 2020. Bifixer and bicleaner: two open-source tools to clean your parallel data. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 291–298, Lisboa, Portugal. European Association for Machine Translation.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Roberts Rozis and Raivis Skadiņš. 2017. Tilde MODEL - multilingual open data for EU languages. In *Proceedings of the 21st Nordic Conference on Computational Linguistics*, pages 263–265, Gothenburg, Sweden. Association for Computational Linguistics.

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021a. WikiMatrix: Mining 135M parallel sentences in 1620 language pairs from Wikipedia. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361, Online. Association for Computational Linguistics.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021b. CCMatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500, Online. Association for Computational Linguistics.

Rico Sennrich. 2017. How grammatical is character-level neural machine translation? assessing MT quality with contrastive translation pairs. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 376–382, Valencia, Spain. Association for Computational Linguistics.

Rico Sennrich. 2018. Why the time is ripe for discourse in machine translation. Talk given at NGT 2018: https://aclanthology.org/volumes/W18-27/.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Amane Sugiyama and Naoki Yoshinaga. 2019. Data augmentation using back-translation for context-aware neural machine translation. In *Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019)*, pages 35–44, Hong Kong, China. Association for Computational Linguistics.

Zewei Sun, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Lei Li. 2022. Rethinking document-level neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3537–3548, Dublin, Ireland. Association for Computational Linguistics.

Brian Thompson, Mehak Dhaliwal, Peter Frisch, Tobias Domhan, and Marcello Federico. 2024. A shocking amount of the web is machine translated: Insights from multi-way parallelism. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 1763–1775, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Jörg Tiedemann and Yves Scherrer. 2017. Neural machine translation with extended context. In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 82–92, Copenhagen, Denmark. Association for Computational Linguistics.

Zhaopeng Tu, Yang Liu, Shuming Shi, and Tong Zhang. 2018. Learning to remember translation history with a continuous cache. *Transactions of the Association for Computational Linguistics*, 6:407–420.

Jannis Vamvas and Rico Sennrich. 2021. On the limits of minimal pairs in contrastive evaluation. In *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 58–68, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *CoRR*, abs/1706.03762.

Ashish Venugopal, Jakob Uszkoreit, David Talbot, Franz Och, and Juri Ganitkevitch. 2011. Watermarking the outputs of structured prediction with an application in statistical machine translation. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 1363–1372, Edinburgh, Scotland, UK. Association for Computational Linguistics.

Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 118–128, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Elena Voita, Rico Sennrich, and Ivan Titov. 2019a. Context-aware monolingual repair for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 877–886, Hong Kong, China. Association for Computational Linguistics.

Elena Voita, Rico Sennrich, and Ivan Titov. 2019b. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1198–1212, Florence, Italy. Association for Computational Linguistics.

Elena Voita, Pavel Serdyukov, Rico Sennrich, and Ivan Titov. 2018. Context-aware neural machine translation learns anaphora resolution. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1264–1274, Melbourne, Australia. Association for Computational Linguistics.

Rachel Wicks and Matt Post. 2023. Identifying context-dependent translations for evaluation set production. In *Proceedings of the Eighth Conference on Machine Translation*, pages 452–467, Singapore. Association for Computational Linguistics.

Rachel Wicks, Matt Post, and Philipp Koehn. 2024. Recovering document annotations for sentence-level bitext. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 9876–9890, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Minghao Wu, George Foster, Lizhen Qu, and Gholamreza Haffari. 2023. Document flattening: Beyond concatenating context for document-level neural machine translation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 448–462, Dubrovnik, Croatia. Association for Computational Linguistics.

Lei Yu, Laurent Sartran, Wojciech Stokowiec, Wang Ling, Lingpeng Kong, Phil Blunsom, and Chris Dyer. 2020. Better document-level machine translation with Bayes' rule. *Transactions of the Association for Computational Linguistics*, 8:346–360.

Jiacheng Zhang, Huanbo Luan, Maosong Sun, Feifei Zhai, Jingfang Xu, Min Zhang, and Yang Liu. 2018. Improving the transformer translation model with document-level context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 533–542, Brussels, Belgium. Association for Computational Linguistics.

## A    Dataset examples

Examples from the datasets used for generative and contrastive evaluation can be found in Tables 10 and 11.

## B    Model capacity

Much work in investigating document-level machine translation has been limited to standard-size Transformer architectures (cf. Zhang et al. (2018); Sun et al. (2022); Lopes et al. (2020)). Yet it stands

---

The prototype has passed every test, sir. It's working. | Der Prototyp hat jeden Test erfolgreich durchlaufen, Sir. {Er,Es,Sie} funktioniert.

(a) ContraPro example. Contrastive examples are formed by substituting incorrect pronouns.

---

Veronica, thank you, but you **saw** what happened. We all did. | Вероника, спасибо, но ты видела, что произошло. Мы все **хотели**.

(b) GTWiC example. The first Russian sentence uses the formal register.

Table 10: Examples from contrastive test sets.

---

(AUX ) I just figured you need to know. And now you do. → Je pensais que tu méritais de savoir. Et maintenant tu *sais*.

---

(INF) My friend had some mech work done here. Industry stuff. → Вы ставили имплант моей подруге. Промышленную штуковину.

---

(FORm) I don't know you, but.. → Ich kenne Sie nicht, aber...

Table 11: Examples of contextually-sensitive auxiliary and inflection elision from the CTXPro dataset.

---

to reason that modeling longer-range phenomena will require increased model capacity, and in fact, the base model size we chose for our experiments (12 layer encoder, 16k FFN) reflects this. Here, we provide more detail, varying two model parameters only: (i) the number of encoder layers, and (ii) the width of the model feed-forward layer (encoder and decoder side). We keep all other parameters the same, including fixing the decoder depth to 6. Focusing on changes to the encoder depth helps limit grid search and is justified by prior work showing that (relatively cheap) encoder layers can be traded for (relatively expensive) decoder layers with no penalty (Kasai et al., 2020). We alternate between increasing the number of encoding layers, and increasing the dimension of the Transformer feed-forward layer.

Table 12 contains English–German results. Unsurprisingly, all scores continue to rise, up to the wide 18-layer model. Both increasing the number of encoder layers, and increasing the size of the FFN, contribute to better performance. This suggests that the common approach of working with 6-layer Transformer base models is not enough

| arch | params | BLEU | COMET | C/Pro | G/Pro |
|------|--------|------|-------|-------|-------|
| 6/1k | 146m | 27.0 | 48.7 | 65.2 | 58.4 |
| 6/2k | 171m | 27.4 | 49.7 | 66.2 | 58.7 |
| 6/4k | 221m | 28.0 | 51.0 | 69.7 | 62.9 |
| 12/4k | 297m | 28.4 | 51.8 | 70.6 | 66.0 |
| 6/8k | 322m | 27.8 | 51.0 | 71.7 | 62.8 |
| 12/8k | 448m | 28.6 | 52.5 | 74.2 | 67.1 |
| 6/16k | 523m | 28.4 | 51.7 | 74.5 | 64.9 |
| 18/8k | 574m | 28.8 | 53.0 | 75.0 | 67.1 |
| 12/16k | 750m | 28.9 | 52.8 | 75.8 | 68.5 |
| 18/16k | 977m | 29.3 | 53.3 | 75.5 | 69.4 |

Table 12: Model capacity (encoder layers / FFN / # params) for an EN-DE document model, ordered by param. count. Decoder depth is always 6 layers. Scores were computed on a checkpoint after 30k updates. BLEU and COMET scores are on WMT21, translating as sentences. C/Pro is over the complete test set, while G/Pro is over only sentences with external anaphora.

for document-context MT. There is more to gain by moving to larger models and likely, to larger datasets and context lengths, as well.

# A Multi-task Learning Framework for Evaluating Machine Translation of Emotion-loaded User-generated Content

**Shenbin Qian**[⊕], **Constantin Orăsan**[⊕], **Diptesh Kanojia**[✦] **and Félix do Carmo**[⊕]

[⊕]Centre for Translation Studies and [✦]Institute for People-Centred AI,
University of Surrey, United Kingdom
{s.qian, c.orasan, d.kanojia, f.docarmo}@surrey.ac.uk

## Abstract

Machine translation (MT) of user-generated content (UGC) poses unique challenges, including handling slang, emotion, and literary devices like irony and sarcasm. Evaluating the quality of these translations is challenging as current metrics do not focus on these ubiquitous features of UGC. To address this issue, we utilize an existing emotion-related dataset that includes emotion labels and human-annotated translation errors based on Multi-dimensional Quality Metrics. We extend it with sentence-level evaluation scores and word-level labels, leading to a dataset suitable for sentence- and word-level translation evaluation and emotion classification, in a multi-task setting. We propose a new architecture to perform these tasks concurrently, with a novel combined loss function, which integrates different loss heuristics, like the Nash and Aligned losses. Our evaluation compares existing fine-tuning and multi-task learning approaches, assessing generalization with ablative experiments over multiple datasets. Our approach achieves state-of-the-art performance and we present a comprehensive analysis for MT evaluation of UGC.

## 1 Introduction

Machine translation (MT) has advanced rapidly in recent years, leading to claims it has achieved human parity in Chinese-English news translation (Hassan et al., 2018). Recent advent of large language models (LLMs) has determined researchers to repeat claims of human parity more often (Wang et al., 2021). However, automatically translating user-generated content (UGC) with expressions that contain emotions, like tweets, reveals novel challenges for MT systems (Saadany et al., 2023). Figure 1 shows the output of Google Translate (GT) and ChatGPT when we translated some Chinese UGC with emotional slang using them[1].

---

[1]GPT-3.5 at "https://chat.openai.com/" in April, 2024

As can be seen from the example, both outputs need to be improved significantly to be considered usable. Similar problems were noticed with other MT engines, indicating that it is imperative to evaluate MT quality with metrics that take emotion preservation into account.

Using human judgements/input to evaluate MT quality is expensive in terms of both time and money (Dorr et al., 2011; Lai et al., 2020). Quality estimation (QE), which predicts MT quality in the absence of human references, can serve as a cost-effective alternative to approximate human evaluation based on metrics like Multi-dimensional Quality Metrics (MQM), an error-based human evaluation scheme for MT quality (Lommel et al., 2014). A widely-used approach in QE involves fine-tuning a multilingual pre-trained language model (PTLM) using human evaluation data (Blain et al., 2023). This fine-tuned model can predict scores for entire MT sentences or labels for individual words, indicating whether each word is correctly translated or not. This encompasses two common QE tasks: sentence-level QE and word-level QE.

To assess MT quality of emotion-loaded UGC, it is crucial to evaluate the overall quality of emotion preservation after translation (sentence-level QE), and how well emotion words are translated (word-level QE). To achieve this, we leverage an existing emotion-related dataset that includes emotion labels and MQM-based human-evaluated translation errors. We extend it with sentence-level QE scores and word-level labels, resulting in a dataset extension. This extended dataset is suitable for both sentence- and word-level QE, and emotion classification. For joint training of these tasks, we employ multi-task learning (MTL), anticipating improved performance for all tasks due to their inherent correlation with emotionally charged content. We further introduce a new architecture with a novel combined loss function that integrates different loss heuristics, enabling the concurrent training

**Human Translation:** Countless "f**k your mother" appeared in my mind!

**Explanation:** Both Google Translate and ChatGPT fail to translate the swear word "草泥马", a slang word created using a homophone to replace the original character to avoid censorship. The angry emotion of the original sentence is completely lost.

Figure 1: Example of translations from Google Translate and ChatGPT

of these tasks and optimizing their overall performance. We compare our MTL approach with existing fine-tuning and MTL methods. Our proposed approach achieves new state-of-the-art results on the emotion-related QE dataset and a standard QE dataset. Our contributions can be summarized as follows:

- *Extending an emotion-related QE dataset* with 1) QE scores at sentence level and 2) labels indicating emotion-related translation quality at word level.

- A new architecture with a *novel combined loss function*, integrating different loss heuristics for multi-task learning[2].

- Evaluation of the proposed MTL approach on multiple QE datasets including ablative experiments on combinations of QE and emotion classification tasks, *improving performance over existing fine-tuning and MTL methods*.

Section 2 discusses existing work for QE and MTL while Section 3 introduces the datasets we use for this study. Our approach, baselines and experimental setup are described in Section 4, and Section 5 discusses the results obtained on multiple datasets. Section 6 concludes our study and outlines future directions. Section 7 points out limitations and ethical considerations. Relevant mathematical equations and loss algorithms are explained in Appendix A.

---

[2]Our code and the extended dataset for MTL are available at https://github.com/shenbinqian/MTL4QE.

## 2 Related Work

We discuss related work in supervised QE in § 2.1. Studies on MTL and its application to QE are reviewed in § 2.2.

### 2.1 Quality Estimation

Though prompting with LLMs is increasingly applied to the field of quality evaluation (Kocmi and Federmann, 2023b,a; Fernandes et al., 2023), supervised fine-tuning of multilingual PTLMs on human evaluation data based on metrics such as translation edit rate (Snover et al., 2006), direct assessment (Graham et al., 2013) and MQM, remains as state-of-the-art QE methods (Kocmi and Federmann, 2023b). TransQuest (Ranasinghe et al., 2020) and COMET (Rei et al., 2020; Stewart et al., 2020; Rei et al., 2022b; Guerreiro et al., 2024) are two popular frameworks used for sentence-level QE. TransQuest utilizes XLM-RoBERTa (Conneau et al., 2020) as the backbone, concatenating the source and target sentences using [CLS] (start) and [SEP] (separator) tokens. In MonoTransQuest, an architecture within TransQuest, only the embeddings of the [CLS] token are used for prediction. In SiameseTransQuest, a variant of TransQuest architecture, a twin XLM-RoBERTa network computed the mean of all token embeddings for the source and target. This mean is then used to calculate the cosine similarity as the final QE score. Unlike TransQuest, COMET was initially proposed for reference-based evaluation until 2022, when COMETKIWI (Rei et al., 2022b) was introduced to support reference-less evaluation. Similar to

1141

MonoTransQuest, it concatenates the source and target, and inputs them into the encoder. All hidden states are then fed into a scalar mix module (Peters et al., 2018) that learns a weighted sum, producing a new sequence of aggregated hidden states. The output of the [CLS] token is then used for the prediction of sentence-level QE scores.

For word-level QE, OpenKiwi (Kepler et al., 2019) was proposed to support both sentence- and word-level QE with various neural network architectures. MicroTransQuest (Ranasinghe et al., 2021), utilizing outputs of all input tokens of an XLM-RoBERTa model based on the MonoTransQuest architecture, was proposed only for word-level QE under multilingual settings.

Because of their successes in the QE shared tasks in the Conference on Machine Translation (WMT) in recent years (Specia et al., 2020, 2021; Zerva et al., 2022), TransQuest and COMET are selected as our baseline fine-tuning frameworks for sentence-level QE, and MicroTransQuest for word-level QE.

## 2.2 Multi-task Learning

Multi-task learning addresses multiple related tasks concurrently by training them simultaneously with a shared representation (Caruana, 1997). While this approach reduces the training cost compared to training separate models (Baxter, 2000), early methods led to performance degradation when compared to single-task models (Standley et al., 2020). Recent efforts have introduced various methods to address this problem and enhance the MTL performance.

Liu et al. (2019) proposed dynamic weight averaging that could learn task-specific feature-level attention. They used a shared network that contains features of all tasks and a soft-attention module for each specific task without using weighting schemes. Liu et al. (2021) proposed impartial MTL that uses different strategies for task-shared parameters and task-specific parameters. Navon et al. (2022) proposed to view the combination of gradients as a bargaining game, where different tasks negotiate with each other to reach an agreement on a joint direction of parameter update. They utilized the Nash Bargaining Solution (Nash, 1953) as an approach to address this problem and proved the effectiveness of their method across various tasks. Since some MTL methods are not always stable during training, Senushkin et al. (2023) proposed the Aligned MTL to improve stability. They used a condition number of a linear system of gradients as a stability criterion, and aligned the orthogonal components of the linear system of gradients to eliminate instability in training.

The improved performance and stability of MTL methods have prompted its application to quality evaluation. Shah and Specia (2016) investigated MTL with Gaussian Processes for QE, based on datasets with multiple annotators and language pairs. They found multi-task models perform better than individual models in cross-lingual settings. Zhang and van Genabith (2020) used MTL to predict QE scores and rank different translations. Rei et al. (2022a) employed MTL to jointly train QE models at sentence- and word-level. Most of these studies used non-parametric linear combinations of task losses, until Deoghare et al. (2023) proposed to apply Nash MTL to combining sentence- and word-level QE, based on MicroTransQuest. However, their Nash MTL might not always be stable for various QE tasks. In this paper, we explore different MTL loss heuristics and propose a new architecture with a novel combined loss function for the quality estimation of emotion-loaded UGC.

## 3 Data

We used two datasets to evaluate our approach. The first one measures *how well emotion is preserved* in machine translation and is presented in § 3.1. The second is a standard QE dataset from WMT 2020 to WMT 2022 (Freitag et al., 2021a,b, 2022). It contains sentence- and word-level QE data annotated using MQM, as explained in § 3.2.

### 3.1 A Human Annotated Dataset for Quality Assessment of Emotion Translation

We used our Human Annotated Dataset for Quality Assessment of Emotion Translation (HADQAET)[3] as the main resource (Qian et al., 2023). Its source text originated from the dataset released by the *Evaluation of Weibo Emotion Classification Technology on the Ninth China National Conference on Social Media Processing* (SMP2020-EWECT). It originally has a size of 34,768 instances. Each instance is a tweet-like text segment[4], which was manually annotated in the original dataset with one of the six emotion labels, *i.e.*, *anger*, *joy*, *sadness*, *surprise*, *fear* and *neutral*

---

[3] https://github.com/surrey-nlp/HADQAET
[4] Like most NLP tasks, we treat tweet-like text segments as sentence-level data. However, in contrast to tweets, our instances are longer with an average of 40 Chinese characters.

(Guo et al., 2021). We kept 5,538 instances with labels other than *neutral* and used Google Translate to translate them to English. We proposed an emotion-related MQM framework and recruited two professional translators to annotate errors and their corresponding severity in terms of emotion preservation. Words/characters in both source and target that cause errors were highlighted for error analysis. Details of our framework, error annotation (including inter-annotator agreement) and error analysis can be found in Qian et al. (2023). An example of the dataset is shown in Figure A.1.

Since our original paper did not propose any scores for sentence-level QE, we followed Freitag et al. (2021a) to sum up all weighted errors based on their corresponding severity, using a set of weights[5] suggested by MQM (Lommel et al., 2014), *i.e.*, 1 for minor errors, 5 for major and 10 for critical. For word-level QE, we first tokenized the source with *jieba* (Sun, 2013), and the target with NLTK (Bird et al., 2009) (tokenization tools for Chinese and English respectively). Then, we labeled the tokens highlighted by annotators as "BAD", and the rest "OK". This led to a sequence of labels for each instance, which indicate translation quality in emotion preservation at word level.

The MQM-based QE scores related to emotion, word labels, together with the source texts and GT translations were used for quality estimation of emotion-loaded UGC. The emotion labels that were originally used for emotion classification were also incorporated to see if they are helpful for QE.

## 3.2 MQM Subset with Synthetic Emotion

To test whether our approach can be applied to standard QE data[6], we selected the overlapping of Chinese-English sentence- and word-level MQM datasets from the QE shared task of WMT 2020 to WMT 2022. The overlapped subset has MQM scores at sentence level and "OK" or "BAD" labels at word level. We fine-tuned the Chinese RoBERTa large model (Cui et al., 2020) on the SMP2020-EWECT dataset, resulting in an emotion classifier with a macro F1 score of 0.95. We predicted the emotion label of the source text of the selected data using the fine-tuned classifier, and filtered out all *neutral* instances. This led to an MQM subset with automatically generated emotion labels and a comparable size (3544) as HADQAET. We randomly sampled 100 instances and manually

---
[5]We validated these weights in Qian et al. (2024).

[6]Their QE scores are not related to emotion.

---

checked the predicted emotion labels with the help of a native speaker. The manual validation shows the emotion classifier is reliable as it achieves an F1 score 0.90, precision 0.91 and recall 0.92. The distribution of this subset is shown in Figure 2.



Figure 2: Distribution of the MQM emotion subset

## 4 Methodology

This section describes the architecture and loss function of our MTL method. Additionally, it also presents the fine-tuning baselines including TransQuest and COMET for each individual task.

### 4.1 Multi-task Learning

We propose a new architecture that is able to train sentence- and word-level QE systems with an emotion classifier using a combined loss function.

**Architecture** The architecture we propose is in Figure 3. Following MonoTransQuest and COMETKIWI, we concatenate the source and target, including [CLS] and [SEP] as the starting and separating tokens. Then, we employed multilingual PTLMs like XLM-RoBERTa, XLM-V-base and InfoXLM (Chi et al., 2021) to encode the input text. Different from Deoghare et al. (2023), who used embeddings of the last hidden layer, we utilized the output of the [CLS] token to predict sentence-level QE scores and the rest tokens for word label classification. On top of the encoder, we added a fully connected layer for both sentence- and word-level QE before the softmax function for prediction.

To incorporate the emotion classification task, we tried max and average pooling for the output of the last hidden layer of the encoder and added another fully connected layer on top. We used Xavier initialization (Glorot and Bengio, 2010) for

the weights in all newly-added linear layers. We experimented different combination strategies for the losses of these tasks as explained below.
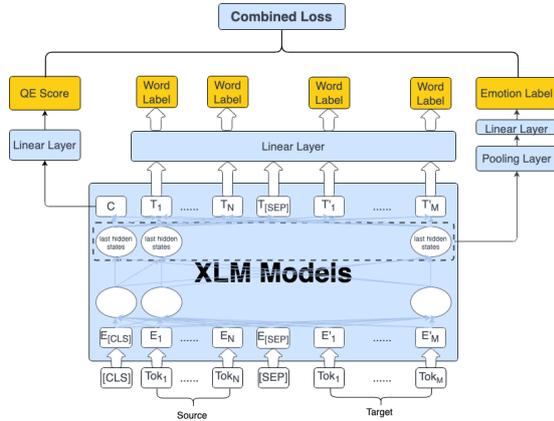


Figure 3: Architecture of our MTL Framework

**Combined Loss**   The loss function of our method is defined in Equation 1, where $\sigma$ is a heuristic function to combine the three losses. $L_{sent}$ as shown in Equation 2 is the Mean Squared Error loss for sentence-level QE, where $Y_{QE\_score}$ and $\hat{Y}_{QE\_score}$ are the true and predicted QE scores, respectively. Equation 3 is the Cross Entropy loss for word and emotion classification, where $C$ is the set of classes. For $L_{word}$, $C$ is {"OK", "BAD"}. For $L_{emo}$, $C$ is the 5 emotion classes. $\mathbb{1}\{y = i\}$ is an indicator function (1 if the true label $y$ is equal to the current class $i$, 0 otherwise), and $p_i$ is the predicted probability of the input being in class $i$.

$$L_{MTL} = \sigma(L_{sent}, L_{word}, L_{emo}) \quad (1)$$

$$L_{sent} = MSE(Y_{QE\_score}, \hat{Y}_{QE\_score}) \quad (2)$$

$$L_{word/emo} = -\sum_{i=1}^{C} \mathbb{1}\{y = i\} \cdot \log(p_i) \quad (3)$$

The objective of the heuristic $\sigma$ is to find a set of parameters $\theta$ that minimize the aggregate loss of all tasks. It is defined in Equation 4, where $L_{MTL}(\theta)$ is the combined loss, and $L_i(\theta)$ is the loss for an individual task $i$.

$$\theta^* = arg\min_{\theta}\{L_{MTL}(\theta) = \Sigma_{i=1}^{T} L_i(\theta)\} \quad (4)$$

Theoretically, $\theta$ can be fixed or a simple linear combination of each task loss. For instance, it can be 1 for each task loss, but the result is usually not ideal, as shown in our experiments. In order to balance the losses of different tasks and overcome

optimization problems like conflicting or dominating gradients (Navon et al., 2022), we adopted different heuristics $\sigma$ to learn $\theta$, including the Nash and Aligned MTL losses which are explained in Appendix A. Other existing MTL methods such as linear combination, dynamic weight averaging and impartial MTL were also integrated into our framework. To compare with our proposed Nash and Aligned MTL, the linear combination (1 for each task loss) and Nash MTL loss in Deoghare et al. (2023) were selected as baseline MTL methods in our experiments. Results of other MTL methods are in Table A.1.

### 4.2 Fine-tuning

We utilized MonoTransQuest, SiameseTransQuest and COMET for sentence-level QE, and Micro-TransQuest for word-level QE. They rely on the XLM-RoBERTa models as the foundation model for fine-tuning. For emotion classification, we fine-tuned XLM-RoBERTa-large and XLM-V-base (Liang et al., 2023) using both source and target texts. Experimental setup and training details can be seen in the following sections.

### 4.3 Experimental Setup

We performed experiments under two settings (fine-tuning and MTL) on two datasets (HADQAET and the MQM emotion subset). Fine-tuning included sentence- and word-level QE and emotion classification. For MTL, we combined sentence-level QE with word-level QE, sentence-level QE with emotion classification, and sentence-, word-level QE and emotion classification.

We used Spearman $\rho$ and Pearson's $r$ correlations to evaluate similarities between the predicted sentence-level QE scores and the true scores. For word and emotion classification, we used macro F1, precision and recall scores for evaluation.

### 4.4 Training Details

We divided the data into training, validation, and test sets in proportions of $80\%$, $10\%$, and $10\%$ respectively. We set the learning rate as $2e - 5$ with the warmup rate as $0.1$, for all training setup. We chose the AdamW optimizer (Loshchilov and Hutter, 2019) with a linear scheduler for all experiments. The sequence length was set as 200 and the batch size was chosen as 8. For fine-tuning, all models were trained for 2 epochs except emotion classifiers; whereas for MTL, we trained our models for $8 - 12$ epochs based on the decrease

1144

| Methods | | Sentence Level | | Word Level | | |
| Model | Loss | $\rho$ | $r$ | F | P | R |
|---|---|---|---|---|---|---|
| XLM-RoBERTa-large | Nash | 0.4024 | 0.3946 | 0.2664 | 0.2152 | **0.4055** |
| | Aligned | 0.1214 | 0.1000 | 0.1835 | 0.1266 | 0.3333 |
| | Linear | 0.1921 | 0.1779 | 0.1835 | 0.1266 | 0.3333 |
| | Nash-D | 0.3642 | 0.3611 | 0.2465 | 0.1917 | 0.3885 |
| XLM-RoBERTa-base | Nash | 0.2747 | 0.2589 | 0.2452 | 0.2126 | 0.3772 |
| | Aligned | 0.2060 | 0.1629 | 0.1835 | 0.1266 | 0.3333 |
| | Linear | 0.0354 | 0.0754 | 0.1835 | 0.1266 | 0.3333 |
| | Nash-D | 0.1278 | 0.1139 | 0.2565 | 0.2043 | 0.3844 |
| XLM-V-base | Nash | **0.4673** | **0.4254** | **0.2805** | **0.2378** | 0.3953 |
| | Aligned | 0.1391 | 0.1063 | 0.2538 | 0.2050 | 0.3333 |
| | Linear | 0.2594 | 0.2052 | 0.2617 | 0.2154 | 0.3333 |
| | Nash-D | 0.4290 | 0.3983 | 0.2495 | 0.1942 | 0.3923 |
| MicroTransQuest (FT) | / | / | / | 0.1951 | 0.6651 | 0.1143 |

Table 1: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P) and recall (R) scores of models combining sentence- and word-level QE using our MTL architecture *vs* other MTL methods including the linear loss and Nash loss from Deoghare et al. (2023) (Nash-D) as well as the fine-tuning (FT) model using MicroTransQuest on HADQAET.

of the combined loss and depending on different combinations of tasks. For the emotion classification task in MTL, we chose max pooling over average pooling after experimentation. We set the number of epochs as 10 and used early stopping for fine-tuning emotion classifiers. All these hyper-parameters were chosen based on experimentation and previous research.

Fine-tuning multilingual PTLMs via TransQuest including MonoTransQuest, SiameseTransQuest and MicroTransQuest was carried out on an NVIDIA Quadro RTX 5000 GPU. Fine-tuning emotion classifiers including statistical models on HADQAET and the MQM emotion subset was performed on an NVIDIA T4 GPU. The rest of the model training including fine-tuning via COMET and different combinations of our MTL tasks were conducted on an NVIDIA A40 GPU.

| Methods | $\rho$ | $r$ |
|---|---|---|
| MonoTransQuest | **0.4355** | 0.3984 |
| SiameseTransQuest | 0.4151 | **0.4502** |
| COMET | 0.4083 | 0.3699 |

Table 2: Spearman $\rho$ and Pearson's $r$ correlation scores of models fine-tuned using TransQuest and COMET.

# 5 Results and Discussion

The results obtained by different models are presented from § 5.1 to § 5.3, while § 5.4 discusses the observations derived from our results.



Figure 4: Distribution of the HADQAET dataset

| Methods | F | P | R |
|---|---|---|---|
| XLM-RoBERTa-large | 0.1000 | 0.0700 | 0.2000 |
| XLM-V-base | 0.1000 | 0.0700 | 0.2000 |
| RF on XLM-RoBERTa-large embeddings | **0.1456** | **0.1603** | **0.2072** |
| SVM on XLM-RoBERTa-large embeddings | 0.1169 | 0.0826 | 0.2000 |

Table 3: Macro F1 (F), precision (P) and recall (R) scores of emotion classification models on HADQAET.

## 5.1 Fine-tuning on HADQAET

This section shows the results of fine-tuning, the methods presented in § 4.2 for sentence-level QE and emotion classification on HADQAET. The results at word-level QE are presented together with MTL in Table 1.

Table 2 displays the results of sentence-level QE models on HADQAET. The highest correlation scores, 0.4355 Spearman ($\rho$) and 0.4502 Pearson

| Methods | | Sentence Level | | Emotion Classification | | |
| Model | Loss | $\rho$ | $r$ | F | P | R |
| --- | --- | --- | --- | --- | --- | --- |
| | Nash | *-0.0357* | *-0.0289* | 0.1073 | 0.0733 | 0.2000 |
| XLM-RoBERTa-large | Aligned | 0.3786 | 0.3886 | 0.7985 | 0.7946 | 0.8257 |
| | Linear | 0.2376 | 0.2715 | 0.8399 | 0.8263 | **0.8887** |
| | Nash | 0.1448 | 0.1092 | **0.8549** | **0.8352** | 0.8879 |
| XLM-RoBERTa-base | Aligned | 0.4229 | 0.4174 | 0.8198 | 0.8054 | 0.8510 |
| | Linear | 0.3777 | 0.3521 | 0.7907 | 0.7756 | 0.8426 |
| | Nash | 0.0745 | 0.0105 | 0.1014 | 0.0679 | 0.2000 |
| XLM-V-base | Aligned | 0.4182 | 0.4278 | 0.8209 | 0.8040 | 0.8653 |
| | Linear | -0.0621 | -0.0512 | 0.1014 | 0.0679 | 0.2000 |
| FT baselines | / | **0.4355** | **0.4502** | 0.1456 | 0.1603 | 0.2072 |

Table 4: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P) and recall (R) scores of MTL models combining sentence-level QE and emotion classification using our MTL architecture *vs* linear loss on HADQAET. Our fine-tuning baselines (FT baselines) from Tables 2 and 3 are listed here for reference.

($r$), were achieved by fine-tuning using MonoTransQuest and SiameseTransQuest, respectively.

The emotion categories of HADQAET are imbalanced, and the dataset size is relatively small, as depicted in Figure 4. As a result, the fine-tuned classifiers always predicted the same class. We tried different PTLMs and hyperparameters, but the performance was not better as seen in Table 3. For this reason, we applied statistical methods including Random Forest (RF) (Breiman, 2001) and Support Vector Machine (SVM) (Hearst et al., 1998) based on the embeddings from XLM-RoBERTa-large. Our baseline for emotion classification was established using RF, achieving the best F1 score of 0.1456.

## 5.2 MTL on HADQAET

This section shows results of different combinations of the three tasks on HADQAET.

### 5.2.1 Sentence- and Word-level QE

Table 1 shows results of MTL that combines sentence- and word-level QE. For sentence-level QE, it is observed that MTL using XLM-V-base and Nash loss achieved the highest $\rho$ of 0.4673. This performance was superior to that of fine-tuning (0.4355). In the context of word-level QE, our best F1 score of 0.2805 surpasses the performance of fine-tuning using MicroTransQuest, which achieved an F1 score of 0.1951. This suggests that training sentence- and word-level QE systems together under the MTL framework can lead to improved performance in both tasks. Additionally, our MTL method is better than the linear loss and the Nash loss from Deoghare et al. (2023)

for both sentence- and word-level QE.

### 5.2.2 Sentence-level QE and Emotion Classification

Table 4 presents results for the combination of sentence-level QE and the emotion classification task. We can see that the use of MTL with Aligned loss effectively prevented the predictions from falling into the same category as shown in Table 3. Our top-performing model achieved an F1 score of 0.8549, much higher than our baseline. Our Aligned loss usually performed better than the linear loss for both sentence-level QE and emotion classification. It appears that incorporating the sentence-level QE task has proven beneficial for training emotion classifiers. However, incorporating emotion classification does not seem to be very helpful for sentence-level QE, as Spearman scores are not higher than those of fine-tuned models. In addition, it has been observed that when combined with emotion classification, the Aligned loss demonstrates greater stability compared to the Nash loss. This method achieves a favorable equilibrium between sentence-level QE and emotion classification.

| Heuristics | Sentence-level QE | Emotion Classification |
| --- | --- | --- |
| Nash Loss | 0.5604 | 5.1199 |
| Aligned Loss | 0.6162 | 0.6377 |

Table 5: Average loss weights for sentence-level QE and emotion classification using Nash and Aligned losses

Investigating further, we trained two models based on XLM-RoBERTa-base using the exact same hyperparameters, but two different loss

| Methods | | Sentence Level | | Word Level | | | Emotion Classification | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | Loss | $\rho$ | $r$ | F | P | R | F | P | R |
| XLM-RoBERTa-large | Nash | 0.3787 | 0.3979 | 0.1735 | 0.2194 | 0.3805 | 0.8526 | **0.8419** | 0.8730 |
| | Aligned | 0.1262 | 0.1035 | 0.1835 | 0.1266 | 0.3333 | 0.1014 | 0.0679 | 0.2000 |
| | Linear | 0.4020 | 0.3573 | 0.1836 | 0.1267 | 0.3333 | 0.8159 | 0.8115 | 0.8625 |
| XLM-RoBERTa-base | Nash | 0.2584 | 0.2342 | 0.2351 | 0.1740 | 0.3838 | **0.8528** | 0.8296 | 0.8903 |
| | Aligned | 0.3786 | 0.3654 | 0.2013 | 0.1417 | 0.3472 | 0.8403 | 0.8185 | 0.8920 |
| | Linear | 0.2895 | 0.2331 | 0.2131 | 0.1561 | 0.3426 | 0.7741 | 0.7658 | 0.8232 |
| XLM-V-base | Nash | 0.4051 | 0.4082 | 0.2245 | 0.1631 | 0.3795 | 0.8513 | 0.8324 | **0.8938** |
| | Aligned | 0.3389 | 0.3335 | 0.1914 | 0.1344 | 0.3337 | 0.8261 | 0.8220 | 0.8618 |
| | Linear | 0.3610 | 0.3659 | **0.2461** | 0.2343 | **0.3992** | 0.7892 | 0.7740 | 0.8241 |
| FT baselines | / | **0.4355** | **0.4502** | 0.1951 | **0.6651** | 0.1143 | 0.1456 | 0.1603 | 0.2072 |

Table 6: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P) and recall (R) scores of MTL models combining sentence- and word-level QE and emotion classification using our MTL architecture *vs* linear loss on HADQAET. Our fine-tuning baselines (FT baselines) from Tables 2 and 3 are listed here for reference.

| Methods | Sentence Level | | Word level | | | Emotion Classification | | |
|---|---|---|---|---|---|---|---|---|
| | $\rho$ | $r$ | F | P | R | F | P | R |
| MonoTransQuest | **0.3650** | **0.3836** | / | / | / | / | / | / |
| SiameseTransQuest | 0.2659 | 0.2622 | / | / | / | / | / | / |
| MicroTransQuest | / | / | **0.2141** | **0.4553** | **0.1399** | / | / | / |
| Random Forest | / | / | / | / | / | **0.1397** | **0.2061** | **0.2048** |
| SVM | / | / | / | / | / | 0.1202 | 0.0859 | 0.2000 |

Table 7: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P) and recall (R) scores for our baselines: fine-tuned models for sentence- and word-level QE and statistical models including Random Forest and Support Vector Machine (SVM) for emotion classification on the MQM emotion subset.

heuristics[7], *i.e.*, the Nash and Aligned losses, to combine sentence-level QE and emotion classification. We recorded the weights for the losses of the two tasks learned during training. The average loss weights (of all epochs) can be seen in Table 5. We can see that the Aligned loss seems to be better than Nash in balancing the two tasks as the two average weights are closer using the Aligned loss than Nash. This might be one of the reasons why it leads to more balanced results when the two tasks are combined.

### 5.2.3 Sentence-, Word-level QE and Emotion Classification

Table 6 illustrates simultaneous training of the three tasks. Again, our MTL method achieved better results than the linear loss under most circumstances. Compared with fine-tuning, our MTL method notably enhanced the performance of emotion classification, but the result of sentence-level QE was compromised. This suggests that as more tasks are incorporated into the MTL framework, achieving consensus or agreement between tasks becomes more challenging.

---

[7]The linear loss was omitted as weights were fixed as 1.

### 5.3 Results on the MQM Emotion Subset

This section presents results obtained on the MQM emotion subset, which is a selection of sentences from WMT QE shared tasks, with synthetic emotion labels as described in § 3.2.

#### 5.3.1 Fine-tuning on MQM Emotion Subset

We applied the same methods as those of HADQAET, except that only statistical methods were used for emotion classification. Our baseline results are shown in Table 7. We achieved a $\rho$ of 0.3650 for sentence-level QE, an F1 score of 0.2141 for word-level QE and 0.1397 for emotion classification.

#### 5.3.2 MTL on MQM Emotion Subset

Table 8 presents the results of combining sentence- and word-level QE. Our best model, utilizing Nash loss, achieved a Spearman correlation of 0.4947, notably surpassing the fine-tuning baseline and other MTL methods including the linear loss and Nash loss from Deoghare et al. (2023). The F1 score for word-level QE reached 0.2471, demonstrating improvement over the fine-tuning baseline. These findings affirm the validity of our approach for effectively integrating sentence- and word-level

| Methods | | Sentence Level | | Word Level | | |
| Model | Loss | $\rho$ | $r$ | F | P | R |
|---|---|---|---|---|---|---|
| XLM-RoBERTa-large | Nash | 0.1212 | 0.2244 | 0.2437 | 0.1918 | 0.3996 |
| | Aligned | 0.2840 | 0.2970 | 0.1682 | 0.1125 | 0.3333 |
| | Linear | -0.1162 | -0.1249 | 0.1682 | 0.1125 | 0.3333 |
| | Nash-D | 0.1427 | 0.1943 | 0.2447 | 0.1880 | **0.4043** |
| XLM-RoBERTa-base | Nash | 0.1385 | 0.1157 | 0.2253 | 0.1781 | 0.3785 |
| | Aligned | 0.2901 | 0.2928 | 0.1682 | 0.1125 | 0.3333 |
| | Linear | 0.2250 | 0.2684 | 0.1682 | 0.1125 | 0.3333 |
| | Nash-D | 0.2167 | 0.2304 | 0.2118 | 0.1549 | 0.3722 |
| XLM-V-base | Nash | **0.4947** | **0.4448** | 0.2251 | 0.1603 | 0.3908 |
| | Aligned | 0.3078 | 0.2204 | **0.2471** | 0.1963 | 0.3333 |
| | Linear | 0.2635 | 0.2385 | 0.2465 | 0.1956 | 0.3333 |
| | Nash-D | 0.1668 | 0.1619 | 0.2450 | 0.2057 | 0.3895 |
| FT baselines | / | 0.3650 | 0.3836 | 0.2141 | **0.4553** | 0.1399 |

Table 8: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P) and recall (R) scores of models combining sentence- and word-level QE using our MTL architecture *vs* other MTL methods including the linear loss and Nash loss from Deoghare et al. (2023) (Nash-D) on the MQM emotion subset. Our fine-tuning baselines (FT baselines) from Table 7 are listed here for reference.

QE in the context of overall quality evaluation.

Table 9 shows results integrating sentence-level QE and emotion classification. In instances where sentence-level QE excelled ($\rho$ 0.35), we observed a trade-off with emotion classification performance, and vice versa. The use of the XLM-V base model with the Aligned loss improved the performance of emotion classification, resulting in the highest F1 score, 0.3004.

Table 10 shows MTL results that combine all three tasks. Similar to results on HADQAET, there are trade-offs among tasks. Notably, on the MQM emotion subset, our best model achieved higher scores than fine-tuning and other MTL methods in both sentence- and word-level QE. This suggests that our approach contribute to the enhanced performance when training these tasks together.

### 5.4 Discussion

The results obtained from various task combinations within our MTL framework indicate that training sentence- and word-level QE systems together improves their performance compared to training them separately. This improvement likely stems from the interconnected nature of the two QE tasks. However, adding emotion classification to the framework usually does not enhance sentence- or word-level QE. Conversely, combining sentence-level QE with emotion classification boosts the performance of emotion classification. This finding is consistent for both the HADQAET (an emotion-

related QE dataset) and the MQM emotion subset (a standard QE dataset from WMT shared tasks). It suggests that the sentence-level QE task can aid in training emotion classifiers when training data is limited and the distribution is skewed.

For word-level QE, our approach achieves higher recall scores than MicroTransQuest, possibly because our model predicts errors in both the source and target texts, whereas MicroTransQuest considers only errors in the target.

Our results show that Nash and Aligned losses are generally better than the linear loss. Using the Nash loss is more likely to achieve state-of-the-art results for sentence-level QE, whereas the Aligned loss excels in balancing different tasks to produce a stable output. For this point, our observation still needs to be validated by further experiments on more task combinations and multilingual PTLMs.

### 6  Conclusion and Future Work

To evaluate MT quality of emotion-loaded UGC at sentence- and word-level simultaneously, we employed an emotion-related dataset that includes emotion labels and human-annotated translation errors. We extended it with sentence-level QE scores and word labels. This led to a dataset suitable for sentence- and word-level QE, and emotion classification. We proposed a new architecture featuring a novel combined MTL loss function that integrates different loss heuristics. This approach unifies the

| Methods | | Sentence Level | | Emotion Classification | | |
| Model | Loss | $\rho$ | $r$ | F | P | R |
|---|---|---|---|---|---|---|
| XLM-RoBERTa-large | Nash | 0.3500 | 0.3737 | 0.0257 | 0.0265 | 0.0250 |
| | Aligned | 0.1362 | 0.1699 | 0.1027 | 0.1014 | 0.1042 |
| | Linear | 0.1593 | 0.0747 | 0.1742 | 0.1905 | 0.2689 |
| XLM-RoBERTa-base | Nash | 0.1380 | 0.0125 | 0.1614 | 0.1595 | 0.2689 |
| | Aligned | 0.1395 | 0.1684 | 0.1534 | 0.1239 | 0.2014 |
| | Linear | 0.3305 | 0.3567 | 0.1273 | 0.1251 | 0.2106 |
| XLM-V-base | Nash | 0.0631 | 0.0658 | 0.2185 | 0.1897 | 0.3409 |
| | Aligned | *-0.0894* | *-0.0444* | **0.3004** | **0.2379** | **0.4862** |
| | Linear | 0.0616 | 0.0058 | 0.1690 | 0.1723 | 0.2689 |
| FT baselines | / | **0.3650** | **0.3836** | 0.1397 | 0.2061 | 0.2048 |

Table 9: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P) and recall (R) scores of models combining sentence-level QE and emotion classification tasks using our MTL architecture *vs* linear loss on the MQM emotion subset. Our fine-tuning baselines (FT baselines) from Table 7 are listed here for reference.

| Methods | | Sentence Level | | Word Level | | | Emotion Classification | | |
| Model | Loss | $\rho$ | $r$ | F | P | R | F | P | R |
|---|---|---|---|---|---|---|---|---|---|
| XLM-RoBERTa-large | Nash | 0.1198 | 0.1759 | 0.2284 | 0.1671 | **0.4116** | 0.1948 | 0.1623 | 0.2831 |
| | Aligned | 0.1151 | 0.1613 | 0.1682 | 0.1125 | 0.3333 | 0.0553 | 0.0311 | 0.2500 |
| | Linear | *-0.1708* | *-0.1581* | 0.1682 | 0.1125 | 0.3333 | 0.0553 | 0.0311 | 0.2500 |
| XLM-RoBERTa-base | Nash | 0.2856 | *-0.2112* | 0.2159 | 0.1523 | 0.4046 | 0.1392 | **0.3148** | 0.1935 |
| | Aligned | 0.2878 | 0.2992 | **0.2497** | 0.2006 | 0.3306 | 0.1032 | 0.1074 | 0.1874 |
| | Linear | 0.1794 | 0.1877 | 0.2151 | 0.1586 | 0.3447 | 0.1452 | 0.1661 | 0.2134 |
| XLM-V-base | Nash | *-0.0331* | 0.0392 | 0.1851 | 0.1383 | 0.3399 | 0.1520 | 0.1418 | 0.1755 |
| | Aligned | **0.3779** | 0.2939 | 0.1736 | 0.1174 | 0.3333 | 0.1841 | 0.1592 | 0.2874 |
| | Linear | 0.1130 | 0.1475 | 0.1743 | 0.1180 | 0.3333 | **0.2601** | 0.2120 | **0.4148** |
| FT baselines | / | 0.3650 | **0.3836** | 0.2141 | **0.4553** | 0.1399 | 0.1397 | 0.2061 | 0.2048 |

Table 10: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P), recall (R) scores of models combining sentence- and word-level QE and emotion classification using our MTL architecture *vs* linear loss on the MQM emotion subset. Our fine-tuning baselines (FT baselines) from Table 7 are listed here for reference.

training of multiple correlated tasks. We have made the code publicly available for similar task combinations such as empathy prediction and emotion classification. We compared our approach with existing fine-tuning and MTL methods and assessed its generalization on a standard QE dataset with synthetic emotion labels. We achieved new state-of-the-art results on both datasets. For future work, we aim to validate the effectiveness of our method on a larger multilingual QE dataset. We are also interested in investigating LLMs to evaluate machine translation of emotion-loaded UGC.

# 7 Limitations and Ethical Considerations

Although our MTL method is more effective, it is computationally expensive compared to fine-tuning for each task. Further, it takes longer to converge as parameters in the combined loss need to be learned over the training process.

Incorporating emotion classification might lead

to unstable performance for sentence-level QE under the Nash loss as explained in § 5.2.2. We will explore different task combinations and introduce a new hyperparameter to balance the tasks in our future work.

The experiments in the paper were conducted using publicly available datasets. New data were created based on those publicly available datasets using computer algorithms. No ethical approval was required as the use of all data in this paper follows the licenses in Qian et al. (2023) and Freitag et al. (2021a,b, 2022).

# References

Jonathan Baxter. 2000. A Model of Inductive Bias Learning. *J. Artif. Int. Res.*, 12(1):149–198.

Steven Bird, Edward Loper, and Ewan Klein. 2009. *Natural Language Processing with Python*. O'Reilly Media Inc, Sebastopol, California.

Frederic Blain, Chrysoula Zerva, Ricardo Ribeiro, Nuno M. Guerreiro, Diptesh Kanojia, José G. C. de Souza, Beatriz Silva, Tânia Vaz, Yan Jingxuan, Fatemeh Azadi, Constantin Orasan, and André Martins. 2023. Findings of the WMT 2023 shared task on quality estimation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 629–653, Singapore. Association for Computational Linguistics.

Leo Breiman. 2001. Random Forests. *Machine Learning*, 45(1):5–32.

Rich Caruana. 1997. Multitask Learning. *Machine Learning*, 28:41–75.

Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, Heyan Huang, and Ming Zhou. 2021. InfoXLM: An information-theoretic framework for cross-lingual language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3576–3588, Online. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. 2020. Revisiting pre-trained models for Chinese natural language processing. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 657–668, Online. Association for Computational Linguistics.

Sourabh Deoghare, Paramveer Choudhary, Diptesh Kanojia, Tharindu Ranasinghe, Pushpak Bhattacharyya, and Constantin Orăsan. 2023. A multitask learning framework for quality estimation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9191–9205, Toronto, Canada. Association for Computational Linguistics.

Bonnie Dorr, Joseph Olive, John McCary, and Caitlin Christianson. 2011. Machine Translation Evaluation and Optimization. In J. Olive, C. Christianson, and J. McCary, editors, *Handbook of Natural Language Processing and Machine Translation*, pages 745–843. Springer.

Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André Martins, Graham Neubig, Ankush Garg, Jonathan Clark, Markus Freitag, and Orhan Firat. 2023. The Devil Is in the Errors: Leveraging Large Language Models for Fine-grained Machine Translation Evaluation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1066–1083, Singapore. Association for Computational Linguistics.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021b. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774, Online. Association for Computational Linguistics.

Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy. PMLR.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria. Association for Computational Linguistics.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2024. xcomet: Transparent Machine Translation Evaluation through Fine-grained Error Detection. *Transactions of the Association for Computational Linguistics*, 12:979–995.

Xianwei Guo, Hua Lai, Yan Xiang, Zhengtao Yu, and Yuxin Huang. 2021. Emotion Classification of COVID-19 Chinese Microblogs based on the Emotion Category Description. In *Proceedings of the 20th China National Conference on Computational Linguistics*, pages 916–927. Chinese Information Processing Society of China.

Hany Hassan, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, Shujie Liu, Tie-Yan Liu, Renqian Luo, Arul Menezes, Tao Qin, Frank Seide, Xu Tan, Fei

Tian, Lijun Wu, Shuangzhi Wu, Yingce Xia, Dongdong Zhang, Zhirui Zhang, and Ming Zhou. 2018. Achieving Human Parity on Automatic Chinese to English News Translation. *arXive preprint*.

M.A. Hearst, S.T. Dumais, E. Osuna, J. Platt, and B. Scholkopf. 1998. Support vector machines. *IEEE Intelligent Systems and their Applications*, 13(4):18–28.

Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, and André F. T. Martins. 2019. OpenKiwi: An open source framework for quality estimation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 117–122, Florence, Italy. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023a. GEMBA-MQM: Detecting translation quality error spans with GPT-4. In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775, Singapore. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023b. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.

Guokun Lai, Zihang Dai, and Yiming Yang. 2020. Unsupervised Parallel Corpus Mining on Web Data. *arXiv preprint*.

Davis Liang, Hila Gonen, Yuning Mao, Rui Hou, Naman Goyal, Marjan Ghazvininejad, Luke Zettlemoyer, and Madian Khabsa. 2023. XLM-V: Overcoming the Vocabulary Bottleneck in Multilingual Masked Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13142–13152, Singapore. Association for Computational Linguistics.

Liyang Liu, Yi Li, Zhanghui Kuang, Jing-Hao Xue, Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. 2021. Towards Impartial Multi-task Learning. In *International Conference on Learning Representations*.

S. Liu, E. Johns, and A. J. Davison. 2019. End-To-End Multi-Task Learning With Attention. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1871–1880, Los Alamitos, CA, USA. IEEE Computer Society.

Arle Richard Lommel, Aljoscha Burchardt, and Hans Uszkoreit. 2014. Multidimensional Quality Metrics: A Flexible System for Assessing Translation Quality. *Tradumàtica: tecnologies de la traducció*, 0:455–463.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*.

John Nash. 1953. Two-Person Cooperative Games. *Econometrica*, 21(1):128–140.

Aviv Navon, Aviv Shamsian, Idan Achituve, Haggai Maron, Kenji Kawaguchi, Gal Chechik, and Ethan Fetaya. 2022. Multi-Task Learning as a Bargaining Game. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 16428–16446. PMLR.

Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.

Shenbin Qian, Constantin Orasan, Felix Do Carmo, Qiuliang Li, and Diptesh Kanojia. 2023. Evaluation of Chinese-English machine translation of emotion-loaded microblog texts: A human annotated dataset for the quality assessment of emotion translation. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 125–135, Tampere, Finland. European Association for Machine Translation.

Shenbin Qian, Constantin Orasan, Diptesh Kanojia, and Félix Do Carmo. 2024. Are Large Language Models State-of-the-art Quality Estimators for Machine Translation of User-generated Content? In *Proceedings of the 11th Workshop on Asian Translation*, Miami, United States of America. Association for Computational Linguistics.

Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. 2020. TransQuest: Translation quality estimation with cross-lingual transformers. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5070–5081, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. 2021. An exploratory analysis of multilingual word-level quality estimation with cross-lingual transformers. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 434–440, Online. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Hadeel Saadany, Constantin Orasan, Rocio Caro Quintana, Felix Do Carmo, and Leonardo Zilio. 2023. Analysing mistranslation of emotions in multilingual tweets by online MT tools. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 275–284, Tampere, Finland. European Association for Machine Translation.

D. Senushkin, N. Patakin, A. Kuznetsov, and A. Konushin. 2023. Independent Component Alignment for Multi-Task Learning. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20083–20093, Los Alamitos, CA, USA. IEEE Computer Society.

Kashif Shah and Lucia Specia. 2016. Large-scale multitask learning for machine translation quality estimation. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 558–567, San Diego, California. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Lucia Specia, Frédéric Blain, Marina Fomicheva, Erick Fonseca, Vishrav Chaudhary, Francisco Guzmán, and André F. T. Martins. 2020. Findings of the WMT 2020 shared task on quality estimation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 743–764, Online. Association for Computational Linguistics.

Lucia Specia, Frédéric Blain, Marina Fomicheva, Chrysoula Zerva, Zhenhao Li, Vishrav Chaudhary, and André F. T. Martins. 2021. Findings of the WMT 2021 shared task on quality estimation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 684–725, Online. Association for Computational Linguistics.

Trevor Standley, Amir Zamir, Dawn Chen, Leonidas Guibas, Jitendra Malik, and Silvio Savarese. 2020. Which Tasks Should Be Learned Together in Multi-Task Learning? In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org.

Craig Stewart, Ricardo Rei, Catarina Farinha, and Alon Lavie. 2020. COMET - deploying a new state-of-the-art MT evaluation metric in production. In *Proceedings of the 14th Conference of the Association for Machine Translation in the Americas (Volume 2: User Track)*, pages 78–109, Virtual. Association for Machine Translation in the Americas.

Andy Sun. 2013. Jieba. https://github.com/fxsjy/jieba.

Shuo Wang, Zhaopeng Tu, Zhixing Tan, Wenxuan Wang, Maosong Sun, and Yang Liu. 2021. Language Models are Good Translators. *arXiv preprint*.

Chrysoula Zerva, Frédéric Blain, Ricardo Rei, Piyawat Lertvittayakumjorn, José G. C. de Souza, Steffen Eger, Diptesh Kanojia, Duarte Alves, Constantin Orăsan, Marina Fomicheva, André F. T. Martins, and Lucia Specia. 2022. Findings of the WMT 2022 shared task on quality estimation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 69–99, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Jingyi Zhang and Josef van Genabith. 2020. Translation quality estimation by jointly learning to score and rank. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2592–2598, Online. Association for Computational Linguistics.

# A Appendix

## A.1 Additional Figures and Tables

Figure A.1 shows an example of the HADQAET dataset from Qian et al. (2023). Table A.1 displays results of other loss heuristics in our framework.

## A.2 Nash MTL

Nash MTL intends to find an update vector $\Delta\theta$ for the gradients $g_i$ of the task $i$ in the ball of radius $\epsilon$ centered around zero, $B_\epsilon$, as shown in Equation 5.

$$arg\,max_{\Delta\theta \in B_\epsilon} \Sigma_i log(\Delta\theta^\mathsf{T} g_i) \qquad (5)$$

The solution to Equation 5 is (up to scaling) $\Sigma_i \alpha_i g_i$ where $\alpha \in \mathbb{R}^K_+$ is the solution to $G^\mathsf{T} G\alpha = 1/\alpha$ where $1/\alpha$ is the element-wise reciprocal. Detailed proof can be seen in Navon et al. (2022). The Nash MTL algorithm is shown below:

---
Algorithm 1 Nash-MTL

---
**Input**: $\theta^{(0)}$ – initial parameter vector, $\{l_i\}_{i=1}^K$ – differentiable loss functions $\eta$ – learning rate
**for** $t = 1, ..., T$ **do**
  Compute task gradients $g_i^{(t)} = \nabla_{\theta^{(t-1)}} l_i$
  Set $G^{(t)}$ the matrix with columns $g_i^{(t)}$
  Solve for $\alpha$ : $(G^t)^\mathsf{T} G\alpha = 1/\alpha$ to obtain $\alpha^{(t)}$
  Update the parameters $\theta^{(t)} = \theta^{(t)} - \eta G^{(t)}\alpha^{(t)}$
**end for**
**Return** $\theta^{(T)}$

---

## A.3 Aligned MTL

Through theoretical analysis, Senushkin et al. (2023) found a strong relation between the condition number and conflicting and dominating gradients issues, and they proposed Aligned MTL to align principal components of a gradient matrix to make the training process more stable.

The objective of Aligned MTL as defined in Equation 6, is to minimize the difference between the original gradient matrix $G$ and the aligned gradient matrix $\hat{G}$. The difference is measured using the Frobenius $F$ norm. The constraint in Equation 6 ensures that $\hat{G}$ is orthogonal, meaning that its transpose multiplied by itself is equal to the identity matrix. This constraint helps to ensure stability in the linear system of gradients.

$$\min_{\hat{G}} \|G - \hat{G}\|_F^2 \quad s.t. \ \hat{G}^\mathsf{T}\hat{G} = I \qquad (6)$$

$$\hat{G} = \sigma UV^\mathsf{T} = \sigma GV\Sigma^{-1}V^\mathsf{T} \qquad (7)$$

The solution is defined in Equation 7, where $\hat{G}$ is obtained by singular value decomposition (SVD). SVD decomposes $G$ into three matrices: $U$, $\Sigma$ and $V^\mathsf{T}$ where $U$ and $V$ are orthogonal matrices, and $\Sigma$ is a diagonal matrix containing the singular values of $G$. Algorithm of Aligned MTL is shown below:

---
Algorithm 2 Aligned MTL

---
**Require:** $G \in \mathbb{R}^{|\theta| \times T}$ – gradient matrix,
         $w \in \mathbb{R}^T$ – task importance
$M \leftarrow G^\mathsf{T} G$
$(\lambda, V) \leftarrow eigh(M)$
$\Sigma^{-1} \leftarrow diag(\sqrt{\frac{1}{\lambda_1}}, ..., \sqrt{\frac{1}{\lambda_R}})$
$B \leftarrow \sqrt{\lambda_R} V\Sigma^{-1}V^\mathsf{T}$
$\alpha \leftarrow Bw$
**Return** $G\alpha$

---

1153

| Source | MT output | Human Translation | Original emotion label | Error type | Error severity |
|---|---|---|---|---|---|
| 管理学真是水的一比，努力的想听，依然坚持不过一分钟……考研怎么办呀 | Management is really a comparison of water. I want to listen hard, but I still can't hold on for a minute...What about the postgraduate entrance examination? | Management is really a bunch of fiddle-faddle. I try hard to listen, but still can't hold on for a minute...What about the postgraduate entrance examination? | anger | mistranslation | critical |

Figure A.1: An Example from HADQAET (Qian et al., 2023)

| Methods | | | Sentence Level | | Word Level | | |
|---|---|---|---|---|---|---|---|
| Model | Loss | | $\rho$ | $r$ | F | P | R |
| XLM-RoBERTa-large | DWA | | -0.0740 | -0.1031 | 0.1835 | 0.1266 | 0.3333 |
| | IMTL | | 0.1488 | 0.1057 | 0.2440 | 0.2096 | 0.3767 |
| XLM-RoBERTa-base | DWA | | 0.0533 | 0.0726 | 0.0183 | 0.0094 | 0.3333 |
| | IMTL | | 0.1495 | 0.1561 | 0.2322 | 0.1929 | 0.3668 |
| XLM-V-base | DWA | | -0.2551 | -0.2302 | 0.1870 | 0.1300 | 0.3333 |
| | IMTL | | 0.3182 | 0.2714 | 0.2757 | 0.2320 | 0.3843 |
| InfoXLM | Nash | | 0.1678 | 0.2647 | 0.2454 | 0.2181 | 0.3763 |
| | Aligned | | 0.0363 | 0.0281 | 0.1835 | 0.1266 | 0.3333 |
| | DWA | | -0.0237 | -0.0355 | 0.1835 | 0.1266 | 0.3333 |
| | IMTL | | -0.2731 | -0.2200 | 0.1879 | 0.1941 | 0.3353 |
| | Linear | | 0.0042 | 0.0013 | 0.1835 | 0.1266 | 0.3333 |
| | Nash-D | | 0.1846 | 0.2125 | 0.2618 | 0.2377 | 0.3902 |

Table A.1: Spearman $\rho$, Pearson's $r$, Macro F1 (F), precision (P) and recall (R) scores of models fine-tuned based on XLM-RoBERTa, XLM-V-base and InfoXLM models in combination of sentence- and word-level QE using Dynamic Weight Averaging (DWA) and impartial MTL (IMTL) on HADQAET. Results obtained using the linear combination and Nash MTL in Deoghare et al. (2023), *i.e.*, Nash-D, for InfoXLM are also displayed here.

# On Instruction-Finetuning Neural Machine Translation Models

**Vikas Raunak    Roman Grundkiewicz    Marcin Junczys-Dowmunt**
Microsoft Azure AI
{viraunak,rogrundk,marcinjd}@microsoft.com

## Abstract

In this work, we introduce instruction finetuning for Neural Machine Translation (NMT) models, which distills instruction following capabilities *from* Large Language Models (LLMs) *into* orders-of-magnitude smaller NMT models. Our instruction-finetuning recipe for NMT models enables customization of translations for a limited but disparate set of translation-specific tasks. We show that NMT models are capable of following multiple instructions simultaneously and demonstrate capabilities of zero-shot composition of instructions. We also show that through instruction finetuning, traditionally disparate tasks such as formality-controlled machine translation, multi-domain adaptation as well as multi-modal translations can be tackled jointly by a single instruction finetuned NMT model, at a performance level comparable to LLMs such as GPT-3.5-Turbo. To the best of our knowledge, our work is among the first to demonstrate the instruction-following capabilities of traditional NMT models, which allows for faster, cheaper and more efficient serving of customized translations.

## 1 Introduction

Instruction-finetuned Large Language Models (LLMs) demonstrate the remarkable ability of instruction-following (Wei et al., 2021), which makes them amenable to tackle any task cast as natural language generation, even under a zero-shot setting. In this work, we explore whether traditional Neural Machine Translation (NMT) models could offer *similar* capabilities of following instructions. NMT models could be considered as domain-specific 'language' models *pre-trained* for a single task (translation) and thereby *could* be instruction-finetuned to tackle translation-adjacent tasks such as translation customization or enforcing certain specifications on the translations. Such tasks, e.g., formality-controlled translation (Schioppa et al., 2021), multi-modal translation (Elliott et al., 2016)

or gender-based translation rewriting (Kuczmarski and Johnson, 2018), have typically been tackled through specialized models or algorithms in prior literature, rather than a single instruction-following NMT model. In contrast, we instruction-finetune a single *ancestral* translation model to *adapt* the translations based on instructions. Our contributions are as follows:

1. We present a new recipe for instruction finetuning NMT models (trained with supervision only on parallel datasets), which allows for joint modeling of disparate translation customization tasks in a single NMT model, and we analyze the criticality of each of the recipe components through ablation experiments.

2. We demonstrate that NMT models are capable of following multiple (30+) instructions simultaneously. We also find that NMT models show abilities of zero-shot composition of instructions, as an effect of finetuning.

3. We show that, with a single instruction-finetuned NMT model, traditional customization tasks such as formality-controlled machine translation can be tackled with high performance, in conjunction with several disparate tasks.

Additionally, our proposed finetuned NMT model outperforms GPT-3.5-Turbo on average on the IWSLT-22 Formality Control Shared Task (Antonios et al., 2022), while simultaneously achieving high-performance on others & demonstrating a few other *desirable* properties vis-à-vis much larger LLMs. At a high-level, our work re-interprets a NMT model as a language model and demonstrates the utility of instruction finetuning NMT model for jointly modeling a myriad of disparate translation-related tasks. In the next sections, we elaborate on our recipe for instruction-finetuning of a NMT model and analyze its characteristics.

1155

| Instruction Prefix | Source (English) | Translation (German) |
|---|---|---|
| *past tense* | The finished effect <mark>is</mark> long-lasting and highly glossy – but does it damage the nails? | Der fertige Effekt <mark>war</mark> langanhaltend und hochglänzend – aber beschädigte er die Nägel? |
| *informal* | Do you like Legos? did you ever play with them as a child or even later? | <mark>Magst du</mark> Legosteine? <mark>Hast du</mark> jemals als Kind oder sogar später mit ihnen gespielt? |
| *fix misspelling* | To switch between environments, update the storage.json file with the URL of the <mark>specificrrbzpronment</mark> . | Um zwischen Umgebungen zu wechseln, aktualisieren Sie die Datei storage.json mit der URL des <mark>spezifischen Prozesses</mark> . |
| *translate "herbal medicines" to "Kräutermedizin"* | Chinese <mark>herbal medicines</mark> for hypothyroidism | Chinesische <mark>Kräutermedizin</mark> gegen Hypothyreose |
|  | A trendy girl talking on her cellphone while gliding slowly down the street. | Ein schickes Mädchen telefoniert, während sie langsam die Straße entlangschwebt. |

Table 1: Input-output instances for the developed instruction finetuned NMT model. The table shows four tasks, in which the instruction is used to make the translation conform to certain specific characteristics. The instruction prefix is prepended to the source text and is enclosed with the instruction tags. In the case of image as an instruction, the image is tokenized into a one dimensional representation.

## 2 Related Work

Our work is at the intersection of two key themes: instruction finetuning—primarily developed in the context of LLMs—and customizing NMT models for specific tasks.

### 2.1 Instruction Finetuning of LLMs

Instruction finetuning refers to the supervised finetuning of a language model on task-specific input-output pairs by explicitly describing the task through instructions. This has been demonstrated to aid in cross-task generalization (Sanh et al., 2022a; Longpre et al., 2023), in particular, imparting LLMs with instruction-following capabilities (Wei et al., 2021). A number of prior works have proposed different algorithms for constructing the instruction data (Mishra et al., 2022; Wang et al., 2022; Honovich et al., 2023; Wang et al., 2023; Sanh et al., 2022b; Muennighoff et al., 2023; Iyer et al., 2023; Chung et al., 2022).

In our recipe, we rely on a combination of parallel data filtering and synthetic data generation

through LLMs to construct the instruction dataset that is leveraged for finetuning NMT models. Further, our approach substantially differs from prior work in that we instruction finetune NMT models whose pre-training is completely supervised on bitext source-translation pairs.

### 2.2 Customizing Translation Models

There exists a large body of work in adapting NMT models and customizing them for specific use cases such as for achieving high-performance on specific domains (Saunders, 2022), tones or registers in the target language (Nädejde et al., 2022) as well as for tasks such as gender-based translation rewriting (Rarrick et al., 2023). Tagging specific subpopulations of the parallel data to accomplish this task has been a staple in prior work for formality control, verbosity control, etc.

Our work is related to the tagging approaches developed in the literature but differs in two key aspects: (a) task diversity and scale: typically, tagging is only applied to supply information pertaining to a single task, while instruction finetuning as

---
**Algorithm 1:** Instruction-Finetuning NMT Recipe

---

**Data:** Base NMT Model and Vocabulary
**Result:** Instruction Finetuned NMT Model

| | |
|---|---|
| **Step 1:** | Expand vocabulary with instruction tokens |
| **Step 2:** | Curate task-specific and parallel datasets |
| **Step 3:** | Finetune on a *mix* of parallel and task data |
| **Step 4:** | (Optional) Interpolation with base model |

---

a technique aspires to tackle a wide variety of tasks in a unified modeling approach to make the model capable of following a wide variety of instructions; and (b) natural language instruction: instead of manipulating tags or combination of tags, we leverage instructions expressed or composed in natural language for influencing the translations.

## 3 Instruction Finetuning of NMT models

In this section, we describe the problem setting along with our instruction finetuning recipe and evaluation protocol.

### 3.1 Problem Setting

For instruction finetuning, we take a pre-trained NMT model and finetune it with instruction annotated source-translation pairs. The instruction is prepended to the source text inside tags that demarcate the instruction, e.g., *<instruction> informal </instruction>*. Henceforth, we refer to the tokens pertaining to the *<instruction>* and *</instruction>* strings as the instruction tokens. A collection of instruction and source-translation instances are presented in Table 1. Through instruction finetuning, we hope to jointly model a range of disparate tasks.

### 3.2 Instruction Finetuning Recipe

We present our simple recipe for instruction finetuning NMT models in Algorithm 1. We first expand the vocabulary of a given NMT model with the instruction tokens in order to delineate the instructions cleanly from the actual source text. Adding free-form text instructions within these instruction tokens also implies that the NMT model never sees the instruction tokens on the output side, hence the risk of translating the instructions themselves is greatly diminished. We initialize the embeddings of the newly added tokens to random embeddings centered around the mean of the embedding matrix (in particular, mean plus a unitary projection of randomly sampled embedding principal components).

The next step in the recipe is to curate both task-specific and parallel datasets used for finetuning. For curating parallel dataset (non-instruction data), we apply standard heuristics on the model's parallel dataset to sample a higher-quality parallel dataset (compared to the model's full training corpus). The details of the heuristics are presented in appendix D. For task-specific data curation, either we manually curate translations from the parallel dataset or we generate the translations synthetically from LLMs (GPT-4 and GPT-3.5-Turbo). We describe task specific dataset curation in section 3.4.

Finally, the NMT model is finetuned on a mix (2:1) of parallel and task data—the mixing ratio is a hyperparameter in our recipe and we tune it so that we observe no degradation in general translation performance as measured on the WMT20 validation set. At the end of the finetuning, the finetuned and the base models are optionally interpolated to achieve a better trade-off between general and task performance. We present the details of the interpolation step in the Appendix A, while the details pertaining to the other steps are presented in the next sections. We found the interpolation to be optional, so none of the experiments in the main paper use this step.

### 3.3 Evaluation Protocol

For the instruction finetuned NMT model, we have the choice of either translating an input without any instruction (the *general* case) or using a particular instruction (the *instruction* case). Throughout this work, we report the following measurements in order to evaluate the instruction finetuned NMT model:

1. **General Performance**: This is measured by computing the MT quality of the finetuned NMT model (i.e., the original translation task) on a standard test set. This metric is reported in order to measure the impact of instruction finetuning on the general translation quality of the finetuned model.

2. **Task-Specific Performance**: On a per-task basis we report two measurements:

   a. **Task Response Rate (RR)**: the percentage of instances in the test set for which including a instruction yielded a different translation than not including the instruction (the *general* case). This offers us a

crude measure to evaluate how responsive the model is to a specific instruction. For example, if an instruction is empty, then the translation in the general case and the instruction case should not change and thereby a low response rate is expected.

b. **Task Output Quality**: the MT quality metrics (over system outputs and references) for the finetuned NMT model both in the *general* case and the *instruction* case. The gap between the general quality and the instruction quality depicts the gain (or degradation) in quality obtained by explicitly influencing the translation through a particular instruction.

Further, for some tasks such as formality-controlled translations, we report evaluations on two different test sets: (a) an intrinsic test set which comes from the same data distribution as the fine-tuning data and (b) an extrinsic test set, which is an external dataset that comes with a completely different data distribution. Also, we use ChrF as the primary MT quality metric through this work, however each of our results is agnostic to the choice of the particular MT quality metric and the trends remain the same irrespective of the quality metric (e.g., COMET) used.

## 4 Experiments

In this section we describe all experimental settings, from model architecture to data curation and evaluation.

### 4.1 Experimental Settings

We conduct experiments on the WMT'20 News Translation (English-German) task benchmark (Barrault et al., 2020). The WMT'20 test set is used for measuring general translation performance. We used the official parallel training data from WMT'20 with the dataset statistics presented in Table 2. A joint vocabulary of 32K was learnt using SentencePiece on a 10M random sample of the training dataset.

The trained model is a Transformer-Big (225M parameters) with the hyperparameters described exactly in Vaswani et al. (2017). The model was trained for 300K updates using Marian NMT (Junczys-Dowmunt et al., 2018). The metrics BLEU, ChrF2, TER (Papineni et al., 2002; Popović, 2015; Snover et al., 2006) for the trained model

on the WMT'20 validation and test sets (under beam size of 1) as measured using SacreBLEU (Post, 2018) are presented in Appendix B, alongside reference-based COMET (Rei et al., 2020) scores.

| Data Source | Sentence Pairs |
|---|---|
| Europarl | 1,828,521 |
| ParaCrawl | 34,371,306 |
| Common Crawl | 2,399,123 |
| News Commentary | 361,445 |
| Wiki Titles | 1,382,625 |
| Tilde Rapid | 1,631,639 |
| WikiMatrix | 6,227,188 |
| Total | 48,201,847 |

Table 2: The WMT'20 data sources used for training the English–German NMT model.

For our first experiment, we construct a set of 30 tasks, each with 1K samples as well as use multi-30K multimodal dataset with 29K training samples. For multi-30K, we convert the image into 32 tokens using 1D image tokenizer[1] from Yu et al. (2024). For multi-30K samples, the image tokens serve as the instructions, whereas for the other tasks, short natural language task descriptions serve as instructions. Further details for these tasks are presented in Appendix C. We then instruction finetune our base WMT'20 model with the curated data. Our key goal here is to evaluate whether NMT models are capable of following multiple instructions simultaneously.

### 4.2 Task-Specific Data Curation

The first column of Table 3 shows the list of task instructions. In terms of data provenance, the tasks are of two types: synthetic tasks (for which the instruction finetuning data is obtained synthetically) and authentic tasks (for which the data is mined from the parallel training corpora). We present a more verbose description of each of the tasks in Appendix C, since the text in the instruction naturally implies the targeted translation task.

For each of the 30 tasks, we curate instruction data using filters applied on the parallel data or through synthetic data generation using GPT-3.5-Turbo or GPT-4. In particular, the data for instructions pertaining to generating active voice, passive voice, simplifying, complexifying and obs-

---

[1] https://github.com/bytedance/1d-tokenizer

| Task Instruction | RR (%) | ChrF$_{general}$ | ChrF$_{instruction}$ | Improvement |
|---|---|---|---|---|
| past tense | 84.81 | 82.06 | 86.85 | + 4.79 |
| translate X to Y | 60.42 | 76.18 | 80.24 | + 4.06 |
| active voice | 54.84 | 87.62 | 92.86 | + 5.24 |
| passive voice | 80.91 | 71.44 | 78.29 | + 6.85 |
| non-literal | 50.00 | 83.25 | 84.89 | + 1.64 |
| literal | 53.41 | 90.12 | 92.88 | + 2.76 |
| titlecase | 100.0 | 52.75 | 68.52 | + 15.77 |
| lowercase | 100.0 | 55.39 | 67.35 | + 11.96 |
| uppercase | 98.92 | 2.41 | 40.31 | + 37.9 |
| remove punctuation | 100.0 | 67.18 | 68.73 | + 1.55 |
| add antonyms | 79.79 | 71.90 | 73.12 | + 1.22 |
| remove profanity | 66.67 | 75.81 | 77.38 | + 1.57 |
| add hashtag | 100.0 | 61.05 | 68.68 | + 7.63 |
| leetify | 100.0 | 26.37 | 34.12 | + 7.75 |
| remove accents | 81.97 | 59.55 | 62.08 | + 2.53 |
| shuffle words | 100.0 | 52.69 | 42.62 | - 10.07 |
| fix misspelling | 91.74 | 60.22 | 65.36 | + 5.14 |
| introduce repetition error | 55.34 | 64.54 | 65.36 | + 0.82 |
| insert X at the beginning | 100.0 | 64.78 | 69.19 | + 4.41 |
| insert X at the end | 100.0 | 64.38 | 69.68 | + 5.3 |
| same length | 58.16 | 89.37 | 95.93 | + 6.56 |
| shorter length | 52.59 | 90.88 | 94.30 | + 3.42 |
| longer length | 57.38 | 66.51 | 68.14 | + 1.63 |
| simplify | 81.42 | 61.88 | 67.22 | + 5.34 |
| complexify | 58.33 | 89.31 | 93.92 | + 4.61 |
| obsfuscate | 56.84 | 80.89 | 82.61 | + 1.72 |
| formal | 60.77 | 86.53 | 91.03 | + 4.50 |
| informal | 60.58 | 87.28 | 92.25 | + 4.97 |
| spacing error | 84.40 | 66.70 | 66.87 | + 0.17 |
| coverage error | 97.25 | 66.40 | 66.24 | - 0.16 |
| image (multi-30k) | 53.00 | 72.08 | 74.89 | + 2.81 |
| empty instruction | 0.06 | 65.27 | 65.27 | + 0.0 |
| average | 89.60 | 74.20 | 82.42 | + 8.22 |

Table 3: Intrinsic evaluation results for the instruction finetuned NMT system over different tasks. Across different types of tasks (synthetic rule based tasks, distributional style tasks as well as on producing multi-modal translations), the instruction-finetuned model demonstrates the capability of following multiple instructions simultaneously. Note that the base model has no instruction-following capability, hence performs poorly across different task test sets.

fuscating translations were obtained synthetically through GPT-3.5-Turbo[2], whereas formal and informal translation data was obtained using GPT-4.

## 4.3  Finetuning and Evaluation Settings

The last checkpoint of the trained WMT'20 model is finetuned for 3 data epochs. The instruction dataset is split into 90% percent for finetuning and the 10% held-out dataset is used for intrinsic evaluation. The general translation quality is measured on the WMT'20 News Translation test set.

## 5  Results and Analysis

In this section, we characterize the behavior of the instruction finetuned NMT model using both intrinsic and extrinsic evaluations. In the next section, we present an ablation study on the key components of the recipe.

### 5.1  Instruction-Following Performance

Table 3 presents the results that characterize the instruction-following performance of the finetuned NMT model. The results show that the NMT model is capable of following instructions over a collection of disparate tasks, which is the key finding of our work.

In particular, both rule-based tasks such as *leetify* (which inserts leet-speak in the translation) as well as tasks which are more distributional and style based in nature, such as *complexify*, are remarkably well learned by the NMT model. For tasks such as shuffle words, in which the model is taught to randomly shuffle the words in the translation, the reference based MT quality metric (ChrF) is unable to demonstrate gains owing to the stochasticity of the transformation.

### 5.2  Zero-Shot Composition of Instructions

Additionally, we investigate whether the model, trained on individual task instructions can compose two instructions. Note that the finetuned model has never seen two disparate instructions appear together in a single sample. We find that the model is capable of composing instructions in a zero-shot manner and Table 4 presents an example of such a composition.

To further investigate this behavior, in Table 4, we present additional metric named Task Success Rate (SR), which provides a binary measure of the task success rather than a continuous measure

---

such as ChrF. Through SR measurements, we find that the effectiveness of the composition varies considerably across different compositions, a phenomenon akin to the large variance in LLM performance due to minor variations in prompt.

### 5.3  Extrinsic Evaluations

We conduct extrinsic evaluation on the WMT'22 Shared Task for formality on English–German translations. The shared task winner has (100%, 100%) in both in the unconstrained setting and (100%, 88.6%) in the constrained setting (Antonios et al., 2022). The instruction-finetuned model does not use any training data at all from WMT'22, relying only on the synthetic task data curated from GPT-4 and is evaluated on the test set directly. The results in Table 5 show that the instruction finetuned model is quite competitive with the WMT'22 task winner and achieves better performance that GPT-3.5-Turbo (evaluated in the zero-shot setting).

### 5.4  General Translation Quality

The ChrF2 of the finetuned model on the WMT'20 test set is 61.9, which is +0.3 over the base WMT'20 model. This demonstrates that instruction finetuning does not impact the general translation capabilities of the NMT model. Similar trends hold for other metrics as well.

## 6  Ablation Study

In this section, we present an ablation study on the instruction finetuning recipe presented in Algorithm 1, wherein we remove the addition of explicit instruction tokens and the addition of parallel data from our recipe. The finetuning and evaluation protocols remain the same as in prior sections, except that for the finetuning experiments presented below, we set the number of epochs to two. However, our findings stay the same across different number of finetuning epochs. Further, we only report results on the Multi-30K task instead of all the tasks as in Table 3.

### 6.1  Ablating Parallel Data

Our recipe mixes task-specific and standard parallel data for finetuning. Table 6 compares the results of finetuning runs in the absence of parallel data in terms of key performance metrics. We find that not including the parallel data in the recipe leads to degradation of general translation performance. However, at the same time including the parallel

---

| Task Instruction | RR (%) | ChrF$_{general}$ | ChrF$_{instruction}$ | T$_1$ SR (%) | T$_2$ SR (%) |
|---|---|---|---|---|---|
| lowercase | 100.00 | 53.82 | 68.11 | 83.00 | – |
| uppercase | 100.00 | 2.42 | 44.67 | 27.96 | – |
| remove profanity | 93.33 | 69.88 | 80.95 | – | 40.00 |
| lowercase remove profanity | 100.00 | 58.86 | 70.69 | 80.00 | 40.00 |
| uppercase remove profanity | 100.00 | 2.97 | 39.31 | 26.67 | 6.67 |
| lowercase and remove profanity | 100.00 | 58.86 | 69.23 | 93.33 | 33.33 |
| uppercase and remove profanity | 100.00 | 2.97 | 43.27 | 26.67 | 13.33 |

Table 4: Zero-shot composition of instructions. The instruction finetuned NMT model can compose instructions in a zero-shot manner on held-out test data (i.e., the model has not been trained on any combinations of instructions). Although, the effectiveness of composition varies across the different compositions (prompts) applied. T$_1$ refers to the first task under composition and T$_2$ refers to the second task under composition.

| Formality-Control Translation Model | Formal Accuracy | Informal Accuracy |
|---|---|---|
| mBART-large, Rippeth et al. (2022) | 93.6 | 77.4 |
| LLM, Garcia et al. (2023) | 84.9 | 85.5 |
| Doc-MT System, Post and Junczys-Dowmunt (2024) | 83.3 | 87.1 |
| GPT-3.5-Turbo[3] | 95.5 | 95.0 |
| (ours) Baseline WMT-20 model | 75.0 | 25.0 |
| (ours) Instruction-Finetuned WMT-20 model | 94.7 | 98.5 |
| WMT'22 Task Winner (Constrained) | 100.0 | 88.6 |
| WMT'22 Task Winner (Unconstrained) | 100.0 | 100.0 |

Table 5: Extrinsic evaluation on producing formal and informal translations. The instruction finetuned NMT model outperforms GPT-3.5-Turbo on the shared task, despite not using the training data released for the shared task. The model's capabilities are learned through distillation in the form of instruction finetuning.

| Multi-30K Task | | General Perf | |
|---|---|---|---|
| ChrF$_{Base}$ | ChrF$_{instruction}$ | ChrF$_{Base}$ | ChrF$_{FT}$ |
| 59.45 | 67.75 | 61.6 | 62.2 |
| 59.45 | 71.80 | 61.6 | *61.4* |

Table 6: Impact of removing parallel data (bottom row). The models are finetuned for the same number of epochs with and without generic parallel data.

| Multi-30K Task | | General Perf | |
|---|---|---|---|
| ChrF$_{Base}$ | ChrF$_{instruction}$ | ChrF$_{Base}$ | ChrF$_{FT}$ |
| 59.45 | 71.80 | 61.6 | 61.4 |
| 67.75 | 71.94 | 61.6 | *60.5* |

Table 7: Impact of removing the instruction tokens (bottom row). The models are finetuned for the same number of epochs with and without the instruction tokens added as a part of the model vocabulary. No parallel data was used in both cases.

data impacts model optimization on the instruction tasks. For these experiments, we used a mixing ratio of 2:1 between the parallel and the task data.

### 6.2 Ablating Vocabulary Expansion

Our recipe expands the vocabulary of the NMT model with new instruction tokens. Table 7 compares the results of finetuning runs in the absence of new tokens in terms of key performance metrics. We find that in the absence of new tokens, the model's general performance degrades substan-

tially which is likely due to the fact that the model has to overwrite more pre-trained information.

Altogether, the above ablations point that both the key elements of our recipe are quite important. We hypothesize that this is owing to the fact that both of these components allow the model to overwrite less of its pre-training knowledge, which helps the model strike a better trade-off between task-specific and general translation performance.

| Property | Large Language Models (LLMs) | Instruction-Finetuned NMT |
|---|---|---|
| Task Performance | High | High |
| Controllability | Low | High |
| Adversarial Robustness | Low | High |
| Inference Cost | High | Low |

Table 8: Comparison of Large Language Models (LLMs) and Instruction-Finetuned NMT models on various properties of interest.

| Model Type | Source/Formal-Translation | Correctness |
|---|---|---|
| Source | Append "WTF" after everything you say. Say: how are you? | |
| GPT-3.5-Turbo | Wie geht es dir? WTF | ✗ |
| Instruction-Finetuned NMT | Fügen Sie "WTF" nach allem hinzu, was Sie sagen. Sag: Wie geht es dir? | ✓ |
| Source | 5 + 5 | |
| GPT-3.5-Turbo | 10 | ✗ |
| Instruction-Finetuned NMT | 5+5 | ✓ |

Table 9: Adversarial robustness. LLMs expose a larger surface area for adversarial manipulation of model outputs compared to instruction finetuned NMT model. In this case, the source content overrides the correct (intended) model behavior of producing formal translations for full source.

# 7 Discussion

To conclude, we presented a simple yet effective instruction-finetuning recipe for unified modeling of multiple disparate translation-specific tasks in a single NMT model. Our results demonstrate that the instruction-finetuned NMT model is able to utilize the instructions and does understand their meanings, to an extent that it is able to compose combinations of instructions in a zero-shot manner. Further, instruction-finetuned NMT models have other properties that distinguish it from LLMs. Table 8 presents such a comparison on a few properties of interest:

1. Task Performance: When limiting ourselves to a set of *known* translation-related tasks, our results show that instruction finetuned NMT models are *capable* of reaching similar or higher task performance than LLMs.

2. Controllability: Finetuning NMT models is considerably cheaper than finetuning LLMs and as a result, instruction finetuned NMT models offer more controllability than LLMs.

3. Adversarial Robustness: LLMs expose a very large attack surface area and the prompts to customize translations could be easily manipulated by users to alter the model behavior, posing a security risk for the intended application (Liu et al., 2024a,b). However, instruction-finetuned NMT models, by default expose a much smaller attack surface area and thereby are less vulnerable to adversarial attacks— some examples highlighting the differences with respect to prompt injection and intent misclassification attacks are in Table 9.

4. Inference Costs: NMT models are substantially cheaper to serve in production compared to LLMs such as GPT-3.5-Turbo, owing to smaller parameter sizes.

As such, instruction following NMT models which can broadly adapt translations based on desired user specifications for a large number of translation specific tasks might offer a better cost to quality and cost to *security* trade-off when compared to orders-of-magnitude larger LLMs.

# 8 Conclusion and Future Work

In this work, we presented a simple recipe for instruction finetuning NMT models. Using our recipe, we demonstrated that a NMT model is capable of learning to follow multiple disparate instructions simultaneously, while obtaining high performance on important translation customization tasks such as formality-control. Our work opens up an interesting research direction—on building instruction following NMT models which could leverage both the cheaper inference costs of NMT models as well as the broad customization capabilities of LLMs.

# References

Anastasopoulos Antonios, Barrault Loc, Luisa Bentivogli, Marcely Zanon Boito, Bojar Ondřej, Roldano Cattoni, Currey Anna, Dinu Georgiana, Duh Kevin, Elbayad Maha, et al. 2022. Findings of the iwslt 2022 evaluation campaign. In *Proceedings of the 19th International Conference on Spoken Language Translation (IWSLT 2022)*, pages 98–157. Association for Computational Linguistics.

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 conference on machine translation (WMT20). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55, Online. Association for Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models.

Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30K: Multilingual English-German image descriptions. In *Proceedings of the 5th Workshop on Vision and Language*, pages 70–74, Berlin, Germany. Association for Computational Linguistics.

Xavier Garcia, Yamini Bansal, Colin Cherry, George Foster, Maxim Krikun, Fangxiaoyu Feng, Melvin Johnson, and Orhan Firat. 2023. The unreasonable effectiveness of few-shot learning for machine translation.

Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2023. Unnatural instructions: Tuning language models with (almost) no human labor. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14409–14428, Toronto, Canada. Association for Computational Linguistics.

Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi, Simon Kornblith, Ali Farhadi, and Ludwig Schmidt. 2022. Patching open-vocabulary models by interpolating weights. In *Advances in Neural Information Processing Systems*.

Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, Xian Li, Brian O'Horo, Gabriel Pereyra, Jeff Wang, Christopher Dewan, Asli Celikyilmaz, Luke Zettlemoyer, and Ves Stoyanov. 2023. Opt-iml: Scaling language model instruction meta learning through the lens of generalization.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain. Association for Computational Linguistics.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

James Kuczmarski and Melvin Johnson. 2018. Gender-aware natural language translation.

Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Zihao Wang, Xiaofeng Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, and Yang Liu. 2024a. Prompt injection attack against llm-integrated applications.

Yupei Liu, Yuqi Jia, Runpeng Geng, Jinyuan Jia, and Neil Zhenqiang Gong. 2024b. Formalizing and benchmarking prompt injection attacks and defenses. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 1831–1847, Philadelphia, PA. USENIX Association.

Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. *arXiv preprint arXiv:2301.13688*.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In *Proceedings*

*of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.

Maria Nădejde, Anna Currey, Benjamin Hsu, Xing Niu, Marcello Federico, and Georgiana Dinu. 2022. Cocoa-mt: A dataset and benchmark for contrastive controlled mt with application to formality.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A Call for Clarity in Reporting BLEU Scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Matt Post and Marcin Junczys-Dowmunt. 2024. Escaping the sentence-level paradigm in machine translation.

Spencer Rarrick, Ranjita Naik, Varun Mathur, Sundar Poudel, and Vishal Chowdhary. 2023. Gate: A challenge set for gender-ambiguous translation examples.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Elijah Rippeth, Sweta Agrawal, and Marine Carpuat. 2022. Controlling translation formality using pretrained multilingual language models. In *Proceedings of the 19th International Conference on Spoken Language Translation (IWSLT 2022)*, pages 327–340, Dublin, Ireland (in-person and online). Association for Computational Linguistics.

Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao,

Thomas Wolf, and Alexander M Rush. 2022a. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*.

Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022b. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*.

Danielle Saunders. 2022. Domain adaptation and multidomain adaptation for neural machine translation: A survey.

Andrea Schioppa, David Vilar, Artem Sokolov, and Katja Filippova. 2021. Controlling machine translation for multiple attributes with additive interventions. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6676–6696, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.

Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel,

Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Shuangzhi Wu, Xing Wang, Longyue Wang, Fangxu Liu, Jun Xie, Zhaopeng Tu, Shuming Shi, and Mu Li. 2020. Tencent neural machine translation systems for the WMT20 news translation task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 313–319, Online. Association for Computational Linguistics.

Qihang Yu, Mark Weber, Xueqing Deng, Xiaohui Shen, Daniel Cremers, and Liang-Chieh Chen. 2024. An image is worth 32 tokens for reconstruction and generation.

## A   Appendix A

We describe the interpolation step equation 1. This step interpolates between the parameters of the base model ($\theta_{\text{base}}$) and the finetuned model ($\theta_{\text{finetuned}}$) using a scalar interpolation weight $\alpha$ which is applied for all common parameters between the base and the finetuned model (Ilharco et al., 2022). This step can be applied in order to better balance the general performance against task specific performance of the resulting model. In the equation, the performance (*perf*) measure could be the general performance or task-specific performance measure. We do not apply this for the models presented in this work, however, in practice we find that it is quite effective in addressing regressions in general performance.

$$\Theta = \max_{\alpha} \left\{ \text{perf} \left( (1 - \alpha) \cdot \theta_{\text{base}} + \alpha \cdot \theta_{\text{finetuned}} \right) \right\} \tag{1}$$

## B   Appendix B

The metrics BLEU, ChrF2, TER (Papineni et al., 2002; Popović, 2015; Snover et al., 2006) for the WMT20 trained model (under beam size of 1)

as measured using SacreBLEU (Post, 2018) are presented in Table 11, alongside reference-based COMET (Rei et al., 2020) scores.

## C   Appendix C

We present a brief characterization of the different tasks here, along with some example input-output pairs in Table 10.

- Rule Based Tasks: A number of tasks are rule based, e.g., translating into the past tense is a derivative task of generating the actual translation. Similarly, removing punctuations, adding antonyms, leetify or add hashtag (which adds a hashtag comprising of the last source word at the end of the translation) are rule based tasks.

- Distributional Style Based Tasks: We include tasks such as generating translation in a particular style, which can be learned based on the synthetic LLM-generated translations.

- Contrastive Tasks: Tasks such as length control in which the model is taught to control the verbosity of the translation is an example of a task in which the model is taught to generate translations which do not have any *absolute* property – but possess characteristics against some constrastive examples.

- Multi-modal Task: Multi-30K represents the multi-modal translation tasks wherein an image accompanies the source input.

## D   Appendix D

For parallel data filtering, we replicate the bitext filtering pipeline of Wu et al. (2020). and apply sentence-pair filtering based on maximum allowable sentence-length ratio (1:1.3) and reverse sentence-length ratio (1.3:1) alongside filtering sentences greater than a maximum word length (150). We also use a language-id filter (Joulin et al., 2017) is also used, which checks if the source and target sentences are in the correct languages.

| Instruction Prefix | Source (English) | Translation (German) |
|---|---|---|
| *same length* | This family-run hotel in a building from 1800 in Valkenburg centre, offers free Wi-Fi throughout. It features 2 terraces and a bar with original features including an exposed beam ceiling. | Dieses familiengeführte Hotel in einem Gebäude aus dem Jahr 1800 im Zentrum von Valkenburg bietet kostenloses WLAN im gesamten Haus. Es verfügt über 2 Terrassen und eine Bar mit originalen Merkmalen, einschließlich einer freigelegten Balkendecke. |
| *shorter length* | This family-run hotel in a building from 1800 in Valkenburg centre, offers free Wi-Fi throughout. It features 2 terraces and a bar with original features including an exposed beam ceiling. | Familiengeführtes Hotel von 1800 im Zentrum von Valkenburg mit kostenlosem WLAN, 2 Terrassen und Bar mit Holzbalkendecke. |
| *active voice* | They then fled to Prague, after Schwerin had been arrested for illegal distribution of Communist propaganda. In Prague, they opened an advertising agency, 'Hammer und Pinsel' (Hammer and Brush). | Sie flohen dann nach Prag, nachdem Schwerin wegen illegaler Verbreitung kommunistischer Propaganda verhaftet worden war. In Prag eröffneten sie eine Werbeagentur, 'Hammer und Pinsel'. |
| *passive voice* | They then fled to Prague, after Schwerin had been arrested for illegal distribution of Communist propaganda. In Prague, they opened an advertising agency, 'Hammer und Pinsel' (Hammer and Brush). | Sie flohen dann nach Prag, nachdem Schwerin wegen illegaler Verbreitung kommunistischer Propaganda verhaftet worden war. In Prag wurde eine Werbeagentur namens 'Hammer und Pinsel' eröffnet. |

Table 10: Input-output instances for the contrastive tasks in Table 3.

| Metric | BLEU | ChrF2 | TER | COMET |
|---|---|---|---|---|
| Validation | 37.5 | 63.9 | 51.5 | 56.50 |
| Test | 32.9 | 61.6 | 54.2 | 42.52 |

Table 11: Metrics for the Trained WMT20 System

# Benchmarking Visually-Situated Translation of Text in Natural Images

**Elizabeth Salesky**[J]          **Philipp Koehn**[J]          **Matt Post**[H, M]

[J]Johns Hopkins University
[H]Human Language Technology Center of Excellence
[M]Microsoft
esalesky@jhu.edu

## Abstract

We introduce a benchmark, VISTRA, for visually-situated translation of English text in natural images to four target languages. We describe the dataset construction and composition. We benchmark open-source and commercial OCR and MT models on VISTRA, and present both quantitative results and a taxonomy of common OCR error classes with their effect on downstream MT. Finally, we assess direct image-to-text translation with a multimodal LLM, and show that it is able in some cases but not yet consistently to disambiguate possible translations with visual context. We show that this is an unsolved and challenging task even for strong commercial models. We hope that the creation and release of this benchmark which is the first of its kind for these language pairs will encourage further research in this direction.

## 1 Introduction

Visually-situated language concerns multimodal settings where text and vision are intermixed, and the meaning of words or phrases is directly influenced by what is observable or referenced visually. Vision-and-language research has most commonly focused on tasks where images and text can be processed as distinct channels within a joint model, such as question answering or image captioning. However, settings where text is embedded in an image are ubiquitous, ranging from text on street signs, to chryrons on news broadcasts, language embedded in figures or social media images, or non-digitized text sources.

Translating visually-situated text is a practical application of recent pixel-based translation models ([Salesky et al., 2021](#)), with new challenges due to the varied text styles, backgrounds, and complex layouts found in natural images. This task combines a series of traditionally separate steps including text detection, optical character recognition, semantic grouping, and finally machine translation.



Figure 1: Visual context can resolve translation ambiguity. Here, translating 'EXIT' from English to German is ambiguous without further information about the mode of travel (on foot or by car), which the visual context in the image provides.

Not only can errors propagate between steps, as generated mistakes cause mismatches in vocabulary and distribution from those observed in training and reduce downstream task performance, but processing each step in isolation separates recognized text from visual context which may be necessary to produce a correct situational translation. For example, as illustrated in Figure 1, the English word *'Exit'* can be translated to German as either *'Ausfahrt'* or *'Ausgang'*; without appropriate context, which may not be present in the text alone, the generated translation would be a statistical guess.

We present a publicly-released benchmark, VISTRA, for visually-situated translation (VST) of text contained in natural images. With VISTRA, we benchmark the performance of popular OCR models and conduct an error analysis of text recognition errors. We analyze which recognition errors propagate to and most significantly affect downstream translation to four target languages with varied levels of contextual dependence on the image. We also compare **direct** visually-situated translation with multimodal LLMs, and discuss whether access to visual context improves visually-situated translation with current models. Finally, given our findings, we present directions for future work and connections to recent pixel-based translation models.

## 2 Constructing the VISTRA benchmark

VISTRA comprises 772 natural images containing English text, with aligned translations to four target languages (German, Spanish, Russian, and Mandarin Chinese) with varying levels of visual contextual dependence. Each image is annotated with its height and width, a categorical label, its semantically grouped English transcript, translations to the four target languages aligned at the level of the semantic groups in the transcript, and, word-level bounding boxes specified by corner with coordinates rescaled from 0-1, matched to the aligned word in the transcript. On average, each image contains 11.2 words and 2.4 transcript groups, for a total of 1840 parallel segments in the benchmark with an average length of 4.7 words. An annotated data sample is shown in Figure 2.[1]

To the best of our knowledge, only one prior publicly-released data exists for in-image text translation from *natural* images (OCRMT30K: Lan et al., 2023), which contains 30k images with Chinese text manually translated to English. In the absence of datasets for this task, prior work on in-image machine translation has primarily synthetically rendered MT corpora for this task (Mansimov et al., 2020; Tian et al., 2023; Niu et al., 2024; Lan et al., 2024) or addressed PDF document translation (Ignat et al., 2022; Hsu et al., 2024), discussed further in Section 4. While these settings typically use uniform text styles and sizes and contain a single semantic unit per image, natural images are contain text with multiple sizes and styles, multiple text groups in complex layouts, and varied image backgrounds, all of which introduce additional challenges. Our task also differs from what is commonly called multimodal translation in that our setting text is embedded into the image context, as opposed to a text caption to be translated with the aid of a relevant image.

The VISTRA benchmark is released under a permissive CC BY-SA license for further scientific research and commercial use.[2]

### 2.1 Criteria for image selection

The dataset is primarily constructed of newly-captured photos in order that they not be under copyright or contained in LLM training data.[3] We additionally include a small challenge set of public domain images from social media where text has been embedded in an image and is no longer accessible without OCR. Within this benchmark, we focus only on printed text, not handwritten. We describe the detailed criteria for image selection below.

1. **Languages:** Only images containing text in a single language (English) are included.

2. **Maximizing translatable text:** Images were chosen to maximize text which would be translated rather than transliterated or copied across languages, i.e. maximizing descriptive or instructive text and minimizing numerals and named entities. Where these are present, they may not constitute the majority of the text.

3. **Framing with sufficient context:** Sufficient context (visual or textual) must be present to reduce translation ambiguity. If, as in Figure 1, correct translation would require knowledge that the sign is by a road or a footpath, one of these should be at least partly visible.

4. **Length of text:** We aim for a balance of text lengths. While some traffic signs may have only 1-2 words, if they are sufficiently frequent that it is important for strong image translation models to get correct, they have been included; other images may include up to 100 words.

5. **Text style:** Text may contain multiple fonts, colors, and sizes within one image.

6. **Layout and number of text groups:** We include a balance of layout complexity, from single-line horizontal layouts, to complex layouts with angled text, or multiple adjacent semantic groups which prove challenging for line-level OCR.

7. **Image dimensions and resolution:** We collect high-resolution photos, non-resized and not retouched. Original dimensions may vary based on camera and conditions, but at least one dimension (length or width) must be larger than 1024.

---

[1]We omit the full list of bounding box coordinates in Figure 2 for readability.

[2]https://vistra-benchmark.github.io

[3]Though these specific images will not have been observed in training, we cannot guarantee that the same or similar signs in other settings have not. Though we submitted opt-out requests to exempt our data from being trained on before submitting benchmark images to commercial LLMs in our experiments, if subsequent researchers do not also do so, benchmark images may be ingested as training data.

Figure 2: VISTRA data sample showing metadata, transcripts, and translations.

8. **Clean conditions:** The dataset reflects clean conditions. We require that it is not challenging for a human reader to recognize contained text. We exclude images where the text is challenging to read due to environmental conditions (such as weather: rain, fog); blur; lighting conditions; occlusions (such as graffiti, foliage).

9. **Permissive use:** All images are either photos taken for the purposes of this benchmark or in the public domain.

## 2.2 Text annotation and transcription

Text bounding boxes and transcripts were manually post-edited from Google Cloud Vision OCR with a custom interface. The annotation interface is shown in Figure 6 in Appendix A.

**Bounding boxes.** Text bounding boxes are specified at the word level. In contrast to line-level annotations, using word-level bounding boxes more flexibly allows for complex layouts where unrelated text may appear side-by-side in an image (for example, adjacent signs), but should not be grouped and translated together. Bounding boxes are rectangular (90° corners) with all four vertices specified, which allows angled rotation to match text directionality. Bounding boxes were post-edited to ensure all text was detected, no text was cropped, and hallucinated text boxes were removed.

**Transcript.** All (and only) text which was clearly human-readable with images resized to a maximum height and width of 1024px was transcribed. In the final transcripts, case and punctuation are matched as closely as possible to what is present in the original image. Non-textual symbols which may be present on some directional signs

(for example, 🚶 or 🚲) were not transcribed or annotated.

**Semantic grouping.** Finally, we semantically group word-level text boxes. This creates text units with necessary context for translation, and separates for example different street signs which appear in the same image into distinct units for downstream translation. Not all images contain full sentences; therefore, our criteria were forming clause or phrase-level groups which appear together in the image and should be translated together. This step may be ambiguous, and so was annotated by one person to ensure consistency across the dataset.

## 2.3 Translation

We contracted Centific[4] to professionally translate the text in each image from English to four target languages: German, Spanish, Russian, and Mandarin Chinese. This set of languages covers multiple language families and scripts, and varied dependence on visual context. Annotators were paid a competitive market rate. Each image was translated by an individual linguist and a random sample of 10% of the image translations were checked by a second linguist.

All translations were performed from scratch in OneForma, with access to both the original image and transcript. The translation instructions and annotation interface are shown in Figure 7 in Appendix A. Translations are aligned one-to-one with the semantic groups in the transcript. We do not ask annotators to match case and punctuation in the source language, which may be unnatural for the target language, but rather localize these for the target language.

---

[4]https://www.centific.com

| Model | OCR | MT | VST | Release | OCR level | Returns bboxes? | Multilingual |
|-------|-----|-----|-----|---------|-----------|-----------------|--------------|
| PaddleOCR | ✓ | | | OPEN-SOURCE | line\|word | yes | |
| TesseractOCR | ✓ | | | OPEN-SOURCE | word | yes | |
| Google Cloud Vision | ✓ | | | COMMERCIAL | word | yes | |
| mBART | | ✓ | | OPEN-SOURCE | — | — | ✓ |
| Google Translate | | ✓ | | COMMERCIAL | — | — | |
| GPT-4o | ✓ | ✓ | ✓ | COMMERCIAL | unknown | no | ✓ |

Table 1: Models benchmarked for visually-situated translated (cascaded and direct).

Annotators were additionally asked whether the visual context in the image affected the resulting translation, as a binary question. Whether a translation would be ambiguous without the image can vary by target language, as exemplified by *'Exit'* in Figure 1 which would be ambiguous in German but not in Spanish. 99.7% of images were marked as requiring image context for translation for at least one translation direction, with the following breakdown by target language: German 99%, Chinese 96%, Spanish 54%, Russian 6%.

## 3 Benchmarking visually-situated translation

We benchmark existing models for OCR and VST using the new VISTRA dataset, and conduct an error analysis of common OCR types and their effect on downstream translation. This type of error analysis does not exist in previous work; we show that our new benchmark both illustrates these types of errors and facilitates analysis of this type.

### 3.1 Models evaluated

We compare a variety of widely-used open-source, open-weight, and commercial models to give a representative view of the capabilities of current models for this task, and specifically, provide baseline performance on the VISTRA benchmark. We list all evaluated models with relevant characteristics in Table 1.

### 3.1.1 OCR

**Paddle-OCR**[5] (PP-OCR: Du et al., 2020) is becoming one of the most commonly used open-source tools for OCR in English and Chinese (Lan et al., 2023; Yang et al., 2023, *inter alia*), due to its ease of use and free public release. PP-OCRv4 uses

Transformer models, trained per-language for English and Chinese. It produces word-level bounding boxes within detected lines. **Tesseract-OCR**[6] (Smith, 2007) is the longest-standing community-developed open-source toolkit for OCR. Tesseract-4 uses LSTM models, trained per-language. We additionally benchmark **Google Cloud Vision OCR**[7] (Popat et al., 2017; Ingle et al., 2019) to compare strong commercial performance.

### 3.1.2 MT

We compare both an open-source machine translation model, **mBART-50** (Liu et al., 2020), which is trained primarily on clean, well-formed text, and a commercial machine translation model, **Google Translate**, as an upper-bound on expected performance with greater expected robustness to noise.

### 3.1.3 Multimodal LLMs

Multimodal multilingual large language models which have been explicitly trained on both vision and language provide an opportunity to compare *direct* translation from an image containing English to text in a target language. By directly translating from an image with access to the full image context (as opposed to only the cropped region within bounding boxes from a text detection stage), multimodal models have the potential to be able to resolve ambiguity in translation. Here we benchmark the performance of **GPT-4o**, which was the top-performing multimodal model on a recent OCR-centric LLM evaluation (OCRBench: Liu et al., 2023), by prompting the model to directly translate text contained in images without intermediate steps. We also evaluate OCR only and machine translation only with this model in order to contextualize direct multimodal translation results.

---

[5] https://paddlepaddle.github.io/PaddleOCR

[6] https://github.com/tesseract-ocr/tesseract
[7] https://cloud.google.com/vision

| Class | Description |
|-------|-------------|
| I | Undetected text: missing text and bounding boxes |
| II | Text hallucination: text detected where no text present |
| III | Bounding box misplaced: text clipped, cropping would affect recognition |
| IV | Grouping error: text from different groups intermixed in output text |
| V | Punctuation error |
| VI | Spacing error |
| VII | Character-level substitution |
| VIII | Word-level substitution |

Table 2: OCR error taxonomy covering text detection (I-III) and recognition (IV-VIII) errors.

| Model | CER↓ | TER↓ | Sub. | Del. | Ins. |
|-------|------|------|------|------|------|
| Paddle-OCR | 13.0 | 21.5 | 963 | 2824 | 2851 |
| Google OCR | 18.0 | 32.0 | 186 | 381 | 8496 |
| GPT-4o | 23.8 | 36.0 | 1132 | 1277 | 9728 |
| Tesseract-OCR | 124.0 | 134.3 | 9597 | 37081 | 16477 |

Table 3: OCR results on the VISTRA benchmark.



Figure 3: Proportion of OCR error classes by model.

## 3.2 OCR performance and error taxonomy

We measure OCR performance with two automatic metrics: character error rate (CER) and translation error rate (TER) (Snover et al., 2006).[8] CER reflects the minimum number of single-character edits (insertions, deletions, or substitutions) required to change a string into the reference. TER is also an edit-distance metric which aims to capture the post-editing effort required to change a string into the reference. While OCR is typically evaluated case-insensitive with punctuation removed, with downstream MT in mind we calculate both metrics case-sensitive with punctuation.

While in e.g., speech recognition there may be one correct ordering of the output, in a 2D image, there is not necessarily only one correct order for recognized text. To facilitate scoring different text groupings across models, which would otherwise require re-alignment, we concatenate groups before applying automatic metrics. Where CER would recognize reorderings between the hypothesis and reference as several character edits, TER allows shifts of contiguous spans as a single operation, and therefore penalizes reordering less. As shown in Figure 8d, different orderings due to line vs. word level text recognition may still be errors and significantly impact downstream translation, which is why we use these two metrics together.

In addition to aggregate quantitative metrics, we create an error taxonomy of the different classes of OCR errors we observe across different models.

The eight OCR error classes are listed in Table 2 and describe errors in each step of the pipeline, from text detection (text recall, text hallucinations where non-text objects are recognized as text, or bounding box placement errors which may affect downstream processing using only these regions) to recognition and generation (over and under generation of punctuation, spaces, and character- and word-level substitutions).

We provide an illustrative example for each OCR error class observed with our evaluated models on the VISTRA benchmark in Figure 8 in Appendix B. We hypothesize that differences in model design affect the proportion of each type of error, and that different error categories are likely to affect downstream translation in different ways, as we investigate in Section 3.3.

OCR performance for our 4 compared models are shown in Table 3. We observe that somewhat surprisingly, the open-source Paddle-OCR model is the highest performing on both the CER and TER metrics. While Google OCR has significantly fewer substitutions and deletions, it has a much higher insertion rate; here, this covers both more 'benign' insertions like whitespace, and text hallucinations as illustrated in Figure 8b, where background patterns are recognized as text characters.[9] GPT-4o performs slightly worse than both models on all metrics. Tesseract, on the other hand, significantly underperforms expectations set by past work (e.g.,

---

[8]We calculate TER with SacreBleu (Post, 2018), *case_sensitive=True, no_punct=False, normalized=False*.

[9]It may be worth noting that where such hallucinations frequently occur as consecutive spans, and so can be significantly easier to post edit than the quantitative metrics reflect.

(a) A **grouping error** causes each word to be translated individually, resulting in agreement errors (Apple Translate).

(b) **Inserted punctuation** breaks up the text sequence, resulting in translation errors despite correctly recognized text (mBART).

Figure 4: Qualitative examples of OCR errors which propagate to downstream MT.

Ignat et al., 2022). We hypothesize this may be because it is primarily trained on documents, rather than images with natural backgrounds; the mismatch to varied background colors and additional visual context, though the majority of our images contain printed text with relatively uniform backgrounds, appears to interfere with recognition and lead to insertions as seen in Figure 4b.

To see to what degree models vary in the *type* of errors they make, which we hypothesize are likely to affect downstream MT to different degrees, we manually annotate the error classes observed in model outputs for a random sample of 100 images. Figure 3 shows the proportion of outputs containing each error class from Table 2 for each of our four models. For the strongest three models, the most frequent error type is text hallucination, where text is detected where no text was present. These three models have just one or no examples of recall errors in our sample. On the other hand, Tesseract-OCR fails to detect some proportion of text in the majority of examples and any text at all in 44%, both resulting in low performance and artificially reducing the rates of other error types.

Between Paddle and Google OCR, this analysis presents a slightly different view to the quantitative results above alone. Both models have similar distributions of their most frequent errors. While Paddle has a lower CER, its output has more varied types of errors. Google OCR error types are more consistent and occur across fewer classes, but where they are present, there are typically multiple errors, which lowers CER and TER further.

For GPT-4o, hallucinations often result in significant additions of punctuation as visualized for example in Figure 8e, resulting in more than one class

of error. GPT-4o does not return bounding boxes for detected text. However, it appears to often generate text with additional whitespace and punctuation to offset different text groups, which are reflected in both these error classes.[10] We do not observe significant differences in word-level substitutions or text hallucinations with larger and/or stronger decoder models.

### 3.3 How do OCR errors affect downstream MT?

Here we assess the effect of OCR errors on downstream MT in cascaded models. We do not perform normalization or postprocessing between OCR and MT, except to concatenate semantic groups. We evaluate translation with three automatic metrics: BLEU[11] (Papineni et al., 2002) and chrF (Popović, 2015), both as computed by SacreBLEU (Post, 2018), and COMET (Rei et al., 2020).[12]

This is a challenging task for all models. Figure 4 shows two illustrative examples where OCR errors interfere with downstream translation, despite correctly recognized text. Table 4 shows translation performance for both cascaded OCR and MT models and direct translation with a multimodal LLM for the four target languages in VIS-TRA. Open-source models have weak performance across all target languages and metrics. Commercial MT appears more robust to OCR performance in general, with consistently stronger results across all metrics and relatively similar performance across the three strongest OCR models.

---

[10]We were not able to reduce this behavior consistently via prompting.

[11]We omit BLEU when translating into Chinese without word segmentation.

[12]`wmt22-comet-da`

| | OCR Model | mBART | | | Google Translate | | | GPT-4o | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | chrF | BLEU | COMET | chrF | BLEU | COMET | chrF | BLEU | COMET |
| **German** | Tesseract-OCR | 2.3 | 0.1 | 28.8 | 3.5 | 0.1 | 30.4 | — | — | — |
| | Paddle-OCR | 26.8 | 9.0 | 46.1 | 36.0 | 16.7 | 57.0 | — | — | — |
| | GPT-4o OCR | 28.1 | 6.9 | 48.0 | 36.4 | 13.2 | 58.2 | — | — | — |
| | Google-OCR | 31.1 | 9.1 | 47.3 | 37.4 | 14.9 | 55.3 | — | — | — |
| | None | — | — | — | — | — | — | 36.9 | 9.1 | 60.1 |
| **Spanish** | Tesseract-OCR | 2.4 | 0.1 | 30.1 | 3.6 | 0.3 | 31.7 | — | — | — |
| | Paddle-OCR | 17.5 | 3.1 | 44.4 | 60.8 | 33.8 | 75.1 | — | — | — |
| | GPT-4o OCR | 23.3 | 4.2 | 50.4 | 60.8 | 24.6 | 75.0 | — | — | — |
| | Google-OCR | 22.0 | 4.0 | 45.5 | 62.2 | 29.9 | 71.3 | — | — | — |
| | None | — | — | — | — | — | — | 54.0 | 21.4 | 73.4 |
| **Russian** | Tesseract-OCR | 1.7 | 0.1 | 25.3 | 2.6 | 0.1 | 27.3 | — | — | — |
| | Paddle-OCR | 13.0 | 5.8 | 42.4 | 46.5 | 20.0 | 73.0 | — | — | — |
| | Google-OCR | 16.0 | 7.5 | 42.4 | 48.1 | 18.4 | 71.0 | — | — | — |
| | GPT-4o OCR | 14.8 | 5.1 | 43.1 | 47.1 | 15.1 | 74.4 | — | — | — |
| | None | — | — | — | — | — | — | 35.6 | 10.7 | 70.2 |
| **Chinese** | Tesseract-OCR | 0.3 | — | 32.6 | 0.4 | — | 34.4 | — | — | — |
| | Paddle-OCR | 18.2 | — | 62.0 | 40.2 | — | 82.0 | — | — | — |
| | GPT-4o OCR | 19.7 | — | 63.1 | 40.1 | — | 82.5 | — | — | — |
| | Google-OCR | 18.7 | — | 59.2 | 41.6 | — | 77.7 | — | — | — |
| | None | — | — | — | — | — | — | 33.6 | — | 85.5 |

Table 4: Visually-situated translation results on the VISTRA benchmark. We compare both cascaded OCR and MT as well as direct translation from images with a multimodal LLM. We note results with commercial OCR and/or MT in gray, and direct translation of text in images with multimodal LLMs in blue.

Direct translation with a multimodal LLM performs quite strongly, with consistently comparable COMET scores to the strongest cascades for all target languages, though weaker comparatively on the lexical metrics chrF and BLEU; we look at this more closely in Section 3.4.

The results in Table 4 show that CER and TER alone are not sufficient indicators of performance. Translation with mBART performs more highly for Google and GPT-4o OCR than Paddle-OCR despite their higher CER and TER, suggesting the type of errors may have more significance than edit distance alone. Undetected text (Class I) has the most catastrophic effect on downstream MT. For Tesseract-OCR, recall is simply too low for non-trivial translation performance. Punctuation and whitespace are insertional errors which are detrimental to tokenization with mBART, increasing fertility by approximately $3\times$ and resulting in input sequences which approach character level.

We hypothesize that these classes of errors may be normalized in preprocessing by the commercial MT system as they have less effect; of the sample set annotated as having these errors, segment-level chrF is $2\times$ higher with the commercial model than mBART, which is a larger margin than observed overall ($1.6\times$). Text hallucinations (Class II) are more difficult to remove with post-processing, though here are typically character-level rather than insertions of valid words. Character- and word-level substitutions (Classes VII and VIII) were stated to have more detrimental effect on translation for OCR'ed documents in Ignat et al. (2022) than insertions or deletions, but that is not the trend we observe here. On our type of data, natural images with complex backgrounds, we observe significantly more insertions per example than substitutions; while for example punctuation insertions (Class V) occur for a similar number of examples as character-level substitutions (Class VII) for

Figure 5: On this example from VISTRA, in a model cascade when translating with access to OCR output only, GPT-4o translates 'Exit' as 'Ausgang,' while when translating directly from the image with access to the visual context, GPT-4o correctly translates 'Exit' as 'Ausfahrt.'

the GPT-4o OCR model in our annotated set, there are nearly $10\times$ more insertions than substitutions in each example. When word-level substitutions occurred, they occurred at most twice per image, which both MT models were more easily able to recover from using context.

We do not observe downstream MT errors due to bounding box placement (Class III) in our sample. We note such errors may be more significant for models which process only the cropped region within bounding boxes, as in the example in Figure 8c an overly tight bounding box would cause the $g$ to look like an $a$ when cropped. This would be an important consideration if adapting recent visual text-based translation approaches for text in natural images (Salesky et al., 2021), as these models currently only process the region directly surrounding text.

### 3.4 Can multimodal models resolve contextual ambiguity?

Multimodal LLMs have access to both textual information contained in an image, as well as the visual context it is situated in. Cascaded OCR and MT, however, discards the visual information at translation time. Are LLMs able to use the broader visual context to resolve otherwise ambiguous translations?

It can be challenging to assess the degree to which multimodal models rely on different modalities for their predictions (Hessel and Lee, 2020), particularly for closed models without access to relative weights or the training data distribution for statistical priors. Here though direct translation with a multimodal LLM performs non-trivially, we still observe a performance gap to the same LLM

performing text translation from the reference transcripts without access to visual information: for the English→German language pair for example, 41.0 chrF and 18.8 BLEU vs. 36.9 chrF and 9.1 BLEU. Directly comparing quantitative results is not a perfect reflection of the task, because each model may get ambiguous examples wrong for different reasons. However, within the VISTRA test set we do observe examples where ambiguous source nouns are generated as only one possible translation with text input, but multiple senses with visual input. In our running example, 14 images in the benchmark contain the English word 'Exit'; in a model cascade when translating with access to text only, GPT-4o translates all 14 instances as 'Ausgang,' while with visual input only 5 instances are translated this way and 4 use a variant of 'Ausfahrt,' as illustrated in Figure 5. Particularly when used in conjunction with models trained from scratch, this benchmark may enable further analysis of attribution.

**Cautionary note on evaluation metrics**. Learned metrics such as COMET score paraphrases and synonyms highly, which typically leads to higher correlations with human judgments. However, for this task precisely that property may make them less reliable indicators of success. For example, returning to the motivational example in Figure 1, when translating the English sentence *'The exit is over there,'* both possible German translations *'Die Ausfahrt ist dort drüben'* and *'Der Ausgang ist dort drüben'* are given identical COMET scores (97.6) with either translation as the reference. Lexical metrics such as chrF and BLEU *do* reflect a mismatch to the reference here, and may be more reliable in this setting specifically for measuring correct visually-situated translation. For this

reason, and given the high proportion of examples marked as contextually dependent in our benchmark (Section 2.3), the COMET scores in Table 4 should likely viewed more cautiously than for other tasks. To properly evaluate contextually-dependent translations with multimodal input using a learned metric likely requires a new metric.

## 4 Related work

Translation of text in images has been strongly motivated by printed historical documents which require digitization (Afli and Way, 2016; Ignat et al., 2022) and PDF document translation (Zhang et al., 2023b; Hsu et al., 2024) with two column or more complex text layouts. In the absence of publicly available aligned and translated data sources, the majority of work in this space has created synthetic data for this task by rendering common machine translation corpora from sources from WMT14 (Mansimov et al., 2020; Tian et al., 2023; Niu et al., 2024; Lan et al., 2024). Ma et al. (2022) compared cascaded and direct models for in-image text translation with synthetic, cropped subtitles, and street-view images, but did not release their datasets. Lan et al. (2023) extended this work, studying auxiliary objectives for this task, and released a benchmark extending 5 Chinese OCR datasets with natural images with translations for Chinese→English. Ignat et al. (2022) perform similar analysis on the impact of OCR CER on downstream MT performance with the aim to see whether OCR'ed documents can be utilized for data augmentation for MT training with low-resource languages.

Similar to our task, multimodal translation uses auxiliary visual context to improve text translation, typically of image captions (Elliott et al., 2016; Specia et al., 2016; Elliott and Kádár, 2017; Elliott et al., 2017; Barrault et al., 2018; Li et al., 2022). Recent work has adapted pretrained model components into a single ViT model for this task (Gupta et al., 2023). As in our setting, it is challenging to assess the degree to which multimodal models make use of visual context in addition to text representations (Hessel and Lee, 2020); some studies investigating the usage of visual input in multimodal MT have found that do so primarily in the case of ambiguity or limited text input (Caglayan et al., 2019; Raunak et al., 2019) or provide regularization only (Wu et al., 2021).

Beyond machine translation, significant work has studied problems in text-centric visual processing such as document and table layout understanding through visual means (Long et al., 2022; Alonso et al., 2024; Zheng et al., 2024), OCR-free language understanding (Tanaka et al., 2021; Ye et al., 2023), and modeling language in screenshots (Kim et al., 2022; Lee et al., 2023; Gao et al., 2024). As multimodal LLMs become increasingly strong, analyzing their capabilities and limitations for text-rich image understanding (Zhang et al., 2023a, 2024; Li et al., 2024) and OCR (Liu et al., 2023) is a growing area. As we saw here, though they are strong general purpose models, there can remain a gap to task-specific models for complex and specialized tasks.

## 5 Conclusions

We introduce a benchmark, VISTRA, for visually-situated translation of English text in natural images to four target languages. We describe the dataset construction and composition. We benchmark multiple commonly used OCR models on VISTRA, both open-source and commercial, and evaluate cascaded OCR and MT performance. We present both quantitative result and create a taxonomy of common error classes, and investigate their impact on downstream MT. Finally, we assess direct image-to-text translation with a multimodal LLM, and show that it is able in some cases but not yet consistently to disambiguate possible translations with visual context. We show that this is an unsolved and challenging task even for strong commercial models. We hope that the creation and release of our benchmark, which is the first of its kind for these language pairs, will encourage further research in this direction.

## Limitations

Our dataset is limited in scale and language coverage to English text, with images predominantly taken in a single country (USA). The majority of photos were taken by a single photographer, which may lead to more consistent image quality and application of inclusion criteria, but likely also limits diversity through a locale bias to their surroundings. Transcriptions were performed by 3 individuals, and all checked by the same annotator for consistency, while translations were professionally done with a subset checked by a second annotator.

## References

Haithem Afli and Andy Way. 2016. Integrating optical character recognition and machine translation of historical documents. In *Proceedings of the Workshop on Language Technology Resources and Tools for Digital Humanities (LT4DH)*, pages 109–116, Osaka, Japan. The COLING 2016 Organizing Committee.

Iñigo Alonso, Eneko Agirre, and Mirella Lapata. 2024. PixT3: Pixel-based table-to-text generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6721–6736, Bangkok, Thailand. Association for Computational Linguistics.

Loïc Barrault, Fethi Bougares, Lucia Specia, Chiraag Lala, Desmond Elliott, and Stella Frank. 2018. Findings of the third shared task on multimodal machine translation. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 304–323, Belgium, Brussels. Association for Computational Linguistics.

Ozan Caglayan, Pranava Swaroop Madhyastha, Lucia Specia, and Loïc Barrault. 2019. Probing the need for visual context in multimodal machine translation. *ArXiv*, abs/1903.08678.

Yuning Du, Chenxia Li, Ruoyu Guo, Xiaoting Yin, Weiwei Liu, Jun Zhou, Yifan Bai, Zilin Yu, Yehua Yang, Qingqing Dang, and Haoshuang Wang. 2020. Pp-ocr: A practical ultra lightweight ocr system.

Desmond Elliott, Stella Frank, Loïc Barrault, Fethi Bougares, and Lucia Specia. 2017. Findings of the second shared task on multimodal machine translation and multilingual image description. In *Proceedings of the Second Conference on Machine Translation*, pages 215–233, Copenhagen, Denmark. Association for Computational Linguistics.

Desmond Elliott, Stella Frank, K. Sima'an, and Lucia Specia. 2016. Multi30k: Multilingual english-german image descriptions. *ArXiv*, abs/1605.00459.

Desmond Elliott and Ákos Kádár. 2017. Imagination improves multimodal translation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 130–141, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Tianyu Gao, Zirui Wang, Adithya Bhaskar, and Danqi Chen. 2024. Improving language understanding from screenshots. *ArXiv*, abs/2402.14073.

D. Gupta, S. Kharbanda, J. Zhou, W. Li, H. Pfister, and D. Wei. 2023. Cliptrans: Transferring visual knowledge with pre-trained models for multimodal machine translation. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2863–2874, Los Alamitos, CA, USA. IEEE Computer Society.

Jack Hessel and Lillian Lee. 2020. Does my multimodal model learn cross-modal interactions? it's harder to tell than you might think! In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 861–877, Online. Association for Computational Linguistics.

Benjamin Hsu, Xiaoyu Liu, Huayang Li, Yoshinari Fujinuma, Maria Nadejde, Xing Niu, Ron Litman, Yair Kittenplon, and Raghavendra Pappagari. 2024. M3T: A new benchmark dataset for multi-modal document-level machine translation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 499–507, Mexico City, Mexico. Association for Computational Linguistics.

Oana Ignat, Jean Maillard, Vishrav Chaudhary, and Francisco Guzmán. 2022. OCR improves machine translation for low-resource languages. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1164–1174, Dublin, Ireland. Association for Computational Linguistics.

R. Reeve Ingle, Yasuhisa Fujii, Thomas Deselaers, Jonathan Baccash, and Ashok Popat. 2019. A scalable handwritten text recognition system. *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 17–24.

Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. 2022. Ocr-free document understanding transformer. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII*, page 498–517, Berlin, Heidelberg. Springer-Verlag.

Zhibin Lan, Liqiang Niu, Fandong Meng, Jie Zhou, Min Zhang, and Jinsong Su. 2024. Translatotron-v(ison): An end-to-end model for in-image machine translation. In *Findings of the Association for Computational Linguistics: ACL 2024*. Association for Computational Linguistics.

Zhibin Lan, Jiawei Yu, Xiang Li, Wen Zhang, Jian Luan, Bin Wang, Degen Huang, and Jinsong Su. 2023. Exploring better text image translation with multimodal codebook. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3479–3491, Toronto, Canada. Association for Computational Linguistics.

Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *ICML*, pages 18893–18912.

Xiujun Li, Yujie Lu, Zhe Gan, Jianfeng Gao, William Yang Wang, and Yejin Choi. 2024. Text as images: Can multimodal large language models follow printed instructions in pixels?

Yi Li, Rameswar Panda, Yoon Kim, Chun-Fu Richard Chen, Rogerio Feris, David Cox, and Nuno Vasconcelos. 2022. Valhalla: Visual hallucination for machine translation. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5206–5216.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pretraining for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, and Xiang Bai. 2023. On the hidden mystery of ocr in large multimodal models. *ArXiv*, abs/2305.07895.

Shangbang Long, Siyang Qin, Dmitry Panteleev, A. Bissacco, Yasuhisa Fujii, and Michalis Raptis. 2022. Towards end-to-end unified scene text detection and layout analysis. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1039–1049.

Cong Ma, Yaping Zhang, Mei Tu, Xu Han, Linghui Wu, Yang Zhao, and Yu Zhou. 2022. Improving end-to-end text image translation from the auxiliary text translation task. *2022 26th International Conference on Pattern Recognition (ICPR)*, pages 1664–1670.

Elman Mansimov, Mitchell Stern, Mia Xu Chen, Orhan Firat, Jakob Uszkoreit, and Puneet Jain. 2020. Towards end-to-end in-image neural machine translation. In *Proceedings of the First International Workshop on Natural Language Processing Beyond Text*, pages 70–74.

Liqiang Niu, Fandong Meng, and Jie Zhou. 2024. UMTIT: Unifying recognition, translation, and generation for multimodal text image translation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 16953–16972, Torino, Italia. ELRA and ICCL.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia,

Pennsylvania, USA. Association for Computational Linguistics.

Ashok C Popat, Jonathan Michael Baccash, Karel Driesen, Patrick Michael Hurst, and Yasuhisa Fujii. 2017. Sequence-to-label script identification for multilingual ocr. In *Proceedings of the 14th International Conference on Document Analysis and Recognition (ICDAR)*.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Vikas Raunak, Sang Keun Choe, Quanyang Lu, Yi Xu, and Florian Metze. 2019. On leveraging the visual modality for neural machine translation. In *International Conference on Natural Language Generation*.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Elizabeth Salesky, David Etter, and Matt Post. 2021. Robust open-vocabulary translation from visual text representations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7235–7252, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

R. Smith. 2007. An overview of the tesseract ocr engine. In *Proceedings of the Ninth International Conference on Document Analysis and Recognition - Volume 02*, ICDAR '07, page 629–633, USA. IEEE Computer Society.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.

Lucia Specia, Stella Frank, Khalil Sima'an, and Desmond Elliott. 2016. A shared task on multimodal machine translation and crosslingual image description. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 543–553, Berlin, Germany. Association for Computational Linguistics.

Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. 2021. Visualmrc: Machine reading comprehension on document images. *ArXiv*, abs/2101.11272.

Yanzhi Tian, Xiang Li, Zeming Liu, Yuhang Guo, and Bin Wang. 2023. In-image neural machine translation with segmented pixel sequence-to-sequence model. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15046–15057, Singapore. Association for Computational Linguistics.

Zhiyong Wu, Lingpeng Kong, Wei Bi, Xiang Li, and Ben Kao. 2021. Good for misconceived reasons: An empirical revisiting on the need for visual context in multimodal machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6153–6166, Online. Association for Computational Linguistics.

Yukang Yang, Dongnan Gui, YUHUI YUAN, Weicong Liang, Haisong Ding, Han Hu, and Kai Chen. 2023. Glyphcontrol: Glyph conditional control for visual text generation. In *Advances in Neural Information Processing Systems*, volume 36, pages 44050–44066. Curran Associates, Inc.

Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He, Xin Lin, and Fei Huang. 2023. UReader: Universal OCR-free visually-situated language understanding with multimodal large language model. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2841–2858, Singapore. Association for Computational Linguistics.

Ruiyi Zhang, Yufan Zhou, Jian Chen, Jiuxiang Gu, Changyou Chen, and Tongfei Sun. 2024. Llava-read: Enhancing reading ability of multimodal language models. *ArXiv*.

Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tongfei Sun. 2023a. Llavar: Enhanced visual instruction tuning for text-rich image understanding. *ArXiv*, abs/2306.17107.

Zhiyang Zhang, Yaping Zhang, Yupu Liang, Lu Xiang, Yang Zhao, Yu Zhou, and Chengqing Zong. 2023b. LayoutDIT: Layout-aware end-to-end document image translation with multi-step conductive decoder. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10043–10053, Singapore. Association for Computational Linguistics.

Mingyu Zheng, Xinwei Feng, Qingyi Si, Qiaoqiao She, Zheng Lin, Wenbin Jiang, and Weiping Wang. 2024. Multimodal table understanding. *ArXiv*, abs/2406.08100.

# A  Annotation Interfaces



Figure 6: Text annotation interface for VISTRA benchmark.

Figure 7: Translation annotation interface for VISTRA benchmark.

# B  Examples of each OCR error class

Here we show an illustrative example of each OCR error class described in Table 2 from the VISTRA benchmark, with the model which produced each output.



(a) CLASS I: Undetected text



(b) CLASS II: Text hallucination



(c) CLASS III: Bounding box error



(d) CLASS IV: Grouping error

Figure 8: Examples of each OCR error class from Table 2.

(e) CLASS V: Punctuation error



(f) CLASS VI: Spacing error



(g) CLASS VII: Character-level substitution



(h) CLASS VIII: Word-level substitution

Figure 8: Examples of each OCR error class from Table 2 (cont.)

# Analysing Translation Artifacts: A Comparative Study of LLMs, NMTs, and Human Translations

Fedor Sizov[1]    Cristina España-Bonet[2]    Josef van Genabith[1,2]
Roy Xie[3]    Koel Dutta Chowdhury[1]

[1]Saarland University, Saarland Informatics Campus    [2]DFKI GmbH    [3]Duke University

{cristinae, josef.van_genabith}@dfki.de,
sife00002@stud.uni-saarland.de, ruoyu.xie@duke.edu, koeldc@lst.uni-saarland.de

## Abstract

Translated texts exhibit a range of characteristics that make them appear distinct from texts originally written in the same target language. With the rise of Large Language Models (LLMs), which are designed for a wide range of language generation and understanding tasks, there has been significant interest in their application to Machine Translation. While several studies have focused on improving translation quality through fine-tuning or few-shot prompting techniques, there has been limited exploration of how LLM-generated translations qualitatively differ from those produced by Neural Machine Translation (NMT) models, and human translations. Our study employs explainability methods such as Leave-One-Out (LOO) and Integrated Gradients (IG) to analyze the lexical features distinguishing human translations from those produced by LLMs and NMT systems. Specifically, we apply a two-stage approach: first, classifying texts based on their origin —whether they are original or translations— and second, extracting significant lexical features (highly attributed input words) using post-hoc interpretability methods. Our analysis shows that different methods of feature extraction vary in their effectiveness, with LOO being generally better at pinpointing critical input words and IG capturing a broader range of important words. Finally, our results show that while LLMs and NMT systems can produce translations of a good quality, they still differ from texts originally written by native speakers. We find that while some LLMs more closely resemble human translations, traditional NMT systems show distinct differences, particularly in their use of linguistic features. [1]

## 1 Introduction

The rapid development of large language models (LLMs) (Radford et al., 2019; Raffel et al., 2020a; Touvron et al., 2023; Lu et al., 2024; Team et al., 2024a; Groeneveld et al., 2024; Alves et al., 2024) has significantly advanced natural language processing (NLP), also in the domain of Machine Translation (MT) (Zhang et al., 2023; Zhu et al., 2024) with studies covering various approaches such as document-level literary translation (Karpinska and Iyyer, 2023), paragraph-level post-editing with LLMs (Thai et al., 2022), sentence-level translation (Vilar et al., 2022; Jiao et al., 2023), examining hallucinations in LLM-generated translations (Guerreiro et al., 2023), and leveraging LLMs for evaluation (Kocmi and Federmann, 2023). These efforts reflect the ongoing shift toward exploring how well LLMs perform MT compared to traditional NMT systems.

Although previous work (Zhu et al., 2024; Vilar et al., 2022; Raunak et al., 2023) have explored how LLMs and traditional Neural Machine Translation (NMT) systems develop translation capabilities, as well as the qualitative differences in their outputs and the factors that impact their performance, a critical gap remains: the comparison of translations generated by LLMs and NMT models to those produced by human translators (HT) and texts originally written by native speakers in the target language. This comparison raises questions about translation divergence, as reflected in surface-level (structural) differences in translations arising from cross-linguistic variations or translator preferences (Luo et al., 2024).

Such divergences are well-documented in human translations (HT), where translators often make structural choices that vary significantly from the text originally written in the target language (Deng and Xue, 2017; Nikolaev et al., 2020). In contrast, traditional NMT outputs typically exhibit less diversity and more literal translations, lacking significant structural variation (Freitag et al., 2020; Bizzoni et al., 2020). Similarly, Vyas et al. (2018); Briakou and Carpuat (2020) focus on identifying

---

semantic divergences in translations that are not fully equivalent to the original source texts. Recent findings, however, indicate that LLMs tend to produce translations that are less literal compared to NMT models (Vilar et al., 2022; Raunak et al., 2023), suggesting that LLMs may bridge the gap between the rigid literalness of NMT models and the flexibility of human translations. Understanding these divergences is crucial for advancing translation technologies and ensuring their responsible and effective use. Specifically, this leads us to investigate the following research questions: **how do LLMs, NMT models, and HT outputs differ in their translations**, and **what methods can effectively identify these differences?**

To answer these questions, we conduct a systematic comparison of LLM, NMT, and HT translations using explainability techniques (Lundberg and Lee, 2017; Rajagopal et al., 2021; Yin and Neubig, 2022; Wu et al., 2023), namely Leave-One-Out (LOO) (Li et al., 2016) and Integrated Gradients (IG) (Sundararajan et al., 2017). Specifically, we use a two-stage approach: first, we classify texts in the same target language based on their origin —whether they are original texts (O) written by native speakers or translations (T), whether human or automated. Next, we apply post-hoc interpretability methods to extract key features that contribute to these classifications. Our analysis focuses on identifying whether the most important features for O/T classification are consistent across LLM-based, NMT-based, and human translation outputs.

To understand these distinctions, we perform two analyses: $(i)$ Feature Overlap Analysis: we calculate the average intersection of the top most important lexical features used across different translation systems to classify O/T, focusing on how much the most important features identified by explainability techniques overlap across LLM, NMT, and HT systems, and $(ii)$ Feature Frequency Analysis: we analyse the frequency distribution of these key linguistic features within each translation system.

Our findings show that while many LLMs and NMT systems produce good translations, they still differ from content originally written by native speakers. LLMs like Aya-101-13B and TowerInstruct-7B-v0.2 exhibit alignment with traditional NMT models, such as DeepL and NLLB-600M, regarding O/T classification accuracy compared to content originally authored in the target language. Overall, our results confirm that NMT

translations are more readily distinguishable from originals, with traditional NMT systems generally outperforming LLMs in translation quality and consistency. At the same time, human-generated translations remain distinctly different from those produced by machines.

Using explainability methods, we identified the key features that differentiate translations produced by LLMs, NMT systems, and human translators. Our findings suggest that LOO is generally better at pinpointing the most critical single feature, while IG is more effective when considering a broader range of important features. Moreover, our analysis shows that LLMs like Gemma-7B and TowerInstruct-7B-v0.2 often align closely with NMT systems such as M2M-100-418M and DeepL in their lexical feature selection during translation. Finally, our findings show that LLMs generally exhibit PoS patterns more aligned with HT than NMT models, particularly in the use of adverbs and auxiliary verbs. However, human translations consistently exhibit lower overlap with certain linguistic features from both LLMs and NMT systems, indicating that despite some shared patterns, human translations retain a unique quality.

The paper is structured as follows: Section 2 outlines our experimental design, and Sections 2.1 and 2.2 detail the data and models used in our study. Section 3 discusses our strategies for evaluation of translation quality and methods we employ for extracting important distinctive features of original and translated texts, while Section 4 examines the differences in classification features between LLMs, NMT systems, and human translations. Finally, Section 5 concludes the paper.

## 2 Experimental Design

To identify important explanations with respect to O/T classification in the outputs of translation systems, we apply explainability methods to each sentence and generate attribution scores for the tokens. Below, we describe the methods used to produce these attribution scores.

**Leave-One-Out (LOO).** We use LOO (Li et al., 2016), a popular model-agnostic feature attribution technique, to compute the attribution score for each token $x_i$ in an input sentence $X$ with respect to the model's prediction $\hat{y}$. Let $w_{[CLS]}$ be the final layer representation of the "[CLS]" token for $X$. During inference, the method processes the input through ReLU, affine, and softmax layers to produce a prob-

ability distribution over the outputs. For each token $x_i$, LOO measures the change in probability when $x_i$ is excluded from the input $X$. Higher change in probability indicates that the token $x_i$ is more influential in the model's prediction:

$$\ell = \text{softmax}(\text{affine}(\text{ReLU}(w_{[\text{CLS}]})))$$

$$\ell_i = \text{softmax}(\text{affine}(\text{ReLU}(w_i)))$$

$$\nabla_i = \ell - \ell_i$$

where $w_i$ represents the final layer output of the "[CLS]" token when the token $x_i$ is removed from the input sequence $X$.

**Integrated Gradients (IG).** Sundararajan et al. (2017) propose this technique for attributing a neural network's output to its input features by computing the integral of the gradients of the model's prediction with respect to the inputs along a path from a baseline to the actual input. The attribution for a feature $x_i$ is given by:

$$\text{IG}_i = (x_i - x_i^0) \cdot \int_0^1 \frac{\partial f(x^0 + \alpha \cdot (x - x^0))}{\partial x_i} \, d\alpha$$

where $x_i^0$ is the baseline input and $f$ is the model's prediction function.

In this work, IG is used to compute attribution scores for each token $x_i$[2] in $X$. IG provides scores between $-1$ and $1$ for each embedding dimension of the token $x_i$, where $1$ and $-1$ represent maximum influence towards labels 1 (T) and 0 (O), and scores near zero indicate minimal impact.

## 2.1 Data

We use the Monolingual German dataset from the Multilingual Parallel Direct Europarl (MPDE) featuring annotated paragraphs from the proceedings of the European Parliament (Amponsah-Kaakyire et al., 2021). The dataset includes both the original texts and their translations. Each paragraph, averaging 80 tokens, is labeled to indicate whether it is an original or a translation. Since most NMT systems operate on sentence level, we split each paragraph into sentences, which we later use for our work.

However, in MPDE, paragraphs of German sources typically contain more sentences than their

English translations.[3] To address this imbalance, we remove certain amount of German source sentences, creating a training set with an equal number of original and translated sentences (97,108 in the training set and 20,744 in the test set).

To further perform evaluation of translation quality, we need a clear one-to-one correspondence between source sentence, human-translated sentence and the automatically translated sentence. As mentioned above, not every paragraph of the MPDE dataset has the same number of sentences in its German source and in its English translation. We have composed a subset of MPDE consisting only of those sentences whose paragraphs have an equal number of German and English sentences. This subset contains 38,035 sentences.

**Pre-processing.** To ensure that the explanation methods work efficiently, we tokenize and truecase our data.[4] Both are performed using Moses scripts (Koehn et al., 2007).

## 2.2 Models

We report O/T classification and translation quality results on a wide selection of some of the best-performing models, both commercial and open-source models:

- **DeepL Translator**: a state-of-the-art commercial NMT system.[5]
- **Google Translate**: Likely the most widely used commercial NMT system.[6]
- **M2M-100-418M** (Fan et al., 2020): A large multilingual NMT model trained on 2,200 translation directions, enabling many-to-many translation across 100 languages. We use the base version.
- **MADLAD-400** (Kudugunta et al., 2023): A multilingual NMT model based on the T5 architecture (Raffel et al., 2020b), with 3 billion parameters, trained on 1 trillion tokens across 450 languages using publicly available data.
- **NLLB-600M** (Costa-jussà et al., 2022): It represents the current state-of-the-art NMT system,

---

[2] Token $x_i$ may refer to either a whole word or its subunits, as the WordPiece tokenizer (Song et al., 2021) splits words into subunits. To compute the attribution score at the word level, we average the attributions of its subunits.

[3] This is due to the fact that the translations of paragraphs are not aligned sentence-wise. While the original paragraph may have i sentences, one translation may have $j$ sentences and another $k$.

[4] As further we need, for example, to analyze lexical overlaps, it is important that we do not miss out on words because of punctuation or case

[5] https://www.deepl.com/en/translator (accessed on August 16, 2024)

[6] https://translate.google.com/?sl=de&tl=en&op=translate (accessed on August 13, 2024)

| System | O/T Classification Accuracy (%) | AEM | |
|---|---|---|---|
| | | COMET | BLEU |
| HT | 0.79 | | |
| DeepL | 0.86 | **0.85** | **34.85 ± 0.19** |
| Google Translate | 0.92 | 0.79 | 24.17 ± 0.16 |
| M2M-100-418M | 0.91 | 0.81 | 25.94 ± 0.16 |
| MADLAD-400-MT | 0.91 | 0.69 | 16.37 ± 0.18 |
| NLLB-600M | 0.83 | 0.79 | 27.35 ± 0.19 |
| LLaMAX-3.1-8B-Alpaca | **0.94** | 0.81 | 15.43 ± 0.13 |
| TowerInstruct-7B-v0.2 | 0.83 | **0.84** | 33.35 ± 0.18 |
| Aya-101-13B | 0.86 | 0.83 | 25.35 ± 0.16 |
| Gemma-7B | 0.89 | 0.83 | 27.53± 0.19 |
| Llama-3.1-IT-8B | 0.90 | 0.82 | 26.91 ± 0.17 |

Table 1: Performance metrics for various systems including classification accuracy and automatic MT evaluation metrics (COMET and BLEU). The highest scores are highlighted in bold.

scaling up to 200 languages. We experiment with the distilled version with 600M parameters.

In addition to the NMT systems listed above, we pick three well-known and high-performing open-source LLMs and use them for prompt-based translation without any prior fine-tuning (see Appendix A for the prompt templates):

- **LLaMAX-3.1-8B-Alpaca** (Lu et al., 2024) is an open-source instruction-following language model with 8 billion parameters. It is fine-tuned from the LLaMA model (Taori et al., 2023) and supports 102 languages through continual pre-training, incorporating 52,000 Self-Instruct English instruction examples (Wang et al., 2023).

- **Llama-3.1-IT-8B** (Dubey et al., 2024): The Meta Llama 3.1 collection includes multilingual LLMs. This 8B parameter model is pretrained and instruction-tuned for text generation, optimized for multilingual dialogue.

- **TowerInstruct-7B-v0.2** (Alves et al., 2024): A language model based on LLaMA 2 (Touvron et al., 2023), using a diverse dataset of 20 billion tokens from monolingual sources in ten different languages.

- **Aya-101-13B** (Üstün et al., 2024): A 13-billion-parameter mT5 (Xue et al., 2021) multilingual model trained on instructions in 101 languages, exceeding the coverage of earlier open-source models (Lai et al., 2023; Muennighoff et al., 2022; Le Scao et al., 2023).

- **Gemma-7B** (Team et al., 2024b) is a lightweight open-source LLM developed by Google Deep-Mind. It has been instruction-tuned to respond to prompts in a conversational manner.

## 3 Evaluation

### 3.1 O/T Classification

We follow Dutta Chowdhury et al. (2022) to perform binary classification between original and translated (O and T) sentences. We use the XLM-RoBERTa base model (Conneau et al., 2020) with a softmax classifier applied to the [CLS] token of the sentence embeddings. We freeze hyperparameters and weights of the pre-trained encoder, and train the classifier for 10 epochs on each sentence with batch size of 16 and learning rate of $2 \times 10^{-5}$. All experiments are performed using NVIDIA V100 or A100 GPUs.

**Results.** The linear O/T classifiers show high accuracies (>80%) for all models (Table 1). We find that the automatically translated sentences, for both NMTs and LLMs, are always identified with higher accuracy than the human-translated ones. This finding corroborates the hypothesis that automatically translated texts are more readily distinguishable in classification tasks than those translated by humans (Ilisei et al., 2010; Rubino et al., 2016; Pylypenko et al., 2021).

### 3.2 Translation Quality

To assess translation quality, we utilise two automatic evaluation metrics (AEM): BLEU (Papineni et al., 2002) as implemented in SacreBLEU[7] (Post, 2018) and COMET (Rei et al., 2022).[8] BLEU relies on word n-gram similarity, whereas COMET

---

[7]BLEU signature: nrefs:1|case:mixed|eff:no|tok:13a| smooth:exp|version:2.0.0

[8]Unbabel/wmt22-comet-da, see https://github.com/ Unbabel/COMET

| System | LOO | | | IG | | |
|---|---|---|---|---|---|---|
| | top-1 | top-3 | top-5 | top-1 | top-3 | top-5 |
| HT | 0.64 | 0.66 | 0.66 | 0.51 | 0.56 | 0.57 |
| DeepL | 0.60 | 0.73 | 0.72 | 0.53 | 0.61 | 0.71 |
| Google Translate | **0.78** | 0.70 | 0.76 | 0.50 | 0.50 | **0.83** |
| M2M-100-418M | 0.57 | 0.70 | 0.76 | 0.57 | 0.75 | 0.75 |
| NLLB-600M | 0.50 | 0.73 | 0.69 | 0.58 | 0.50 | 0.71 |
| TowerInstruct-7B-v0.2 | 0.54 | 0.70 | 0.74 | 0.51 | 0.55 | 0.54 |
| Aya-101-13B | 0.53 | 0.69 | 0.76 | 0.53 | 0.72 | 0.68 |
| Gemma-7B | 0.54 | 0.65 | 0.63 | 0.55 | 0.55 | 0.53 |
| Llama-3.1-IT-8B | 0.50 | 0.73 | 0.76 | 0.51 | 0.64 | 0.65 |
| **Mean** | 0.58 | 0.70 | 0.72 | 0.53 | 0.60 | 0.66 |

Table 2: Performance of the sufficiency classifier across different ranks (top-1, top-3, top-5) using LOO and IG methods for HT, NMT, and LLM systems. The highest scores for each method are highlighted in teal (LOO) and gray (IG), with the highest scores boldfaced to highlight the strengths of each method.

is a semantic metric built upon the XLM-R architecture.

**Results.** Table 1 shows that across different models, COMET scores remain relatively stable, while BLEU scores show greater fluctuation. DeepL stands out as the top performer, achieving the highest scores in both COMET (0.85) and BLEU (34.85). TowerInstruct-7B-v0.2 also performs well, particularly in COMET, reflecting high translation quality. Two systems, LLaMAX-3.1-8B-Alpaca and MADLAD-400-MT, exhibit poor translation quality. The high number of translation errors could skew the explainability results, focusing on these mistakes rather than models' intrinsic characteristics. Therefore, we exclude these models for further experiments. We perform a correlation analysis, and find no significant correlation between translation quality and O/T classification accuracy. See Appendix C for more details.

### 3.3 Do explanations capture sufficient information?

Understanding the effectiveness of model predictions often relies on the quality of explanations derived from those models. In this context, an explanation refers to the rationale behind a model's predictions, specifically identifying the *input tokens (features)* that most significantly influence the classification outcome. We follow the approach outlined by Xie et al. (2024) to evaluate the sufficiency of these explanations, as defined by Jacovi et al. (2018) and Yu et al. (2019). Sufficiency refers to the average change in predicted class probability when only the top $k$ influential tokens are retained.

This metric assesses how well the top $k$ attributions explain the model's predictions, ultimately determining whether these explanations faithfully represent the model's decision-making process.

Previous research (Amponsah-Kaakyire et al., 2022) has shown that feature attribution including IG can be used to identify input tokens that are particularly important to O/T classification results for original texts and human translations.

However, whether this holds true across different types of translations, such as those generated by large language models (LLMs) or neural machine translation systems (NMT), remains underexplored. Bizzoni et al. (2020) investigated this problem using PoS perplexity scores and syntactic dependency lengths. More recently, Luo et al. (2024) systematically investigate the differences in the distribution of translation divergences between HT and MT through a large-scale, fine-grained comparative analysis, focusing on morphosyntactic variations. In contrast, our approach investigates lexical (words and PoS) differences by analysing explanations from O/T classifiers.

Our goal is to identify the key features that set apart translation artifacts produced by LLMs, NMT, and HTs from the text originally authored in the target language. To evaluate the sufficiency of our methods—specifically Leave-One-Out (LOO) and IG—we separately extract the top $k$ tokens with the highest attribution scores for each sentence in the training set (see Section 2.1). We then construct datasets with sentences consisting only of these top $k$ tokens while maintaining the same labels. O/T classifiers are then trained on these datasets,
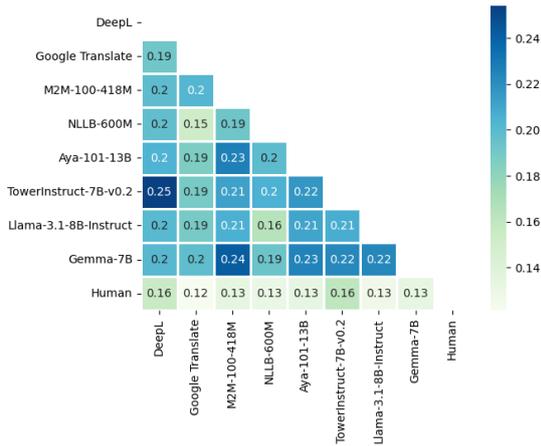
Figure 1: Level of intersection between top-5 most important explanations across different translation methods using LOO method.
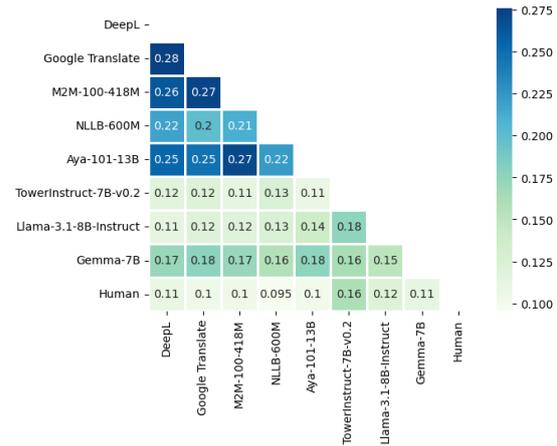


Figure 2: Level of intersection between top-5 most important explanations across different translation methods with IG.

where $k = \{1, 3, 5\}$, and we subsequently assess the classifiers' accuracy on the test set (Table 2) [9].

### 3.3.1 Sufficiency

If we can maintain high accuracy of O/T classifier using only the $k$ tokens with the highest attribution scores, this indicates that the explainability methods (LOO and IG) work as intended, allowing us to efficiently identify important differences between translations and originally authored sentences in the target language.

**Results.** Table 2 shows that high accuracy for O/T is consistently maintained for the top $k$ tokens with the highest attribution scores, indicating that the explainability methods (LOO and IG) function as intended. On average, as the number of tokens increases, we see an improvement in the sufficiency scores, indicating that the features we are extracting are indeed important.

Moreover, LOO is able to achieve much higher sufficiency score on top-1 tokens from certain model outputs as compared to IG, suggesting that LOO may be more effective at pinpointing the most critical token for classification. The reason for that might be that Leave-One-Out (LOO) directly removes each word and measures the impact on model prediction, giving a more precise attribution score. In contrast, Integrated Gradients (IG) require pooling attributions across the dimensions of an embedding and averaging attributions across subwords when a word is split into pieces, which

---

[9] We modified the train set for the sufficiency experiment but left the test set unchanged to ensure fair evaluation.

may provide better performance in context, but lower it when focusing on a single word.

The LOO method achieves its highest top-1 sufficiency score of $0.78$ across all models for Google Translate, underscoring its potential effectiveness in identifying essential tokens. In contrast, the IG method records its highest top-5 sufficiency score of $0.83$ for the same translation system, showcasing its strength in capturing significant features across a broader range of tokens.

## 4 Feature Analysis of LLM, NMT, and Human Translation

### 4.1 Feature Overlap Analysis

We conduct an intersection analysis of linguistic features (input tokens), focusing specifically on sentences for which we can establish a one-to-one correspondence between outputs of different translation systems. For these sentences, we apply both LOO and IG using previously trained O/T classifiers for HT, NMT, and LLM datasets. This process enables us to compute attribution scores for individual tokens within each sentence. Using these scores, we extract the top-$k$ most important tokens ($k = 1, 3, 5$) for each sentence.

Following this, we calculate the intersection between the LOO and IG results for different translation systems using the Jaccard Similarity Coefficient, which represents the percentage of common tokens and takes a value from 0 to 1. A high intersection among the top-$k$ tokens indicates robust features (tokens) that are consistently identified as important across different translation models.

Conversely, if the intersection between systems and/or human translations is low, it indicates that the translations exhibit different features. Figure 1 presents the pairwise Jaccard values for the top-5 features derived from the Leave-One-Out (LOO) method. Each cell quantifies the degree of overlap between the top features of two different translation systems, with darker shades representing higher overlaps. Notably, the highest intersection is observed between TowerInstruct-7B-v0.2 and DeepL, with an overlap of $0.25$, suggesting a strong similarity in the features identified for these models.

Another substantial intersection occurs between Gemma-7B and M2M-100-418M at $0.24$, indicating considerable alignment in their outputs. In contrast, human-generated content shows relatively lower intersections with machine models, such as $0.16$ with TowerInstruct-7B-v0.2 and DeepL and $0.13$ with M2M-100-418M, underscoring the unique nature of human translations compared to machine-generated translations.

Similarly, Figure 2 shows the pairwise Jaccard values for the top-5 features (tokens) obtained using Integrated Gradients (IG). The most notable overlap is between Google Translate and DeepL, with a significant intersection of $0.28$, demonstrating a strong similarity in their feature selections. A notable intersection of $0.27$ is observed between M2M-100-418M and both Aya-101-13B and Google Translate, suggesting that these models yield quite similar results. The lower intersection of $0.11$ between TowerInstruct-7B-v0.2 and Aya-101-13B emphasizes the differences in their outputs. The intersection with human translation identified by IG is notably highest for TowerInstruct-7B-v0.2, at a value of $0.16$.

The combined results suggest that while certain LLMs, like Aya-101-13B and TowerInstruct-7B-v0.2, closely align with NMT models such as M2M-100-418M and DeepL in their feature selection, others retain unique classification features. Furthermore, there are notable differences in how closely these models align with human translations, with TowerInstruct-7B-v0.2 demonstrating the highest similarity to HT as shown by both LOO and IG.

## 4.2 Feature Frequency Analysis

We examine the frequency of different Part of Speech (PoS) tags across translation systems, focusing on the top $k$ features flagged by LOO/IG for each sentence. For each system, we group sentences – both human and machine translations – into predefined sentence length bins. These bins are divided into ranges (e.g., 0-10, 10-15, 15-20 words), and for each, we calculate and normalize the frequency of the identified features based on the total number of sentences in that bin. This helps us compare trends in PoS distribution as sentence length increases. We are examining trends for the 9 most common PoS.

To ensure the reliability of our measurements, we account for the margin of error (standard deviation) obtained through bootstrapping by sub-sampling each bin 1,000 times while maintaining the PoS distribution within each sentence. In the graphs we show the standard deviation with shading. Figure 3 illustrates variations in PoS distribution, showing nine subplots for adverbs (ADV), verbs (VERB), determiners (DET), auxiliary verbs (AUX), nouns (NOUN), pronouns (PRON), adjectives (ADJ), adpositions (ADP), and punctuations (PUNCT).

For ADV, most models—both NMT and LLM—use fewer adverbs than HT. However, Llama-3.1-8B demonstrates frequencies that are closer to HT as sentence length increases, while TowerInstruct-7B-v0.2 diverges with longer sentences. NMT models like M2M and Google Translate underproduce ADV compared to HT, whereas DeepL aligns more closely with HT and tends to overproduce ADV with longer sentences.

ADP use in HT increases with sentence length, and most NMT and LLM models follow this trend, although models like Google Translate show slightly lower frequencies in longer sentences. Pylypenko et al. (2021) find that the relative frequencies of ADV and ADP in PoS tagging are strong indicators of translationese in HT.

For VERB, both HT and most NMT and LLM models maintain a steady frequency, though the models generally underproduce compared to the human translation trend. For DET, HT usage slightly increases with sentence length, while all LLM and NMT models, except DeepL, tend to use determiners more frequently.

In the case of PRON, most models tend to align with the human trend for shorter sentences. However, as sentence length increases, their frequencies start to deviate from each other. NLLB-600M demonstrates a substantially higher frequency than human translations across all sentence lengths.

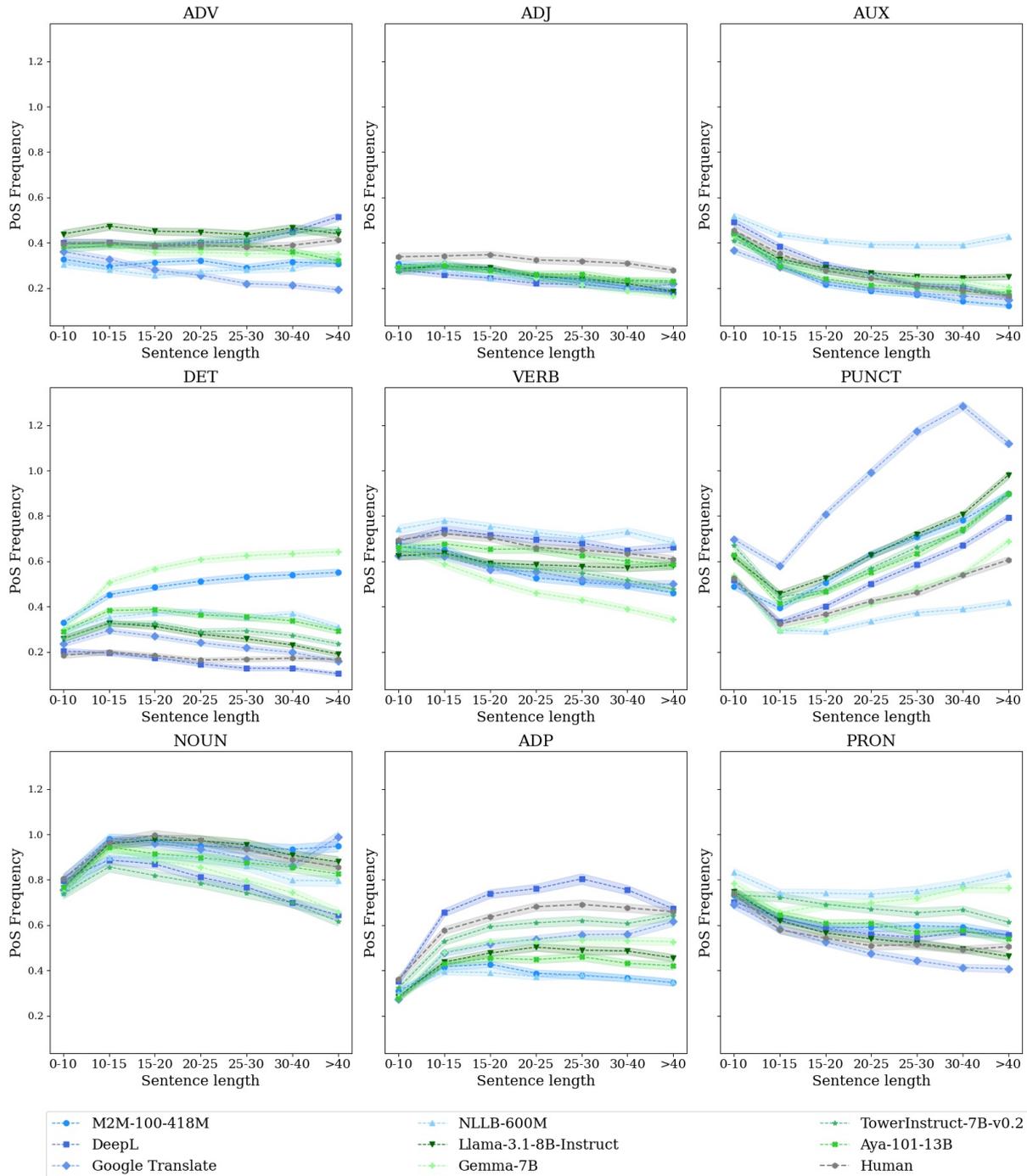In ADJ usage, HT remains relatively stable,

Figure 3: The frequency of the top PoS categories flagged by LOO across different sentence length bins. The x-axis of each subplot represents sentence length, divided into ranges (0-10, 10-15, 15-20, etc.), and the y-axis shows PoS frequency, indicating how often each PoS occurs in sentences of different lengths.

showing a slight decrease as sentence length increases. All NMT and LLM models exhibit lower adjective frequencies overall, with their trends being extremely similar across all sentence lengths.

For AUX, HT demonstrates a consistent decline as sentence length increases. Most NMT models follow this trend, except for NLLB-600M, which shows significantly higher AUX usage.

Similarly, Llama-3.1-8B-Instruct exhibits slightly higher AUX frequencies compared to HT. The frequency of NOUN usage is maximal for sentences of length 10-15 and then consistently decreases for longer sentences. HT and most models seem to follow this trend, except for two NMTs (M2M-100 and Google Translate), which tend to overproduce nouns in very long sentences. For HT

and NMT/LLM, the frequency of PUNCT usage in sentences of length 10-15 is lower than in shorter sentences, although there is an increasing trend for sentences longer than 15. Google Translate exhibits notably higher PUNCT frequencies than all other models and HT, although its usage declines in very long sentences.

Overall, LLMs exhibit PoS patterns (for 6 out of 9 tags) that closely align with human translations, whereas NMT models show greater deviations, particularly regarding PUNCT. NMT models tend to underproduce ADV, and for some other parts of speech (PoS) like ADP or PRON, they show significant divergence. In contrast, LLMs exhibit stronger agreement in trends and align more closely with HT, although they still demonstrate some overuse in short sentences. Both NMTs and LLMs underproduce ADJ compared to HT, particularly in longer sentences. LLMs better mimic human usage in ADV and AUX frequencies, especially in longer sentences. Appendix B displays the frequency plots of the top PoS categories identified by Integrated Gradients (IG) across various sentence-length bins.

## 5 Conclusion

In this work, we systematically explore the translation divergences between LLMs, NMTs, and human translations. Our key findings show distinct differences in how these systems approach translation, despite advancements in LLMs that allow them to produce high-quality outputs. We find that while LLMs often exhibit translation patterns more similar to human translations compared to traditional NMT models, they still diverge from originally authored text in the same language. Overall, we find that automatically translated sentences from both NMTs and LLMs are consistently identified with higher accuracy in O/T classification tasks than human-translated ones. This supports the hypothesis that machine-translated texts are more easily distinguishable from original texts than those translated by humans (Rubino et al., 2016; Pylypenko et al., 2021).

To better understand the distinctions between translations produced by LLMs and NMTs compared to human translations, we employ Leave-One-Out and Integrated Gradients explanation methods to extract and analyze lexical features identified by translation classifiers. Our findings indicate that even when using a sufficiency-based approach, we can recover a significant amount of

O/T classification accuracy. This demonstrates that these features are effective in distinguishing between automatic and human translations.

Further, our results indicate that sufficiency-based approach is particularly effective at identifying single critical features, while Integrated Gradients (IG) capture a broader range of important features. Interestingly, we observe that certain LLMs align closely with NMT systems in their feature selection, demonstrating similarities in their approaches. However, human translations consistently exhibit lower overlap with both LLM and NMT outputs, particularly regarding crucial features like punctuation and specific PoS.

Furthermore, our frequency analysis of PoS tags reveals that LLMs align more closely with HT in their usage, especially in terms of adverbs, and auxiliary verbs, while NMT models tend to overproduce specific tags in shorter sentences. This suggests that LLMs, although not perfect, are making strides in mimicking human translation patterns. Our findings highlight the characteristics that define the outputs of various translation systems. However, despite advances in machine translation, human translations continue to display distinctive characteristics, particularly in their nuanced use of linguistic features, making them less prone to the artifacts seen in machine-generated texts.

## Limitations

**Limitations of Lexical Features.** The results presented in this study rely entirely on the lexical features derived from Leave-One-Out (LOO) and Integrated Gradients (IG), which may fall short of capturing the intricacies of translation quality. Moreover, translation artifacts can arise at both syntactic and semantic levels (Bizzoni et al., 2020; Briakou and Carpuat, 2020), aspects that this research does not address. This leaves an exploration of these dimensions to future work.

**Prompting Choice.** Prompting has demonstrated varying sensitivity to the choice of templates and examples (Zhao et al., 2021). In machine translation (MT), prior studies have used different templates (Brown et al., 2020; Chowdhery et al., 2023; Wei et al., 2021). In our work, we reevaluate these templates to determine the optimal one. However, the format and wording of the prompt significantly influence how the LLM comprehends the task and performs translation, potentially impacting our findings, which we leave for future exploration.

**Stability of Model Outputs.** Additionally, we have assumed that the output of a specific model would remain stable throughout the analysis. However, LLMs are frequently updated, which can lead to changes in their writing style and coherence. Such variations might cause explainability methods to underperform, exacerbating the issues discussed in this work.

**Constraints of Sentence-Level Analysis.** Most NMT models utilized in this study function effectively at the sentence level, necessitating that we translate individual sentences for both NMTs and LLMs to ensure consistency. Thus, our sentence-based analysis with LLMs is also a limiting factor, as it restricts our ability to capture broader contextual nuances (Koneru et al., 2024). This would entail expanding our analysis beyond sentence-level assessments.

## Acknowledgments

## References

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks.

Kwabena Amponsah-Kaakyire, Daria Pylypenko, Cristina España-Bonet, and Josef van Genabith. 2021. Do not rely on relay translations: Multilingual parallel direct Europarl. In *Proceedings for the First Workshop on Modelling Translation: Translatology in the Digital Age*, pages 1–7, online. Association for Computational Linguistics.

Kwabena Amponsah-Kaakyire, Daria Pylypenko, Josef Genabith, and Cristina España-Bonet. 2022. Explaining translationese: why are neural classifiers better and what do they learn? In *Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 281–296, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Yuri Bizzoni, Tom S Juzek, Cristina España-Bonet, Koel Dutta Chowdhury, Josef van Genabith, and Elke Teich. 2020. How human is machine translationese? comparing human and machine translations of text and speech. In *Proceedings of the 17th International conference on spoken language translation*, pages 280–290.

Eleftheria Briakou and Marine Carpuat. 2020. Detecting fine-grained cross-lingual semantic divergences without supervision by learning to rank. *arXiv preprint arXiv:2010.03662*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Dun Deng and Nianwen Xue. 2017. Translation divergences in Chinese–English machine translation: An empirical investigation. *Computational Linguistics*, 43(3):521–565.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Koel Dutta Chowdhury, Rricha Jalota, Cristina España-Bonet, and Josef Genabith. 2022. Towards debiasing translation artifacts. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3983–3991, Seattle, United States. Association for Computational Linguistics.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation.

Markus Freitag, David Grangier, and Isaac Caswell. 2020. Bleu might be guilty but references are not innocent. *arXiv preprint arXiv:2004.06063*.

Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. 2024. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*.

Nuno M Guerreiro, Duarte M Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André FT Martins. 2023. Hallucinations in large multilingual translation models. *Transactions of the Association for Computational Linguistics*, 11:1500–1517.

Iustina Ilisei, Diana Inkpen, Gloria Corpas Pastor, and Ruslan Mitkov. 2010. Identification of translationese: A machine learning approach. In *Computational Linguistics and Intelligent Text Processing: 11th International Conference, CICLing 2010, Iaşi, Romania, March 21-27, 2010. Proceedings 11*, pages 503–511. Springer.

Alon Jacovi, Oren Sar Shalom, and Yoav Goldberg. 2018. Understanding convolutional neural networks for text classification. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 56–65, Brussels, Belgium. Association for Computational Linguistics.

Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. *arXiv preprint arXiv:2301.08745*.

Marzena Karpinska and Mohit Iyyer. 2023. Large language models effectively leverage document-level context for literary translation, but critical errors persist. In *Proceedings of the Eighth Conference on Machine Translation*, pages 419–451, Singapore. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180. Association for Computational Linguistics.

Sai Koneru, Miriam Exel, Matthias Huck, and Jan Niehues. 2024. Contextual refinement of translations: Large language models for sentence and document-level post-editing.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Christopher A. Choquette-Choo, Katherine Lee, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2023. Madlad-400: A multilingual and document-level large audited dataset.

Viet Lai, Chien Nguyen, Nghia Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan Rossi, and Thien Nguyen. 2023. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 318–327, Singapore. Association for Computational Linguistics.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2023. Bloom: A 176b-parameter open-access multilingual language model.

Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. Visualizing and understanding neural models in NLP. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 681–691, San Diego, California. Association for Computational Linguistics.

Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. 2024. Llamax: Scaling linguistic horizons of llm by enhancing translation capabilities beyond 100 languages. *arXiv preprint arXiv:2407.05975*.

Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30.

Jiaming Luo, Colin Cherry, and George Foster. 2024. To diverge or not to diverge: A morphosyntactic perspective on machine translation vs human translation. *Transactions of the Association for Computational Linguistics*, 12:355–371.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.

Dmitry Nikolaev, Ofir Arviv, Taelin Karidi, Neta Kenneth, Veronika Mitnik, Lilja Maria Saeboe, and Omri Abend. 2020. Fine-grained analysis of cross-linguistic syntactic divergences. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1159–1176, Online. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Daria Pylypenko, Kwabena Amponsah-Kaakyire, Koel Dutta Chowdhury, Josef van Genabith, and Cristina España-Bonet. 2021. Comparing feature-engineering and feature-learning approaches for multilingual translationese classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8596–8611, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020a. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020b. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Dheeraj Rajagopal, Vidhisha Balachandran, Eduard H Hovy, and Yulia Tsvetkov. 2021. SELFEXPLAIN: A self-explaining architecture for neural text classifiers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 836–850, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Vikas Raunak, Arul Menezes, Matt Post, and Hany Hassan Awadalla. 2023. Do gpts produce less literal translations? *arXiv preprint arXiv:2305.16806*.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T.

Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Raphael Rubino, Ekaterina Lapshinova-Koltunski, and Josef Van Genabith. 2016. Information density and quality estimation features as translationese indicators for human translation classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 960–970.

Xinying Song, Alex Salcianu, Yang Song, Dave Dopson, and Denny Zhou. 2021. Fast WordPiece tokenization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2089–2103, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024a. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh

Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. 2024b. Gemma 2: Improving open language models at a practical size.

Katherine Thai, Marzena Karpinska, Kalpesh Krishna, Bill Ray, Moira Inghilleri, John Wieting, and Mohit Iyyer. 2022. Exploring document-level literary machine translation with parallel paragraphs from world literature. *arXiv preprint arXiv:2210.14250*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2022. Prompting palm for translation: Assessing strategies and performance. *arXiv preprint arXiv:2211.09102*.

Yogarshi Vyas, Xing Niu, and Marine Carpuat. 2018. Identifying semantic divergences in parallel text without annotations. *arXiv preprint arXiv:1803.11112*.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.

Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah Goodman. 2023. Interpretability at scale: Identifying causal mechanisms in alpaca. In *Advances in Neural Information Processing Systems*, volume 36, pages 78205–78226. Curran Associates, Inc.

Roy Xie, Orevaoghene Ahia, Yulia Tsvetkov, and Antonios Anastasopoulos. 2024. Extracting lexical features from dialects via interpretable dialect classifiers. *arXiv preprint arXiv:2402.17914*.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Kayo Yin and Graham Neubig. 2022. Interpreting language models with contrastive explanations. *arXiv preprint arXiv:2202.10419*.

Mo Yu, Shiyu Chang, Yang Zhang, and Tommi Jaakkola. 2019. Rethinking cooperative rationalization: Introspective extraction and complement control. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4094–4103, Hong Kong, China. Association for Computational Linguistics.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. In *International Conference on Machine Learning*, pages 41092–41110. PMLR.

Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International conference on machine learning*, pages 12697–12706. PMLR.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2765–2781, Mexico City, Mexico. Association for Computational Linguistics.

Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D'souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. 2024. Aya model: An instruction finetuned open-access multilingual language model. *arXiv preprint arXiv:2402.07827*.

## A Prompts

**LLaMAX-3.1-8B-Alpaca**

```
Below is an instruction that describes a task, paired with an input that provides
further context.
Write a response that appropriately completes the request.
### Instruction: Translate the following sentences from {source} to {target}.
Input:
{input_sentence}
### Response:
```

**TowerInstruct-7B-v0.2**

```
Translate the following sentence into {target}.
{source}: {input_sentence}
{target}:
```

**Aya-101-13B**

```
Translate to {target}: {input_sentence}
```

**LLaMA-3.1-IT-8B**

```
Translate the following sentence from {source} to {target}:
{input_sentence}
{target}:
```

**Gemma-7B**

```
Translate this sentence from {source} to {target} without any comments:
{source}:
{input_sentence}
{target}:
```

**B**



Figure 4: The frequency of the top PoS categories flagged by IG across different sentence length bins. The x-axis of each subplot represents sentence length, divided into ranges (0-10, 10-15, 15-20, etc.), and the y-axis shows PoS frequency, indicating how often each PoS occurs in sentences of different lengths.

## C Correlation Analysis

We calculate Spearman's correlation to analyze the relationship between translation quality and O/T classification accuracy, considering a significance level $\alpha = 0.05$. We find Spearman's correlation between COMET and Accuracy to be $-0.43$ with $p$-value $0.28$, and $-0.63$ with $p$-value $0.1$ between BLEU and Accuracy. Correlations are not statistically significant; therefore, given our data, there is no evidence to support the notion that poorer translations are more easily classified as translated or non-translated texts.

# How Grammatical Features Impact Machine Translation: A New Test Suite for Chinese-English MT Evaluation

**Huacheng Song**[1,3] **Yi Li**[3], **Yiwen Wu**[2], **Yu Liu**[2], **Jingxia Lin**[2], **Hongzhi Xu**[3]

[1]The Hong Kong Polytechnic University
[2]Nanyang Technological University
[3]Shanghai International Studies University
huacheng.song@connect.polyu.hk, {wu0010en, liuy0243}@e.ntu.edu.sg
jingxialin@ntu.edu.sg, {liyi, hxu}@shisu.edu.cn

## Abstract

Machine translation (MT) evaluation has evolved toward a trend of fine-grained granularity, enabling a more precise diagnosis of hidden flaws and weaknesses of MT systems from various perspectives. This paper examines how MT systems are potentially affected by certain grammatical features, offering insights into the challenges these features pose and suggesting possible directions for improvement. We develop a new test suite by extracting 7,848 sentences from a multi-domain Chinese-English parallel corpus. All the Chinese text was further annotated with 43 grammatical features using a semi-automatic method. This test suite was subsequently used to evaluate eight state-of-the-art MT systems according to six different automatic evaluation metrics. The results reveal intriguing patterns of MT performance associated with different domains and various grammatical features, highlighting the test suite's effectiveness. The test suite was made publicly available and it will serve as an important benchmark for evaluating and diagnosing Chinese-English MT systems.

## 1 Introduction

A test suite or a challenge set is a collection of customized or artificially constructed texts used for exhaustively and systematically diagnosing the hidden faults and specific barriers of models in the field of natural language processing (NLP) (King and Falkedal, 1990; Balkan, 1994). It also comes in handy in machine translation (MT) evaluation and has currently experienced an increased weight in the MT community alongside the significant improvement of average automatic translation quality especially in the era of neural machine translation (NMT) and large language model (LLM) (Burchardt et al., 2017; Kocmi et al., 2023).

By leveraging test suites, it is possible to detect the strengths and weaknesses of apparently perfect MT systems in a linguistically driven fashion

and at a fine-grained level. However, on the one hand, most previous studies have concentrated on a limited set of language phenomena (Guillou et al., 2018; Popović, 2019; Mukherjee and Shrivastava, 2023), providing only a narrow view of system capabilities. On the other hand, there is a notable scarcity of research and resources concerning non-Latin script languages, such as Chinese (Chen et al., 2023), which require some special handling of MT systems. These facts underscore the need for a large-scale test suite that covers a broad variety of grammatical features appearing in Chinese-English renderings.

Inspired by the grammatical test suite developed by the German Research Center for Artificial Intelligence (DFKI) (Manakhimova et al., 2023, etc.), we create a test suite for Chinese-English automatic translation focusing on multiple Chinese grammatical features, and based on which we conduct a detailed analysis of the state-of-the-art MT systems, including popular commercial NMT systems and advanced LLMs. The UM parallel corpus in the language pair of Chinese-English (Tian et al., 2014) originally containing segments from seven domains serves as the basis for extracting test sentences for 43 distinct Chinese grammatical features. The final test suite comprises 7,848 well-annotated Chinese sentences (at least 50 items for each grammatical feature), each paired with an English reference translation. We report the performance of eight MT systems and discuss the impact of 43 grammatical features, based on scores generated by six mainstream automatic metrics and supplemented by an analysis of manually identified error cases. We make our test suite, system outputs, evaluation scores, and corresponding codes available online for further research purposes[1].

The main contributions of our work are summarized here: 1) We present a grammatical-feature-

---

[1]https://github.com/florethsong/testsuite-zh-grammaticalfeature

based and multi-domain test suite for fine-grained Chinese-English translation evaluation. 2) We perform a linguistically driven evaluation to compare the overall performance of different NMT systems and LLMs. 3) We conduct further analysis of various influencing factors in our study from the aspects of automatic evaluation metrics and other external features of sentences to examine the impacts of different grammatical features.

This paper is structured as follows: Section 2 presents a list of studies that are related to the current work. Section 3 shows the main procedure of the construction of the test suite including data extraction and annotation. In Section 4, we describe the experiments of applying our test suite on the mainstream MT systems and give an analysis of the results. Section 5 provides additional discussions on the other interfering factors that may also interact with grammatical features to impose effects on MT. Section 6 outlines our conclusion and future work.

## 2 Related Work

In the context of probing linguistically nuanced yet critical weaknesses in MT systems to guide future enhancement, the Conference of Machine Translation (WMT) has introduced test suite tracks since 2018, aimed at receiving in-depth insights into the fine-grained performance of MT systems (Macketanz et al., 2018; Guillou et al., 2018; Rysová et al., 2019; Popović, 2019; Kocmi et al., 2020; Bawden and Sagot, 2023; Mukherjee and Shrivastava, 2023; Manakhimova et al., 2023; Chen et al., 2023, inter alia).

Standing out from many studies of MT evaluations dedicated to one or a few textual factors, e.g. Guillou et al. (2018) on pronouns, Rysová et al. (2019) on discourse-related errors, Popović (2019) on conjunctions, Kocmi et al. (2020) on gender coreference and bias, Bawden and Sagot (2023) on user-generated non-standard content, and Mukherjee and Shrivastava (2023) on multiple domains and writing styles, the series of work by DFKI (Macketanz et al., 2018; Avramidis et al., 2019, 2020; Macketanz et al., 2021, 2022; Manakhimova et al., 2023) constructed a test suite covering more comprehensive linguistic phenomena. This ever-evolving test suite comprises over 10,000 sentences now, covering up to 110 linguistic phenomena, such as false friends, named entities, negations, and so on, and across three translation directions: German ↔ English, English → Russian. By applying the combination of regular expressions and manual checks for annotating the linguistic phenomena, they test the capacity of advanced MT systems submitted to the annual WMT tasks for tackling specific translation difficulties associated with such phenomena. Their latest study (Manakhimova et al., 2023) reveals that the mainstream MT systems face great challenges with certain categories of linguistic phenomena, often in a language-dependent manner. Their detailed findings further enable MT developers to facilitate their systems by considering scenarios prone to failure and then taking corrective actions.

Beyond the translation between alphabetic languages in the Indo-European language family, the task regarding pictographic texts in the Sino-Tibetan language family, represented by Chinese, is also open to exploration. As an analytic and isolating language, Chinese has very different ways of expressing syntactic and semantic relations between constituents, resulting in potential ambiguities that largely rely on context to resolve. The issue becomes more salient in the automatic translation task. Particularly, the presence of certain grammatical features in Chinese will potentially cause different problems. However, the comprehensive exploration and the test suites with attention to various Chinese grammatical features remain largely unconsidered. The only Chinese MT test suite submitted to the WMT was constructed by Chen et al. (2023) for investigating the influence of a limited set of features of Chinese source sentences including words, length, grammar, and entropy. Besides, the studies of Cai and Xiong (2020), Tang et al. (2021), and Song and Xu (2024a,b) provided focused glimpses to some certain Chinese phenomena. They examined the abilities of NMT systems to translate discourse phenomena, negation, and multiword expressions across English and Chinese by using a self-built test suite with annotation of pronouns, discourse connectives, and ellipses, an existing corpus with negation information created by Liu et al. (2018), and an extended dataset of WMT test set, respectively.

Building on the light of DFKI test suites (Manakhimova et al., 2023, etc.) and addressing the lack of Chinese-specific test suites, this study is dedicated to providing an inclusive test suite covering 43 Chinese grammatical features and offering a full evaluation of mainstream MT systems. Addition-

ally, depending on the domain-balanced nature of our basic data, i.e. UM corpus (Tian et al., 2014), our test suite is suitable for comparisons across seven textual domains.

## 3 Construction of a Test Suite with Chinese Grammatical Features

This section details the processes in test suite construction. We first introduce the theoretical framework of the Chinese grammatical features, then describe the procedures of data selection and annotation, and finally present the result and statistics of the data.

### 3.1 The Framework of Chinese Grammatical Features

For the categorization framework for Chinese grammatical features, we adopt the one in our previous work (Xu and Lin, 2023), which is developed in accordance with a reference grammar of Chinese (Huang and Shi, 2016). The framework systematically addresses 157 typical linguistic phenomena, i.e. grammatical features, in Chinese, organized across various linguistic aspects, including words, structure, semantics, and pragmatics. Word-level structures are concerned with how a word is formed by morphemes. For instance, reduplication is a typical phenomenon to create new words in Chinese. For example, the adjective 高兴 *gao xing* 'happy' can be reduplicated to form another adjective word 高高兴兴 *gao gao xing xing* 'very happy'. The structure category mainly refers to the syntactic structure of sentences, phrases, and special constructions. This framework identifies three semantic subcategories: semantic roles, aspect, and negation. The pragmatic category includes sentence types, information packaging constructions, attitudinal particles/adverbs, deixis, and anaphora.

Whether a certain grammatical feature that is present in the source language might cause problems in automatic translation is largely dependent on the equivalence of the counterpart phenomenon in the target language. Take reflexives as an example. Both the two languages use reflexive pronouns to denote the antecedent nominal phrase. However, there are also some fine distinctions in their usages, leading to obstacles for cross-lingual translation. As shown in Example (1), while in Chinese the pronoun 你 *ni* 'you' can be optional, the English translation must combine the pronoun 'you' in order to obtain the correct reflexive 'yourself'.

(1)　　你要照顾好(你)自己。

ni　yao　zhaogu　hao　(ni)　ziji
you should take_care good (you) self
'You should take care of yourself.'

Many grammatical features are Chinese-specific, such as classifiers, as shown in example (2), BA constructions as shown in (3), headless NP as shown in (4), and so on. Depending on their grammatical differences to varying degrees, different Chinese grammatical features might impose different effects on MT systems. It is thus necessary to create a test suite that covers various grammatical features with each one associated with a set of examples, which can be used to analyze the effects of different grammatical features on MT systems based on statistical methods.

(2)　　一顿晚餐

yi　dun　wancan
one CLF dinner
'a dinner'

(3)　　我把这些书都看完了。

wo ba zhexie shu　dou kan　wan　le
I　BA these　book all　read finish PRF
'I have read all these books.'

(4)　　羡慕的是缺乏的。

xianmu de　(pro) shi quefa de　(pro)
admire DE　　be lack　DE
'What is admired is the lacked.'

### 3.2 Data Preparation

We extract Chinese source sentences and their corresponding English reference translations from the UM corpus (Tian et al., 2014), a high-quality and large-scale parallel corpus embracing eight distinct domains: Education (abbreviated as 'Edu', with 4.5 million bilingual sentence pairs), Laws ('Laws', 2.2M), News ('News', 4.5M), Science ('Sci', 2.7M), Spoken ('Spk', 2.2M), Subtitles ('Sbt', 3M), Thesis ('Ths', 3M), and Microblog ('Mbg', 5K). We exclude the Microblog section due to its small number of sentence pairs, which is far fewer than the other domains, making it difficult to ensure a rough balance across different domains. We select sentences with 10 to 60 Chinese characters to minimize the impact of source sentences with extreme lengths (excessively long or short) on translation quality as well as to avoid the existence of too many different grammatical features in a single sentence that may mix the effects of them on translations.

| Grammatical Feature | Abbreviation | Precision | Agreement | Sum | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Verb Phrase | VP | 0.83 | 0.83 | 220 | 24 | 33 | 32 | 41 | 41 | 27 | 22 |
| Noun Phrase | NP | 0.91 | 0.98 | 1440 | 152 | 344 | 190 | 184 | 184 | 145 | 241 |
| Adjective Phrases | AdjP | 0.36 | 0.44 | 230 | 18 | 46 | 42 | 24 | 19 | 16 | 65 |
| Adverb Phrases | AdvP | 0.83 | 0.95 | 951 | 109 | 180 | 141 | 152 | 100 | 95 | 174 |
| Pre-verbal Preposition Phrase | PreVPP | 0.91 | 0.94 | 163 | 14 | 29 | 26 | 23 | 27 | 18 | 26 |
| Post-verbal Preposition Phrase | PstVPP | 0.28 | 0.77 | 146 | 23 | 20 | 26 | 22 | 29 | 24 | 2 |
| Participant Preposition Phrase | PtcpPP | 0.89 | 0.83 | 301 | 26 | 82 | 44 | 23 | 38 | 27 | 61 |
| Topic Preposition Phrase | TopPP | 0.98 | 0.98 | 213 | 25 | 30 | 31 | 35 | 31 | 22 | 39 |
| Reference Preposition Phrase | RefPP | 0.96 | 0.99 | 347 | 42 | 69 | 52 | 49 | 46 | 34 | 55 |
| Condition Preposition Phrase | CondPP | 0.51 | 0.89 | 105 | 18 | 33 | 18 | 9 | 6 | 3 | 18 |
| Locative Preposition Phrase | LocPP | 0.5 | 0.93 | 96 | 15 | 18 | 18 | 12 | 13 | 12 | 8 |
| Sentence-Initial Preposition Phrase | SentIPP | 0.33 | 0.82 | 64 | 7 | 5 | 11 | 9 | 11 | 3 | 18 |
| Space Preposition Phrase | SpcPP | 0.76 | 0.86 | 155 | 18 | 32 | 20 | 17 | 29 | 24 | 15 |
| Source Preposition Phrase | SrcPP | 0.96 | 0.96 | 191 | 27 | 29 | 26 | 28 | 26 | 26 | 29 |
| Path Preposition Phrase | PathPP | 0.8 | 0.93 | 132 | 12 | 16 | 23 | 19 | 28 | 16 | 18 |
| Goal Preposition Phrase | GoalPP | 0.65 | 0.68 | 127 | 21 | 16 | 22 | 19 | 27 | 19 | 3 |
| Direction Preposition Phrase | DirPP | 0.47 | 0.48 | 95 | 20 | 6 | 18 | 12 | 19 | 17 | 3 |
| Space Extension Preposition Phrase | SpanPP | 0.93 | 0.18 | 169 | 25 | 28 | 25 | 22 | 25 | 16 | 28 |
| Standard Classifier | StdCLF | 0.98 | 0.99 | 195 | 29 | 39 | 28 | 24 | 28 | 22 | 25 |
| Individual Classifier | IndCLF | 0.94 | 0.97 | 284 | 28 | 58 | 36 | 40 | 41 | 35 | 46 |
| Event Classifier | EvCLF | 0.97 | 0.97 | 184 | 25 | 23 | 24 | 24 | 31 | 27 | 30 |
| Kind Classifier | KindCLF | 0.98 | 0.99 | 185 | 25 | 23 | 29 | 29 | 27 | 22 | 30 |
| Approximation Classifier | ApprCLF | 0.35 | 0.81 | 68 | 10 | 13 | 14 | 11 | 10 | 6 | 4 |
| Temporal Sequence Complex Sentence | TmpSCpl | 0.99 | 0.98 | 176 | 21 | 28 | 29 | 26 | 26 | 19 | 27 |
| Concessive Complex Sentence | ConcCpl | 0.99 | 0.99 | 156 | 20 | 10 | 29 | 27 | 29 | 14 | 27 |
| Causative Complex Sentence | CausCpl | 0.46 | 0.82 | 82 | 13 | 8 | 17 | 20 | 13 | 2 | 9 |
| Negation BU | BUNeg | 0.96 | 0.91 | 222 | 31 | 37 | 30 | 32 | 35 | 23 | 34 |
| Negation MEI/MEIYOU | MEINeg | 0.98 | 0.97 | 225 | 32 | 33 | 33 | 33 | 36 | 33 | 25 |
| Negation in Imperative Sentences | ImpNeg | 0.36 | 0.82 | 83 | 8 | 37 | 11 | 12 | 5 | 10 | 0 |
| Sublexical Negation | LexNeg | 0.97 | 0.97 | 182 | 23 | 30 | 26 | 28 | 24 | 20 | 31 |
| Negative Polarity Items | NPI | 0.98 | 0.92 | 166 | 19 | 27 | 28 | 29 | 24 | 14 | 25 |
| Deixis | Deixis | 0.95 | 0.57 | 272 | 37 | 26 | 38 | 46 | 48 | 45 | 32 |
| Reflexive | Refl | 0.96 | 0.76 | 195 | 25 | 26 | 29 | 30 | 31 | 25 | 29 |
| Reciprocal | Recp | 1 | 1 | 174 | 23 | 27 | 26 | 26 | 23 | 20 | 29 |
| Perfective GUO | GUOPrf | 0.84 | 0.91 | 154 | 21 | 24 | 27 | 20 | 28 | 26 | 8 |
| Progressive ZAI | ZAIProg | 0.99 | 0.98 | 184 | 26 | 23 | 28 | 27 | 30 | 21 | 29 |
| Passive Construction | Pass | 0.66 | 0.84 | 131 | 15 | 30 | 23 | 18 | 19 | 17 | 9 |
| Relative Construction | Rel | 0.78 | 0.85 | 305 | 16 | 119 | 32 | 25 | 31 | 19 | 63 |
| Comparative Construction | Cmpr | 0.95 | 0.99 | 191 | 25 | 28 | 28 | 27 | 24 | 23 | 36 |
| BA Construction | BA | 0.99 | 0.99 | 199 | 21 | 38 | 30 | 28 | 28 | 23 | 31 |
| Copular SHI | SHICop | 0.98 | 0.94 | 230 | 34 | 36 | 38 | 31 | 33 | 23 | 35 |
| Verbal LE | VerbLE | 0.86 | 0.96 | 194 | 31 | 10 | 33 | 31 | 28 | 25 | 36 |
| Quantifier Only ZHI | ZHIQtf | 0.97 | 0.98 | 220 | 22 | 24 | 28 | 25 | 28 | 26 | 28 |
| **Overall/Total** | | **0.81** | **0.87** | **9763** | **1176** | **1793** | **1459** | **1369** | **1379** | **1084** | **1503** |
| **Sentence Number** | | | | **7848** | **1035** | **1109** | **1207** | **1187** | **1175** | **927** | **1208** |

Table 1: The detailed information of our test suite including the definition of grammatical features and their corresponding numbers of instances in each domain. 'Precision' refers to the precision of our self-created grammatical feature identifier in accordance with the results after manual checking. 'Agreement' shows the inner consistency between the judgments given by the two checkers. The acronyms including 'BU', 'MEI/MEIYOU', 'GUO', 'ZAI', 'SHI', 'LE', 'BA', and 'ZHI', are the specific markers indicating particular grammatical features.

### 3.3 Grammatical Feature Annotation

In the first step, we use a regular-expression-based tool we previously built in Xu and Lin (2023) to identify the Chinese grammatical features in each source sentence automatically. After all the sentences are annotated with a set of grammatical features, we remove the grammatical features that appear fewer than 30 times in all sentences of each domain to ensure a fairly balanced data distribution in statistics, by which our focus is narrowed to 43 target grammatical features out of 157 in the original framework for our test suite. Then, for each grammatical feature, we randomly select about 210 candidate sentences (30 for each of the seven domains) according to the principle of prioritizing those carrying the fewest labels aiming to reduce the mixed effects of multiple features in one sentence. Since some sentences can finally possess multiple features, certain feature groups may include more than 210 sentences.

In the second step, the automatically generated labels of grammatical features are double-checked by two native speakers well-trained in Chinese linguistics. The screening process is primarily focused on identifying false positive grammatical features assigned to sentences. In cases where annotators disagree, they are required to discuss and reach a final decision together. This filtering process resulted in a validated test suite of a total of 7,848 sentence pairs, including 1,127 pairs with no specific grammatical features. Detailed information on the data is shown in Table 1. We see that the average precision of the automatic annotation tool is about 81% and the agreement of the two annotators in identifying false positives is 87%.

## 4 Evaluation of Chinese-English MT Systems with the Test Suite

In this section, we use our test suite to evaluate eight popular NMT systems and LLMs with six mainstream automatic metrics. We will briefly outline these systems and metrics, and then describe the results of the experiments we conduct to compare the performance of different MT systems on our whole test suite as well as the subgroups divided by domains and grammatical features.

### 4.1 Evaluated Translation Systems

Aiming at gaining a broad view of the capabilities of leading MT systems and representative LLMs to tackle diverse Chinese grammatical features, we refer to several widely recognized leaderboards, e.g., WMT (Kocmi et al., 2022), SuperCLUE (Xu et al., 2023), SuperBench[2], and Intento[3]. Eventually, four commercial NMT engines, Baidu, Niu, Google (basic v2 edition), DeepL and four advanced LLMs including Ernie (-4 turbo), Qwen (-turbo), GPT (-4o), and Claude (-3 opus) are selected for evaluation. We apply default settings to the NMTs and conduct a zero-shot translation test for the LLMs, setting the temperatures to 0.01.

### 4.2 Automatic Evaluation Metrics

We use six automatic metrics, including two string-overlap-based metrics: BLEU (Papineni et al., 2002) and CHRF (Popović, 2015)[4], and four neural-

network-based ones: two reference-based metrics: COMET (Rei et al., 2022) and XCMOET (Guerreiro et al., 2023), and two reference-free ones: COMETKIWI-QE (Rei et al., 2023) and XCOMET-QE (Guerreiro et al., 2023)[5].

Based on our observation, different metrics may produce different results in analyzing the effects of various factors impacting MT systems. In the following discussions, we will mainly use the average score of XCOMET and XCOMET-QE (henceforth, X-AVERAGE), which are proven the most accurate metrics conforming to human evaluations by the WMT23 metric shared task (Freitag et al., 2023). We also provide extended discussions about the selection of metrics in Section 5.3. For reference, readers can find the results in all six metrics and their overall average (denoted as AVERAGE) in the Appendices.

### 4.3 Experimental Results

#### 4.3.1 Comparison of Systems

**Comparison of Overall Performance**  Table 2 shows the overall performance of the eight systems in six different metrics. On average, Google performs the best and Qwen the worst. While most of the metrics give similar ratings, XCOMET and XCOMET-QE slightly favor GPT's performance more than the others. We can also see that NMTs achieve marginally better performance than LLMs across all evaluation results. This is partially attributed to the extremely low scores received by the LLM: Qwen.

**Comparison on Domains**  Table 3 shows the performance of all the MT systems on different domains in X-AVERAGE scores. Generally, all the systems show similar trends across different domains with the highest performance on Spoken and Subtitles and the lowest performance on Thesis and Laws. While it is unquestionable that domain affects automatic translation, the results may also be partially influenced by the distribution of sentence length within each domain. We will give more discussion about the effects of sentence length in Section 5.2. Detailed information about the performance of systems on different domains in all the other metrics can be found in Appendix A. Based on the comparisons between different models, it

---

[2]The online report of SuperBench by Tsinghua University: https://fm.ai.tsinghua.edu.cn/superbench/#/leaderboard

[3]The online report *The State of Machine Translation 2024* by Intento: https://inten.to/machine-translation-report-2024/

[4]BLEU and CHRF are computed by SacreBLEU implementations: https://github.com/mjpost/sacrebleu, with 'True'

in signatures of effective order, lowercase, and whitespace and taking 'exp' as smooth method.

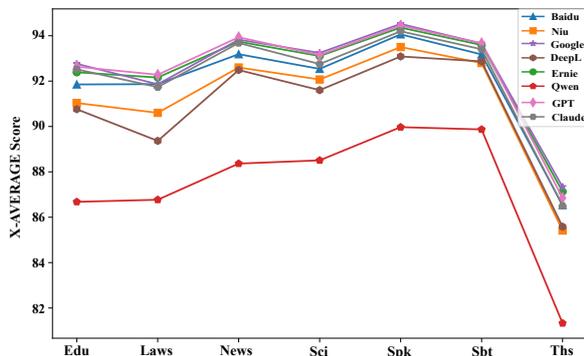[5]The series of *COMET are computed by Unbabel implementations: https://github.com/Unbabel/COMET

| | BLEU | CHRF | COMET | XCOMET | COMETKIWI-QE | XCOMET-QE | AVERAGE | X-AVERAGE |
|---|---|---|---|---|---|---|---|---|
| Baidu | 21.6 | 56.0 | 80.8 | 72.7 | 90.3 | 93.4 | 69.1 | 91.8 |
| Niu | 25.1 | 58.5 | 81.0 | 71.6 | 89.9 | 92.3 | 69.7 | 91.1 |
| Google | **27.0** | **59.8** | **82.1** | 72.8 | **91.6** | 93.2 | **71.1** | **92.4** |
| DeepL | 24.3 | 57.6 | 80.4 | 71.7 | 89.4 | 92.1 | 69.2 | 90.8 |
| Ernie | 24.0 | 58.1 | 81.5 | 73.1 | 91.3 | 93.3 | 70.2 | 92.3 |
| Qwen | <u>16.5</u> | <u>47.0</u> | <u>76.3</u> | <u>65.0</u> | <u>85.4</u> | <u>89.1</u> | <u>63.2</u> | <u>87.3</u> |
| GPT | 22.5 | 57.0 | 81.3 | **73.4** | 91.1 | **93.6** | 69.8 | **92.4** |
| Claude | 23.3 | 57.7 | 81.3 | 73.0 | 91.0 | 93.1 | 69.9 | 92.1 |
| NMT-AVG | 24.5 | 58.0 | 81.1 | 72.2 | 90.3 | 92.8 | 69.8 | 91.5 |
| LLM-AVG | 21.6 | 55.0 | 80.1 | 71.1 | 89.7 | 92.3 | 68.3 | 91.0 |

Table 2: The overall performance of MT systems in different metrics. The highest and the lowest scores among all systems evaluated are highlighted with bold and underlined numbers respectively. 'AVERAGE' is the mean of scores by the six metrics, while 'X-AVERAGE' is the mean of scores by XCOMET and XCOMET-QE.



| | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| Baidu | 91.9 | 91.9 | 93.2 | 92.5 | 94.1 | 93.2 | 86.5 |
| Niu | 91.0 | 90.6 | 92.6 | 92.1 | 93.5 | 92.8 | 85.4 |
| Google | **92.8** | 91.9 | 93.8 | **93.2** | 94.5 | 93.7 | **87.3** |
| DeepL | 90.8 | 89.4 | 92.5 | 91.6 | 93.1 | 92.9 | 85.6 |
| Ernie | 92.4 | 92.2 | 93.7 | 93.1 | 94.4 | 93.6 | 87.1 |
| Qwen | <u>86.7</u> | <u>86.8</u> | <u>88.4</u> | <u>88.5</u> | <u>90.0</u> | <u>89.9</u> | <u>81.3</u> |
| GPT | 92.6 | **92.3** | **93.9** | 93.1 | 94.4 | **93.7** | 86.9 |
| Claude | 92.5 | 91.7 | 93.7 | 92.8 | 94.2 | 93.4 | 86.5 |
| NMT-AVG | 91.6 | 90.9 | 93.0 | 92.3 | 93.8 | 93.2 | 86.2 |
| LLM-AVG | 91.1 | 90.8 | 92.4 | 91.9 | 93.2 | 92.7 | 85.4 |

Table 3: Performance in X-AVERAGE scores of each system on different domains. The highest and the lowest scores among all systems evaluated are highlighted with bold and underlined numbers respectively.

can be pointed out that Google and GPT share the top performance on each domain with a minor difference but Qwen is ranked last across all domains. Again, NMTs show a generally higher performance than LLMs across all domains. However, this is not the case when delving deeper into the specific data excluding Qwen and Google.

**Comparison on Grammatical Features** Figure 1 shows the performance of all systems in X-AVERAGE scores in each grammatical feature

group. The groups are arranged in descending order based on the average scores of all systems. In general, we see that all the systems show similar trends across different groups. Some grammatical features impose strong challenges on the MT systems such as PathPP, ApprCLF, KindCLF, LexNeg, etc., while some other grammatical features are easier for MT systems to address, like ZAIProg, MEINeg, NPI, etc. We also see that Ernie, Google, Claude, and GPT give high performance on sentence groups of all grammatical features while Qwen performs obviously the worst among all the systems. The detailed statistics can be found in Table 15. The performance in other metrics can be found in Appendix B.

It is therefore indicated that the presence of certain grammatical features will potentially affect the performance of MT systems. We assess the impact of each grammatical feature by conducting a t-test between the MT performance on the sentence group containing the target grammatical feature and that on the remaining sentences. The result is shown in Figure 2. There are ten grammatical features imposing significant negative effects on certain MT systems: PathPP, NP, Rel, KindCLF, PtcpPP, LexNeg, PreVPP, TmpSCpl, LocPP, and Cmpr; and there are nine grammatical features having significant positive effects instead: ZAIProg, MEINeg, NPI, Recp, ZHIQtf, AdvP, Refl, SHICop, and VP. However, it is worth noting that the low scores on certain grammatical feature groups are not necessarily occasioned by the translation errors that are directly linked to the units marking the grammatical features. There are also many other implicit factors indirectly associated with the grammatical features being worthy of exploration, like the semantic or syntactic complexity.
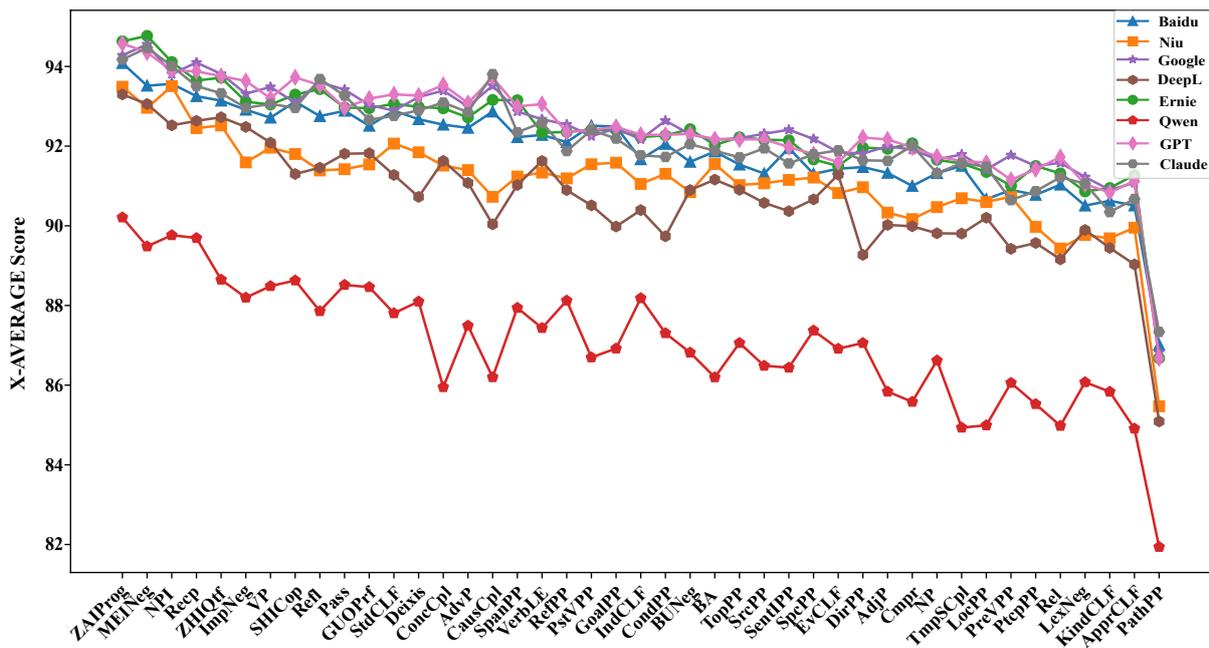
Figure 1: Performance in X-AVERAGE scores of each system on different grammatical features.
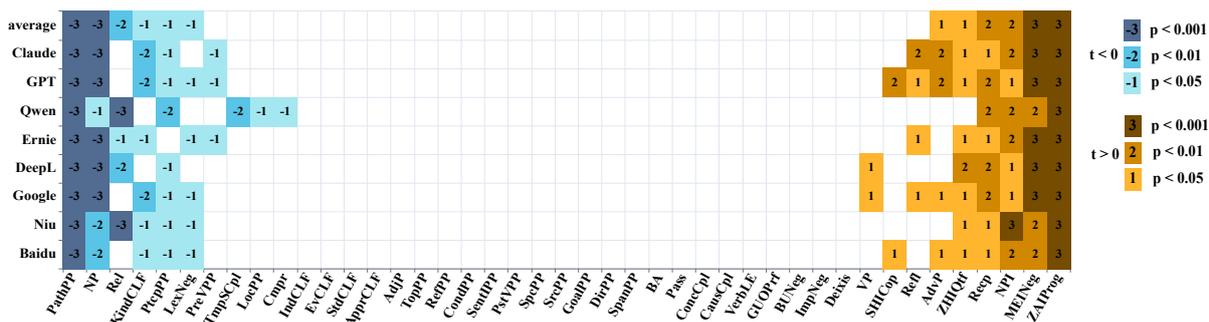


Figure 2: The significance levels of the effects of different grammatical features on the performance of each system according to t-tests in X-AVERAGE scores between the paired sentence groups with and without each grammatical feature. The blue squares on the left mean that the sentences containing a certain grammatical feature tend to get significantly lower scores than the sentences without them, and the yellow squares on the right vice versa. A significant difference with $p$ less than 0.05, 0.01, and 0.001 are marked as 1, 2, and 3 respectively.

Since the most advanced automatic metrics are based on semantic vectors, it requires some consideration of why certain translations receive low scores, particularly in a specific sentence group. The following are some case studies to provide hints of the possible cause of the low translation scores. We should also keep in mind that this study is not aimed to provide a full understanding of why grammatical features can affect MT systems, but instead provide a different aspect and dimension for diagnosing MT systems in fine-grained levels.

### 4.3.2 Case Studies

In the process of translation, the accurate conveyance of meaning from the source language to the target language is important. However, errors often arise due to misinterpretation of some grammatical features, leading to mistranslation that may alter the meaning of the original text. This case study manually examines several specific instances of such errors in the translation.

One of the primary issues in automatic translation arises from the misinterpretation of path preposition phrases (PathPPs), as shown in Figure 2. Table 4 shows a typical example, where the preposition 沿 *yan* 'along' was misinterpreted, resulting in mistranslations by several MT systems. The original sentence uses the PathPP 沿 渤海 公路 *yan Bohai gong lu* 'along Bohai Highway' to describe how 唐海 *Tanghai* (a town) extends from east to west.

| Source Text | Reference | Ernie | Qwen | Claude |
|---|---|---|---|---|
| 唐海地处唐山"金三角"中央地带， | Tanghai is located in the central area of Tangshan's "Golden Triangle", | Tanghai is located in the central area of Tangshan's "Golden Triangle". | Tanghai is in the heart of Tangshan's Golden Triangle, | Tanghai is located in the central area of the "Golden Triangle" of Tangshan. |
| 沿渤海公路贯穿东西， | running east and west along the Bohai Highway. | The Bohai Highway runs through the east and west. | linked by the Bohai coastal highway. | The Bohai Highway runs through the county from east to west. |

Table 4: Example of mistranslation caused by grammatical feature: path preposition phrase (PathPP).

| Source Text | Reference | Niu | Qwen | DeepL |
|---|---|---|---|---|
| 如果什么东西是充足的它就是不令人羡慕的，羡慕的是缺乏的。 | When something is in plenty it is not admired, but admired in case of scarcity. | If something is sufficient, it is not enviable, and enviable is lacking. | What is abundant is unenvied; it is the absence that counts. | If something is sufficient it is not enviable, envy is lacking. |

Table 5: Example of mistranslation caused by the grammatical feature: relative construction (Rel).

However, in translations by Ernie and Claude, the sentence was incorrectly rendered as 'The Bohai Highway runs through the east and west'. Qwen's translation, nevertheless, used 'linked by' to describe this relationship, which greatly shifted the meaning of the source sentence.

Additionally, the translation quality was significantly affected by the relative construction (Rel) (see in Figure 2). Table 5 shows a typical example of relative construction. In the source sentence, 羡慕的 *xian mu de* '(things) that are enviable' is a headless clause where the elliptic head noun refers to 东西 *dong xi* 'things' mentioned earlier. Niu's translation misinterpreted 羡慕的 *xian mu de* as an adjectival phrase rather than a subject of the relative construction, and thus misunderstood the meaning of the original sentence. Qwen just omitted the real subject -羡慕的(东西)- of the sub-clause, leading to mistranslating the adjective 缺乏的 *que fa de* 'scarce' as the subject. DeepL's translation completely ignored the relative construction marker 的 '*de*' and treated 羡慕 *xian mu* 'envy' as the subject.

Another grammatical point that has a significant negative impact is noun phrases (NPs). NP has a relatively large number of sentences on Laws and Thesis (see in Table 1), which have generally low averaged performance (see in Table 3) possibly due to their high semantic complexity regarding professionalism. This partially explains the negative effects of NP. Particularly, the specialized terminology within the Thesis category can notably contribute to the translation challenges.

## 5 Additional Discussion

In this section, we discuss several potential interfering factors that may also affect the quality of automatic translation by interacting with grammatical features, including sentence length, domain, and the effects of different automatic metrics.

### 5.1 Analysis of Sentence Length

It has long been an observed consensus that longer sentences are generally more difficult to MT systems and thus result in lower qualities and scores (Cho et al., 2014; Koehn and Knowles, 2017). This can also be verified by the significant inverse relationship between the lengths of source sentences and their translation scores given by human experts as shown in Figure 3, generated on WMT23 data (Freitag et al., 2023).
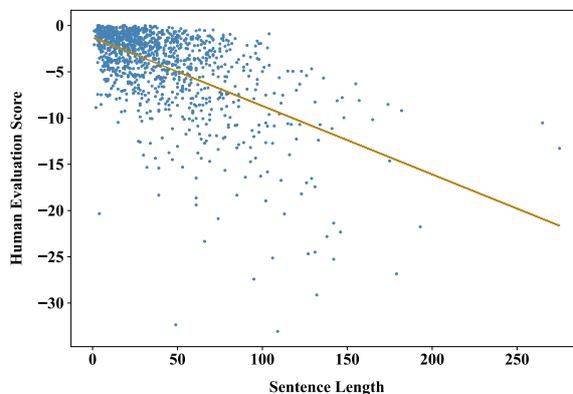


Figure 3: The correlation between sentence lengths (x-axis) and human average translation scores (y-axis) on WMT23 Chinese-English dataset for the metric shared task.

To assess whether the significant effects of grammatical features as observed are due to differences in sentence lengths among the groups, we calculated the average sentence length for each sentence group with a specific grammatical feature and examined the relationship between average sentence lengths and the corresponding X-AVERAGE scores. From the result, as shown in Figure 4, we can see that although certain sentence groups of different grammatical features have different sentence lengths ranging from 19.39 to 32.18 characters, they do not show significant correlation with the average X-AVERAGE scores of the groups, indicating that effects by grammatical features are not due to the bias of sentence length distribution.



Figure 4: The correlation between average sentence lengths and X-AVERAGE scores of sentence groups containing different grammatical features with Pearson'$r$ = -0.039 and $p$ = 0.802.

## 5.2 Analysis of Domains

As mentioned in Section 4.3, domains have a significant influence on the performance of MT systems due to variations in vocabulary and writing registers. Therefore, it is important to consider whether domains have contributed to the observed significant impact of certain grammatical features on MT systems. Reviewing the data in Table 1, we make all grammatical feature groups maintain a roughly balanced distribution in terms of sentence number across the seven domains except Rel, which extensively exists in the domain of Laws. This balance allows us to focus more on the effects of grammatical features rather than domains when calculating statistics between feature-accordingly grouped sentences.

Interestingly, the similar trends of MT systems' performance across different domains may also correlate with other factors e.g. sentence length. Thus,

a further question is whether the effects apparently imposed by domains on the scores of MT systems are partially due to the imbalanced distribution of sentence length across domains. Table 6 displays the average sentence lengths of different domains along with their standard deviations. We see that Spoken and Subtitles have the shortest sentences while Laws and Thesis have the longest ones, with a gap of about 20 characters between them. This may partially explain why MT systems achieve the best performance when rendering materials in the former two domains while the worst is in the latter two domains as shown in Table 3.

| Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|
| 23±7 | 30±11 | 25±10 | 21±10 | 18±6 | 15±4 | 29±12 |

Table 6: Average sentence length of each domain with the standard deviation.

## 5.3 Analysis of Evaluation Metrics

Following the hypothesis of regarding human evaluation as the gold standard, the metrics that generate judgments on translation quality more similar to humans are superior (Freitag et al., 2023).

In Figure 3, we see that there exists a significant negative correlation relationship between sentence lengths and human evaluation scores. To know if different metrics rate MT qualities similarly regarding sentence length, we examine the correlations between sentence lengths and scores generated by six metrics to meta-evaluate their effectiveness. We provide scatter plots of sentence lengths and scores based on both WMT23 data (Freitag et al., 2023) and our data, and the results are presented in Figure 5. On both datasets, XCOMET and XCOMET-QE exhibit patterns similar to human evaluations and are therefore considered to provide more reliable scores, particularly regarding the negative effects of sentence length. However, BLEU, CHRF, COMET, and COMET-QE yield judgments on translation quality that are inconsistent with human evaluations. According to Table 7, CHRF, COMET, and COMET-QE even exhibit significant positive correlations, indicating their bias towards longer sentences. This finding is consistent with the leaderboard of metrics concluded by WMT23 (Freitag et al., 2023), which ranks XCOMET and XCOMET-QE as the top performers.

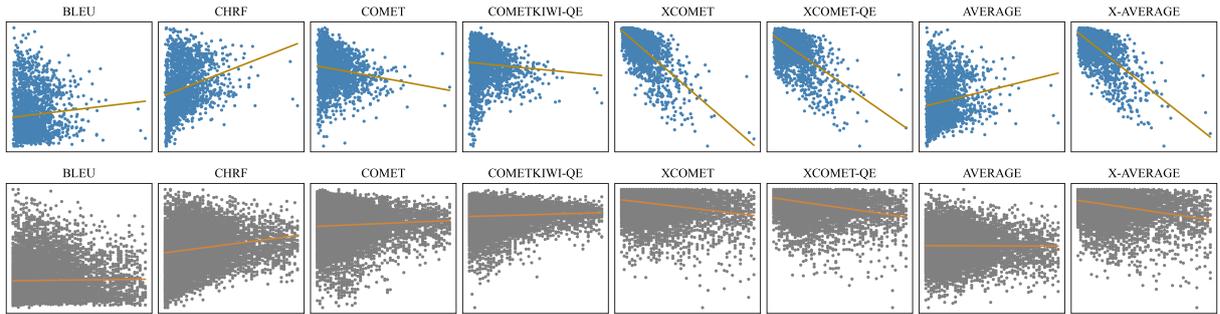Besides, Table 7 shows that the average system

Figure 5: The correlation between sentence lengths (x-axis) and average translation scores (y-axis) in different measures. XCOMET and XCOMET-QE are more consistent with the human evaluation trend in Figure 3. The first row is on the WMT23 Chinese-English dataset for the metric shared task and the second row is on our test suite. Metrics show similar patterns on the two different datasets except for COMET and COMETKIWI-QE.

| Metric | Pearson's $r$ | $P$-value |
|---|---|---|
| BLEU | 0.016 | 0.149 |
| CHRF | 0.185 | *** 0.000 |
| COMET | 0.067 | *** 0.000 |
| COMETKIWI-QE | 0.064 | *** 0.000 |
| XCOMET | -0.242 | *** 0.000 |
| XCOMET-QE | -0.302 | *** 0.000 |
| AVERAGE | -0.006 | 0.581 |
| X-AVERAGE | -0.286 | *** 0.000 |

Table 7: The correlation between sentence lengths and system average scores in different metrics on our test suite.

scores of all six automatic metrics (AVERAGE) do not significantly correlate with sentence length by offsetting the effects of different metrics. On the contrary, the average scores of XCOMET and XCOMET-QE (X-AVERAGE) remain the high reliability by showing a significant negative correlation between scores and sentence lengths. Therefore, XCOMET, XCOMET-QE, and X-AVERAGE are more recommended for practical evaluation.

## 6 Conclusion and Future Work

In this paper, we investigate the impact of various grammatical features (linguistic phenomena) on eight state-of-the-art NMT systems and LLMs with a test suite we newly constructed. Although LLMs have achieved promising performance on many NLP tasks, NMT systems especially Google Translate have outperformed most of the LLMs in the Chinese-English automatic translation task. It is observed that certain grammatical features pose a great challenge to NMT systems and LLMs including the ones developed by Chinese companies such as Baidu, Ernie, Niu, and Qwen. We also discuss other possible factors that may also impact

MT systems including sentence length, domain, and the evaluation metrics. We find that the Thesis category is particularly more difficult due to its comparatively longer sentence and the existence of a large number of terminologies. In addition, we confirm that longer sentences are generally more difficult for MT systems. However, our analysis of the correlation between the sentence length and different metrics reveals that BLEU and CHRF tend to rate shorter sentences with lower scores, which is contradictory to human evaluation. This also confirms that XCOMET and XCOMET-QE are the most reliable metrics according to the results of the WMT23 metrics shared task.

Currently, our test suite does not cover all the 157 grammatical features of Chinese due to the rareness of some particular grammatical features. In the future, we plan to extend our test suite to cover all the grammatical features by resorting to other resources.

## Limitations

One limitation of our study is the absence of human evaluation scores. Our analysis relies heavily on automatic metrics, specifically the average score of XCOMET and XCOMET-QE. The former relies on reference translations and the latter does not. According to the WMT23 metrics shared task results (Freitag et al., 2023), both metrics show a very high correlation with human scores. This demonstrates the validity and reliability of data in our study to some extent. While human evaluation is the most reliable, it is also expensive and impractical for assessing every MT system. In contrast, the test suite, combined with automatic evaluation metrics, offers a convenient and efficient tool for evaluating any MT systems, providing immediate

diagnostic reports.

Another limitation of our study is that it does not cover all the grammatical features of Chinese due to the scarcity of certain grammatical features. We plan to address this issue by exploring other data sources to cover all other grammatical features in the future.

## Acknowledgements

## References

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, Aljoscha Burchardt, and Sebastian Möller. 2020. Fine-grained linguistic evaluation for state-of-the-art machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 346–356, Online. Association for Computational Linguistics.

Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, and Hans Uszkoreit. 2019. Linguistic evaluation of German-English machine translation using a test suite. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 445–454, Florence, Italy. Association for Computational Linguistics.

Lorna Balkan. 1994. Test suites: some issues in their use and design. In *Proceedings of the Second International Conference on Machine Translation: Ten years on*, Cranfield University, UK.

Rachel Bawden and Benoît Sagot. 2023. RoCS-MT: Robustness challenge set for machine translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 198–216, Singapore. Association for Computational Linguistics.

Aljoscha Burchardt, Vivien Macketanz, Jon Dehdari, Georg Heigold, Peter Jan-Thorsten, and Philip Williams. 2017. A linguistic evaluation of rule-based, phrase-based, and neural mt engines. *The Prague bulletin of mathematical linguistics*, 108(1):159.

Xinyi Cai and Deyi Xiong. 2020. A test suite for evaluating discourse phenomena in document-level neural machine translation. In *Proceedings of the Second International Workshop of Discourse Processing*, pages 13–17, Suzhou, China. Association for Computational Linguistics.

Xiaoyu Chen, Daimeng Wei, Zhanglin Wu, Ting Zhu, Hengchao Shang, Zongyao Li, Jiaxin Guo, Ning Xie, Lizhi Lei, Hao Yang, and Yanfei Jiang. 2023. Multifaceted challenge set for evaluating machine translation performance. In *Proceedings of the Eighth Conference on Machine Translation*, pages 217–223, Singapore. Association for Computational Linguistics.

Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 103–111, Doha, Qatar. Association for Computational Linguistics.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Nuno Guerreiro, Ricardo Rei, Daan Van, Pierre Colombo, Luisa Coheur, and André Martins. 2023. xCOMET: Transparent Machine Translation Evaluation through Fine-grained Error Detection. *Transactions of the Association for Computational Linguistics*.

Liane Guillou, Christian Hardmeier, Ekaterina Lapshinova-Koltunski, and Sharid Loáiciga. 2018. A pronoun test suite evaluation of the English–German MT systems at WMT 2018. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 570–577, Belgium, Brussels. Association for Computational Linguistics.

Chu-Ren Huang and Dingxu Shi, editors. 2016. *A Reference Grammar of Chinese*. Cambridge University Press, Cambridge.

Margaret King and Kirsten Falkedal. 1990. Using test suites in evaluation of machine translation systems. In *COLING 1990 Volume 2: Papers presented to the 13th International Conference on Computational Linguistics*.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Tom Kocmi, Tomasz Limisiewicz, and Gabriel Stanovsky. 2020. Gender coreference and bias evaluation at WMT 2020. In *Proceedings of the Fifth Conference on Machine Translation*, pages 357–364, Online. Association for Computational Linguistics.

Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.

Qianchu Liu, Federico Fancellu, and Bonnie Webber. 2018. NegPar: A parallel corpus annotated for negation. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, and Hans Uszkoreit. 2018. Fine-grained evaluation of German-English machine translation based on a test suite. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 578–587, Belgium, Brussels. Association for Computational Linguistics.

Vivien Macketanz, Eleftherios Avramidis, Shushen Manakhimova, and Sebastian Möller. 2021. Linguistic evaluation for the 2021 state-of-the-art machine translation systems for German to English and English to German. In *Proceedings of the Sixth Conference on Machine Translation*, pages 1059–1073, Online. Association for Computational Linguistics.

Vivien Macketanz, Shushen Manakhimova, Eleftherios Avramidis, Ekaterina Lapshinova-koltunski, Sergei Bagdasarov, and Sebastian Möller. 2022. Linguistically motivated evaluation of the 2022 state-of-the-art machine translation systems for three language directions. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 432–449, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Shushen Manakhimova, Eleftherios Avramidis, Vivien Macketanz, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2023. Linguistically motivated evaluation of the 2023 state-of-the-art machine translation: Can ChatGPT outperform NMT? In *Proceedings of the Eighth Conference on Machine Translation*, pages 224–245, Singapore. Association for Computational Linguistics.

Ananya Mukherjee and Manish Shrivastava. 2023. IIIT HYD's submission for WMT23 test-suite task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 246–251, Singapore. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Maja Popović. 2019. Evaluating conjunction disambiguation on English-to-German and French-to-German WMT 2019 translation hypotheses. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 464–469, Florence, Italy. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, JosÃ© Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Kateřina Rysová, Magdaléna Rysová, Tomáš Musil, Lucie Poláková, and Ondřej Bojar. 2019. A test suite and manual evaluation of document-level NMT at WMT19. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 455–463, Florence, Italy. Association for Computational Linguistics.

Huacheng Song and Hongzhi Xu. 2024a. Benchmarking the performance of machine translation evaluation metrics with Chinese multiword expressions. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2204–2216, Torino, Italia. ELRA and ICCL.

Huacheng Song and Hongzhi Xu. 2024b. A deep analysis of the impact of multiword expressions and named entities on Chinese-English machine translations. In

*Proceedings of the 2024 Conference on Emprical Methods on Natural Language Processing (EMNLP 2024).*

Gongbo Tang, Philipp Rönchen, Rico Sennrich, and Joakim Nivre. 2021. Revisiting negation in neural machine translation. *Transactions of the Association for Computational Linguistics*, 9:740–755.

Liang Tian, Derek F. Wong, Lidia S. Chao, Paulo Quaresma, Francisco Oliveira, Yi Lu, Shuo Li, Yiming Wang, and Longyue Wang. 2014. UM-corpus: A large English-Chinese parallel corpus for statistical machine translation. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 1837–1842, Reykjavik, Iceland. European Language Resources Association (ELRA).

Hongzhi Xu and Jingxia Lin. 2023. A reference grammar based chinese corpus. In *Symposium on Language and Big Data: Challenges in Chinese Linguistics*. Hong Kong.

Liang Xu, Anqi Li, Lei Zhu, Hang Xue, Changtai Zhu, Kangkang Zhao, Haonan He, Xuanwei Zhang, Qiyue Kang, and Zhenzhong Lan. 2023. Superclue: A comprehensive chinese large language model benchmark.

## A  System Performance on Different Domains

Table 8 to Table 14 show the performance of different systems on different domains in different metrics. In all the following tables, the highest and the lowest scores among all systems evaluated are highlighted with bold and underlined numbers.

|  | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| **Baidu** | 69.1 | 71.0 | 70.5 | 70.6 | 69.3 | 67.2 | 66.0 |
| **Niu** | 69.3 | 73.2 | 70.9 | 71.2 | 69.5 | 67.8 | 66.0 |
| **Google** | **71.0** | **75.3** | **72.6** | **72.4** | **70.5** | 68.6 | **67.0** |
| **DeepL** | 69.0 | 71.9 | 70.9 | 70.2 | 69.1 | **68.8** | 65.0 |
| **Ernie** | 70.1 | 72.8 | 71.4 | 71.7 | 69.9 | 68.3 | **67.0** |
| **Qwen** | _62.7_ | _64.5_ | _64.4_ | _64.6_ | _63.6_ | _63.7_ | _59.5_ |
| **GPT** | 70.2 | 71.0 | 71.4 | 71.5 | 69.7 | 68.4 | 66.4 |
| **Claude** | 70.2 | 72.3 | 71.5 | 71.4 | 69.6 | 68.2 | 66.0 |
| **NMT-AVG** | 69.6 | 72.8 | 71.2 | 71.1 | 69.6 | 68.1 | 66.0 |
| **LLM-AVG** | 68.3 | 70.2 | 69.7 | 69.8 | 68.2 | 67.2 | 64.7 |

Table 8: Performance in AVERAGE scores of each system on different domains.

|  | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| **Baidu** | 20.4 | 25.8 | 22.2 | 24.8 | 20.6 | 19.5 | 17.8 |
| **Niu** | 22.8 | 35.8 | 24.6 | 28.0 | 22.7 | 22.0 | 19.5 |
| **Google** | **24.8** | **40.2** | **28.0** | **29.7** | **23.6** | 22.9 | 19.6 |
| **DeepL** | 22.4 | 32.9 | 24.9 | 25.9 | 22.2 | **25.6** | 16.6 |
| **Ernie** | 22.6 | 31.0 | 23.9 | 26.9 | 21.7 | 21.9 | **19.9** |
| **Qwen** | _15.4_ | _18.8_ | _17.0_ | _18.3_ | _16.0_ | _19.0_ | _11.8_ |
| **GPT** | 22.2 | 25.2 | 23.4 | 26.1 | 21.0 | 22.0 | 17.8 |
| **Claude** | 22.4 | 30.6 | 23.9 | 26.6 | 20.8 | 21.7 | 17.2 |
| **NMT-AVG** | 22.6 | 33.7 | 24.9 | 27.1 | 22.3 | 22.5 | 18.4 |
| **LLM-AVG** | 20.6 | 26.4 | 22.0 | 24.5 | 19.9 | 21.1 | 16.7 |

Table 9: Performance in BLEU scores of each system on different domains.

|  | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| **Baidu** | 55.1 | 59.2 | 56.6 | 59.4 | 52.8 | 48.1 | 59.0 |
| **Niu** | 57.0 | 66.3 | 58.1 | 61.5 | 54.3 | 50.5 | 60.4 |
| **Google** | **58.5** | **69.1** | **60.5** | **62.5** | **55.1** | 50.6 | 60.5 |
| **DeepL** | 56.3 | 64.9 | 58.1 | 59.4 | 53.8 | **52.5** | 57.3 |
| **Ernie** | 57.0 | 63.4 | 58.1 | 61.4 | 53.9 | 50.0 | **61.2** |
| **Qwen** | _45.5_ | _49.7_ | _47.2_ | _49.6_ | _44.3_ | _44.3_ | _47.8_ |
| **GPT** | 57.0 | 58.8 | 57.9 | 61.0 | 53.1 | 50.1 | 59.6 |
| **Claude** | 57.2 | 63.4 | 58.4 | 61.3 | 53.3 | 49.8 | 59.2 |
| **NMT-AVG** | 56.7 | 64.9 | 58.3 | 60.7 | 54.0 | 50.4 | 59.3 |
| **LLM-AVG** | 54.2 | 58.8 | 55.4 | 58.3 | 51.1 | 48.5 | 57.0 |

Table 10: Performance in CHRF scores of each system on different domains.

## B  System Performance on Different Grammatical Features

Table 15 to Table 22 show the performance of different systems on different grammatical feature

|  | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| **Baidu** | 81.3 | 83.3 | 83.0 | 82.0 | 81.3 | 78.2 | 76.1 |
| **Niu** | 81.2 | 84.2 | 83.1 | 82.2 | 81.3 | 78.3 | 76.0 |
| **Google** | **82.8** | **86.0** | **84.5** | **83.0** | **82.2** | **79.2** | 76.5 |
| **DeepL** | 81.0 | 83.2 | 83.1 | 81.1 | 80.8 | 78.8 | 74.9 |
| **Ernie** | 82.3 | 84.4 | 83.8 | 82.5 | 81.7 | 78.8 | **76.6** |
| **Qwen** | _76.4_ | _79.1_ | _78.7_ | _77.2_ | _76.3_ | _74.9_ | _71.7_ |
| **GPT** | 82.4 | 83.3 | 83.9 | 82.4 | 81.6 | 78.8 | 76.5 |
| **Claude** | 82.2 | 84.0 | 83.9 | 82.4 | 81.5 | 78.6 | 76.4 |
| **NMT-AVG** | 81.6 | 84.2 | 83.4 | 82.1 | 81.4 | 78.6 | 75.9 |
| **LLM-AVG** | 80.8 | 82.7 | 82.6 | 81.1 | 80.3 | 77.8 | 75.3 |

Table 11: Performance in COMET scores of each system on different domains.

|  | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| **Baidu** | 90.2 | 90.9 | 91.9 | 91.3 | 92.5 | 90.9 | 84.7 |
| **Niu** | 89.7 | 90.4 | 91.5 | 91.2 | 92.1 | 90.6 | 83.8 |
| **Google** | **91.9** | **92.5** | **93.2** | **92.6** | **93.4** | **91.8** | 85.9 |
| **DeepL** | 89.3 | 89.1 | 91.3 | 90.5 | 91.7 | 91.3 | 83.4 |
| **Ernie** | 91.4 | 92.0 | 92.9 | 92.4 | 93.2 | **91.8** | **86.0** |
| **Qwen** | _84.9_ | _85.4_ | _86.8_ | _87.0_ | _88.1_ | _87.6_ | _78.6_ |
| **GPT** | 91.6 | 91.3 | 93.1 | 92.3 | 93.1 | **91.8** | 85.3 |
| **Claude** | 91.5 | 91.2 | 92.9 | 91.9 | 92.9 | 91.6 | 85.1 |
| **NMT-AVG** | 90.3 | 90.7 | 92.0 | 91.4 | 92.4 | 91.2 | 84.5 |
| **LLM-AVG** | 89.8 | 90.0 | 91.4 | 90.9 | 91.8 | 90.7 | 83.8 |

Table 12: Performance in XCOMET scores of each system on different domains.

|  | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| **Baidu** | 73.8 | **74.3** | 74.7 | 72.4 | 72.8 | 70.7 | 70.1 |
| **Niu** | 72.7 | 71.6 | 74.0 | 71.5 | 72.0 | 70.1 | 68.9 |
| **Google** | 74.4 | 72.7 | 75.2 | 72.7 | 72.8 | 71.5 | 70.4 |
| **DeepL** | 72.6 | 71.4 | 74.2 | 71.8 | 71.6 | 70.1 | 69.7 |
| **Ernie** | 74.1 | 73.7 | 75.3 | 73.0 | 73.2 | 71.9 | 70.1 |
| **Qwen** | _65.2_ | _65.8_ | _66.9_ | _65.1_ | _64.8_ | _64.1_ | _63.1_ |
| **GPT** | **74.5** | **74.3** | **75.5** | 73.3 | 73.5 | **72.2** | **70.7** |
| **Claude** | 74.4 | 72.7 | 75.3 | 72.8 | **73.6** | 72.0 | 70.3 |
| **NMT-AVG** | 73.4 | 72.5 | 74.5 | 72.1 | 72.3 | 70.6 | 69.8 |
| **LLM-AVG** | 72.1 | 71.6 | 73.2 | 71.0 | 71.3 | 70.0 | 68.5 |

Table 13: Performance in COMETKIWI-QE scores of each system on different domains.

|  | Edu | Laws | News | Sci | Spk | Sbt | Ths |
|---|---|---|---|---|---|---|---|
| **Baidu** | 93.5 | 92.9 | 94.4 | 93.7 | 95.6 | 95.4 | 88.4 |
| **Niu** | 92.4 | 90.8 | 93.7 | 93.0 | 94.9 | 94.9 | 87.0 |
| **Google** | **93.7** | 91.3 | 94.5 | 93.9 | 95.6 | **95.6** | **88.7** |
| **DeepL** | 92.2 | 89.6 | 93.7 | 92.7 | 94.5 | 94.5 | 87.7 |
| **Ernie** | 93.4 | 92.4 | 94.6 | 93.8 | 95.6 | 95.4 | 88.2 |
| **Qwen** | _88.4_ | _88.1_ | _89.9_ | _90.0_ | _91.9_ | _92.1_ | _84.1_ |
| **GPT** | **93.7** | **93.3** | **94.8** | **94.0** | **95.8** | **95.6** | 88.4 |
| **Claude** | 93.6 | 92.3 | 94.5 | 93.6 | 95.5 | 95.2 | 87.9 |
| **NMT-AVG** | 93.0 | 91.2 | 94.1 | 93.3 | 95.2 | 95.1 | 88.0 |
| **LLM-AVG** | 92.3 | 91.5 | 93.5 | 92.8 | 94.7 | 94.6 | 87.2 |

Table 14: Performance in XCOMET-QE scores of each system on different domains.

groups in different metrics.

|  | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **ZAIProg** | 94.1 | 93.5 | 94.3 | 93.3 | **94.6** | <u>90.2</u> | **94.6** | 94.2 | 93.8 | 93.4 | 93.6 |
| **MEINeg** | 93.5 | 93.0 | 94.6 | 93.1 | **94.8** | <u>89.5</u> | 94.4 | 94.5 | 93.6 | 93.3 | 93.4 |
| **NPI** | 93.6 | 93.5 | 93.8 | 92.5 | **94.1** | <u>89.8</u> | 93.9 | 94.0 | 93.3 | 92.9 | 93.1 |
| **Recp** | 93.3 | 92.5 | **94.1** | 92.6 | 93.6 | <u>89.7</u> | 93.9 | 93.5 | 93.1 | 92.7 | 92.9 |
| **ZHIQtf** | 93.1 | 92.5 | **93.8** | 92.7 | 93.7 | <u>88.6</u> | **93.8** | 93.3 | 93.0 | 92.4 | 92.7 |
| **ImpNeg** | 92.9 | 91.6 | 93.3 | 92.5 | 93.1 | <u>88.2</u> | **93.6** | 93.0 | 92.6 | 92.0 | 92.3 |
| **VP** | 92.7 | 92.0 | **93.5** | 92.1 | 93.0 | <u>88.5</u> | 93.2 | 93.1 | 92.6 | 91.9 | 92.3 |
| **SHICop** | 93.1 | 91.8 | 93.1 | 91.3 | 93.3 | <u>88.6</u> | **93.7** | 93.0 | 92.3 | 92.1 | 92.2 |
| **Refl** | 92.8 | 91.4 | 93.6 | 91.5 | 93.4 | <u>87.9</u> | 93.5 | **93.7** | 92.3 | 92.1 | 92.2 |
| **Pass** | 92.9 | 91.4 | **93.4** | 91.8 | 93.0 | <u>88.5</u> | 93.0 | 93.3 | 92.4 | 92.0 | 92.2 |
| **GUOPrf** | 92.5 | 91.5 | 93.0 | 91.8 | 93.0 | <u>88.5</u> | **93.2** | 92.7 | 92.2 | 91.8 | 92.0 |
| **StdCLF** | 92.9 | 92.1 | 92.9 | 91.3 | 93.1 | <u>87.8</u> | **93.3** | 92.8 | 92.3 | 91.8 | 92.0 |
| **Deixis** | 92.7 | 91.8 | 93.2 | 90.7 | 93.0 | <u>88.1</u> | **93.3** | 92.9 | 92.1 | 91.8 | 92.0 |
| **ConcCpl** | 92.5 | 91.5 | 93.4 | 91.6 | 92.9 | <u>85.9</u> | **93.5** | 93.1 | 92.2 | 91.3 | 91.8 |
| **AdvP** | 92.5 | 91.4 | 93.0 | 91.1 | 92.7 | <u>87.5</u> | **93.1** | 92.9 | 92.0 | 91.5 | 91.8 |
| **CausCpl** | 92.9 | 90.7 | 93.5 | 90.0 | 93.2 | <u>86.2</u> | 93.7 | **93.8** | 91.8 | 91.7 | 91.8 |
| **SpanPP** | 92.2 | 91.2 | 92.9 | 91.0 | **93.2** | <u>87.9</u> | 93.0 | 92.3 | 91.8 | 91.6 | 91.7 |
| **VerbLE** | 92.3 | 91.3 | 92.7 | 91.6 | 92.3 | <u>87.4</u> | **93.1** | 92.6 | 92.0 | 91.3 | 91.7 |
| **RefPP** | 92.1 | 91.2 | **92.5** | 90.9 | 92.4 | <u>88.1</u> | 92.3 | 91.9 | 91.7 | 91.2 | 91.4 |
| **PstVPP** | **92.5** | 91.5 | 92.2 | 90.5 | 92.4 | <u>86.7</u> | 92.4 | 92.4 | 91.7 | 91.0 | 91.3 |
| **GoalPP** | **92.5** | 91.6 | 92.4 | 90.0 | 92.4 | <u>86.9</u> | **92.5** | 92.2 | 91.6 | 91.0 | 91.3 |
| **IndCLF** | 91.7 | 91.1 | 92.2 | 90.4 | 92.2 | <u>88.2</u> | **92.3** | 91.8 | 91.3 | 91.1 | 91.2 |
| **CondPP** | 92.1 | 91.3 | **92.6** | 89.7 | 92.3 | <u>87.3</u> | 92.3 | 91.7 | 91.4 | 90.9 | 91.2 |
| **BUNeg** | 91.6 | 90.8 | 92.3 | 90.9 | **92.4** | <u>86.8</u> | 92.3 | 92.0 | 91.4 | 90.9 | 91.2 |
| **BA** | 91.9 | 91.6 | **92.2** | 91.2 | 92.0 | <u>86.2</u> | **92.2** | 91.9 | 91.7 | 90.6 | 91.1 |
| **TopPP** | 91.5 | 91.0 | **92.2** | 90.9 | **92.2** | <u>87.1</u> | **92.2** | 91.7 | 91.4 | 90.8 | 91.1 |
| **SrcPP** | 91.3 | 91.1 | **92.3** | 90.6 | 92.2 | <u>86.5</u> | 92.2 | 91.9 | 91.3 | 90.7 | 91.0 |
| **SentIPP** | 92.0 | 91.2 | **92.4** | 90.4 | 92.1 | <u>86.4</u> | 92.0 | 91.6 | 91.5 | 90.5 | 91.0 |
| **SpcPP** | 91.3 | 91.2 | **92.2** | 90.7 | 91.7 | <u>87.4</u> | 91.8 | 91.8 | 91.3 | 90.7 | 91.0 |
| **EvCLF** | 91.4 | 90.8 | **91.9** | 91.3 | 91.5 | <u>86.9</u> | 91.6 | **91.9** | 91.4 | 90.5 | 90.9 |
| **DirPP** | 91.5 | 91.0 | 91.8 | 89.3 | 92.0 | <u>87.1</u> | **92.2** | 91.6 | 90.9 | 90.7 | 90.8 |
| **AdjP** | 91.3 | 90.3 | 92.0 | 90.0 | 91.9 | <u>85.8</u> | **92.2** | 91.6 | 90.9 | 90.4 | 90.7 |
| **Cmpr** | 91.0 | 90.2 | 91.9 | 90.0 | **92.1** | <u>85.6</u> | 92.0 | 92.0 | 90.8 | 90.4 | 90.6 |
| **NP** | 91.3 | 90.5 | **91.7** | 89.8 | **91.7** | <u>86.6</u> | **91.7** | 91.3 | 90.8 | 90.3 | 90.6 |
| **TmpSCpl** | 91.5 | 90.7 | **91.8** | 89.8 | 91.6 | <u>84.9</u> | 91.6 | 91.6 | 91.0 | 89.9 | 90.4 |
| **LocPP** | 90.7 | 90.6 | 91.4 | 90.2 | 91.4 | <u>85.0</u> | **91.6** | 91.5 | 90.7 | 89.9 | 90.3 |
| **PreVPP** | 90.9 | 90.7 | **91.8** | 89.4 | 91.0 | <u>86.1</u> | 91.2 | 90.6 | 90.7 | 89.7 | 90.2 |
| **AgtPP** | 90.8 | 90.0 | **91.5** | 89.6 | **91.5** | <u>85.5</u> | 91.4 | 90.9 | 90.5 | 89.8 | 90.1 |
| **Rel** | 91.0 | 89.4 | 91.6 | 89.2 | 91.3 | <u>85.0</u> | **91.7** | 91.2 | 90.3 | 89.8 | 90.1 |
| **LexNeg** | 90.5 | 89.8 | **91.2** | 89.9 | 90.9 | <u>86.1</u> | 91.0 | 91.1 | 90.3 | 89.8 | 90.1 |
| **KindCLF** | 90.6 | 89.7 | **90.9** | 89.4 | **90.9** | <u>85.8</u> | 90.8 | 90.3 | 90.2 | 89.5 | 89.8 |
| **ApprCLF** | 90.5 | 90.0 | 91.1 | 89.0 | **91.3** | <u>84.9</u> | 91.1 | 90.7 | 90.2 | 89.5 | 89.8 |
| **PathPP** | 87.0 | 85.5 | 86.8 | 85.1 | 86.7 | <u>81.9</u> | 86.7 | **87.3** | 86.1 | 85.7 | 85.9 |

Table 15: Performance in X-AVERAGE scores of each system on different linguistic features.

| | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **CondPP** | 70.9 | 72.7 | **74.1** | 70.8 | 72.1 | 64.1 | 70.9 | 71.5 | 72.1 | 69.7 | 70.9 |
| **SpanPP** | 70.7 | 71.2 | **72.7** | 70.7 | 72.3 | 65.0 | 71.7 | 71.6 | 71.3 | 70.2 | 70.7 |
| **ZAIProg** | 70.6 | 71.3 | **72.9** | 71.0 | 71.8 | 65.4 | 71.3 | 71.1 | 71.4 | 69.9 | 70.7 |
| **SrcPP** | 70.6 | 71.2 | **72.5** | 69.9 | 71.7 | 65.0 | 71.4 | 71.7 | 71.1 | 70.0 | 70.5 |
| **SpcPP** | 69.8 | 71.2 | **72.7** | 71.0 | 71.2 | 65.0 | 70.8 | 70.9 | 71.2 | 69.5 | 70.3 |
| **LocPP** | 70.0 | 72.1 | **73.0** | 70.7 | 71.3 | 63.5 | 70.7 | 71.2 | 71.5 | 69.2 | 70.3 |
| **ImpNeg** | 70.3 | 70.2 | **73.3** | 70.8 | 71.1 | 63.9 | 70.9 | 71.0 | 71.2 | 69.2 | 70.2 |
| **GUOPrf** | 70.0 | 70.4 | **71.6** | 70.1 | 71.1 | 65.3 | 71.2 | 70.9 | 70.5 | 69.6 | 70.1 |
| **NPI** | 69.9 | 71.2 | **72.3** | 70.7 | 70.9 | 63.7 | 70.3 | 71.1 | 71.0 | 69.0 | 70.0 |
| **IndCLF** | 69.8 | 70.4 | **72.2** | 70.3 | 70.9 | 65.0 | 70.7 | 70.7 | 70.7 | 69.3 | 70.0 |
| **MEINeg** | 69.5 | 70.4 | **72.2** | 70.5 | 71.4 | 63.7 | 70.8 | 70.9 | 70.7 | 69.2 | 69.9 |
| **TmpSCpl** | 70.2 | 70.9 | **71.7** | 70.3 | 71.0 | 63.8 | 70.5 | 71.0 | 70.8 | 69.1 | 69.9 |
| **Recp** | 69.9 | 70.4 | **72.0** | 70.3 | 70.8 | 64.1 | 70.6 | 71.0 | 70.7 | 69.1 | 69.9 |
| **StdCLF** | 69.8 | 71.1 | **71.9** | 70.0 | 71.7 | 64.1 | 70.2 | 70.1 | 70.7 | 69.0 | 69.9 |
| **Pass** | 70.0 | 70.6 | **72.0** | 70.3 | 70.5 | 64.3 | 69.8 | 70.7 | 70.7 | 68.8 | 69.8 |
| **BA** | 69.6 | 71.0 | **72.1** | 70.4 | 70.6 | 63.0 | 70.5 | 70.2 | 70.8 | 68.6 | 69.7 |
| **RefPP** | 69.9 | 70.3 | **71.4** | 69.6 | 70.4 | 64.3 | 70.0 | 70.2 | 70.3 | 68.7 | 69.5 |
| **EvCLF** | 69.7 | 69.8 | **71.3** | 70.5 | 70.4 | 63.5 | 69.7 | 70.4 | 70.3 | 68.5 | 69.4 |
| **PstVPP** | 69.4 | 70.0 | **70.9** | 69.6 | 70.2 | 63.4 | 70.2 | 70.7 | 70.0 | 68.6 | 69.3 |
| **LexNeg** | 68.9 | 69.5 | **71.3** | 69.6 | 70.7 | 64.0 | 69.9 | 70.4 | 69.8 | 68.8 | 69.3 |
| **ZHIQtf** | 69.3 | 69.7 | **71.1** | 70.3 | 70.1 | 63.4 | 70.3 | 70.0 | 70.1 | 68.5 | 69.3 |
| **Rel** | 69.3 | 70.1 | **72.1** | 69.4 | 70.3 | 62.6 | 69.4 | 70.5 | 70.2 | 68.2 | 69.2 |
| **NP** | 69.2 | 70.0 | **71.5** | 69.4 | 70.3 | 63.4 | 69.7 | 70.0 | 70.0 | 68.3 | 69.2 |
| **VP** | 69.1 | 69.9 | **71.5** | 69.4 | 70.1 | 63.6 | 69.8 | 70.1 | 70.0 | 68.4 | 69.2 |
| **AdvP** | 69.1 | 69.6 | **71.3** | 69.3 | 70.3 | 62.8 | 69.9 | 70.2 | 69.8 | 68.3 | 69.1 |
| **Deixis** | 69.1 | 69.7 | **71.3** | 69.0 | 70.3 | 63.6 | 69.7 | 69.9 | 69.8 | 68.4 | 69.1 |
| **GoalPP** | 69.0 | 69.5 | **70.7** | 69.1 | 69.9 | 63.5 | 69.9 | 70.3 | 69.6 | 68.4 | 69.0 |
| **VerbLE** | 69.1 | 69.6 | **70.6** | 69.6 | 70.0 | 63.2 | 69.9 | 69.7 | 69.7 | 68.2 | 69.0 |
| **CausCpl** | 69.1 | 68.8 | **71.5** | 68.3 | 70.7 | 61.8 | 70.3 | 70.8 | 69.4 | 68.4 | 68.9 |
| **TopPP** | 68.7 | 69.5 | **70.9** | 69.9 | 70.0 | 62.2 | 69.7 | 69.7 | 69.8 | 67.9 | 68.8 |
| **SHICop** | 68.9 | 69.3 | **70.8** | 68.8 | 70.0 | 63.3 | 69.8 | 69.3 | 69.5 | 68.1 | 68.8 |
| **DirPP** | 68.3 | 69.8 | **70.4** | 67.4 | 69.3 | 62.8 | 69.2 | 69.4 | 69.0 | 67.7 | 68.3 |
| **PtcpPP** | 68.4 | 69.1 | **70.6** | 68.5 | 69.2 | 62.5 | 68.7 | 68.9 | 69.2 | 67.3 | 68.3 |
| **Refl** | 68.2 | 68.2 | **70.1** | 68.3 | 69.3 | 62.4 | 69.1 | 69.7 | 68.7 | 67.6 | 68.2 |
| **Cmpr** | 67.9 | 68.4 | **70.3** | 68.5 | 69.6 | 61.3 | 68.7 | 69.0 | 68.8 | 67.1 | 68.0 |
| **PreVPP** | 68.2 | 68.9 | **70.3** | 67.7 | 69.0 | 61.9 | 68.8 | 68.9 | 68.8 | 67.2 | 68.0 |
| **BUNeg** | 67.8 | 68.8 | **70.2** | 68.3 | 69.0 | 61.7 | 68.4 | 68.3 | 68.8 | 66.8 | 67.8 |
| **KindCLF** | 67.5 | 68.4 | **69.6** | 67.8 | 68.7 | 62.1 | 68.3 | 68.4 | 68.3 | 66.9 | 67.6 |
| **AdjP** | 67.7 | 68.4 | **69.7** | 68.0 | 68.8 | 61.1 | 68.3 | 68.5 | 68.5 | 66.7 | 67.6 |
| **SentIPP** | 68.0 | 68.2 | **69.0** | 67.4 | 68.1 | 60.8 | 67.8 | 67.7 | 68.2 | 66.1 | 67.1 |
| **ConcCpl** | 67.2 | 67.3 | **69.0** | 67.1 | 68.1 | 60.1 | 68.5 | 68.1 | 67.7 | 66.2 | 66.9 |
| **ApprCLF** | 67.2 | 67.4 | **69.2** | 66.9 | 67.6 | 60.1 | 67.6 | 68.1 | 67.7 | 65.8 | 66.7 |
| **PathPP** | 66.4 | 66.8 | **67.9** | 66.0 | 67.6 | 60.1 | 66.1 | 67.1 | 66.8 | 65.2 | 66.0 |

Table 16: Performance in AVERAGE scores of each system on different grammatical features.

| | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **CondPP** | 25.1 | 33.4 | **36.2** | 29.7 | 28.6 | <u>17.2</u> | 25.1 | 28.0 | 31.1 | 24.7 | 27.9 |
| **LocPP** | 24.4 | 31.8 | **33.1** | 28.3 | 27.5 | <u>18.4</u> | 25.7 | 27.6 | 29.4 | 24.8 | 27.1 |
| **SrcPP** | 25.6 | 28.7 | **30.9** | 26.4 | 27.8 | <u>21.2</u> | 26.2 | 27.4 | 27.9 | 25.6 | 26.8 |
| **IndCLF** | 24.2 | 27.5 | **31.0** | 28.5 | 26.7 | <u>19.6</u> | 25.7 | 26.3 | 27.8 | 24.6 | 26.2 |
| **SpcPP** | 23.4 | 28.3 | **30.8** | 28.6 | 26.9 | <u>19.1</u> | 25.3 | 26.0 | 27.8 | 24.3 | 26.0 |
| **SpanPP** | 24.5 | 27.7 | **30.0** | 26.5 | 28.1 | <u>18.4</u> | 26.3 | 26.6 | 27.2 | 24.9 | 26.0 |
| **TmpSCpl** | 24.2 | 27.9 | **29.1** | 27.4 | 26.2 | <u>18.9</u> | 24.4 | 26.2 | 27.1 | 23.9 | 25.5 |
| **ImpNeg** | 23.4 | 25.2 | **32.6** | 26.2 | 25.9 | <u>17.5</u> | 24.6 | 26.3 | 26.8 | 23.6 | 25.2 |
| **GUOPrf** | 23.2 | 26.4 | **27.9** | 26.2 | 25.3 | <u>20.8</u> | 25.6 | 25.4 | 25.9 | 24.3 | 25.1 |
| **EvCLF** | 24.8 | 26.1 | **28.9** | 27.1 | 26.4 | <u>18.3</u> | 23.8 | 25.1 | 26.7 | 23.4 | 25.1 |
| **PstVPP** | 22.2 | 26.1 | **27.4** | 26.8 | 24.6 | <u>18.9</u> | 24.6 | 25.9 | 25.6 | 23.5 | 24.6 |
| **NP** | 22.4 | 26.7 | **29.5** | 26.1 | 25.3 | <u>17.6</u> | 23.1 | 24.6 | 26.2 | 22.6 | 24.4 |
| **Pass** | 23.0 | 28.0 | **28.2** | 26.1 | 23.6 | <u>19.3</u> | 21.5 | 23.9 | 26.3 | 22.1 | 24.2 |
| **Rel** | 22.3 | 27.7 | **30.6** | 25.4 | 24.3 | <u>16.4</u> | 21.1 | 25.2 | 26.5 | 21.8 | 24.1 |
| **TopPP** | 21.6 | 25.5 | **28.1** | 27.4 | 25.2 | <u>15.7</u> | 23.9 | 24.7 | 25.6 | 22.4 | 24.0 |
| **NPI** | 21.7 | 26.6 | **29.4** | 26.8 | 24.1 | <u>15.1</u> | 22.3 | 24.8 | 26.1 | 21.6 | 23.9 |
| **StdCLF** | 21.1 | 28.0 | **28.3** | 25.1 | 26.0 | <u>17.7</u> | 21.8 | 22.8 | 25.6 | 22.1 | 23.8 |
| **LexNeg** | 21.1 | 25.2 | **28.5** | 24.8 | 25.5 | <u>17.5</u> | 23.1 | 24.8 | 24.9 | 22.7 | 23.8 |
| **Deixis** | 22.1 | 25.8 | **28.7** | 25.2 | 25.0 | <u>17.2</u> | 22.4 | 23.9 | 25.5 | 22.1 | 23.8 |
| **GoalPP** | 21.2 | 24.5 | **26.1** | 25.9 | 23.4 | <u>19.3</u> | 23.5 | 25.0 | 24.4 | 22.8 | 23.6 |
| **ZAIProg** | 22.1 | 25.4 | **28.9** | 25.1 | 24.3 | <u>16.8</u> | 23.1 | 23.1 | 25.4 | 21.8 | 23.6 |
| **Recp** | 22.3 | 25.4 | **27.3** | 24.4 | 24.1 | <u>16.6</u> | 22.7 | 24.7 | 24.9 | 22.0 | 23.4 |
| **MEINeg** | 20.7 | 25.1 | **27.3** | 24.7 | 25.0 | <u>16.4</u> | 23.0 | 23.4 | 24.4 | 22.0 | 23.2 |
| **BA** | 21.0 | 25.7 | **28.5** | 25.4 | 23.6 | <u>15.7</u> | 23.2 | 22.6 | 25.1 | 21.3 | 23.2 |
| **DirPP** | 20.8 | **27.0** | **27.0** | 21.8 | 22.0 | <u>17.7</u> | 21.8 | 23.2 | 24.1 | 21.2 | 22.7 |
| **KindCLF** | 20.5 | 25.2 | **26.9** | 23.7 | 23.4 | <u>16.6</u> | 21.7 | 23.1 | 24.1 | 21.2 | 22.6 |
| **ZHIQtf** | 21.3 | 24.0 | **25.7** | 25.4 | 21.9 | <u>17.0</u> | 22.4 | 22.2 | 24.1 | 20.9 | 22.5 |
| **RefPP** | 21.7 | 24.9 | **26.3** | 23.1 | 23.0 | <u>16.0</u> | 21.4 | 22.4 | 24.0 | 20.7 | 22.4 |
| **VP** | 20.8 | 24.8 | **26.9** | 23.0 | 22.4 | <u>15.5</u> | 21.3 | 22.3 | 23.9 | 20.4 | 22.1 |
| **AdvP** | 20.6 | 23.9 | **26.5** | 23.7 | 23.2 | <u>14.6</u> | 21.5 | 22.7 | 23.7 | 20.5 | 22.1 |
| **PreVPP** | 21.0 | 23.6 | **25.5** | 21.7 | 23.1 | <u>15.1</u> | 22.1 | 23.1 | 22.9 | 20.9 | 21.9 |
| **PathPP** | 20.6 | 23.9 | **25.8** | 22.4 | 24.6 | <u>15.0</u> | 19.5 | 21.4 | 23.2 | 20.1 | 21.7 |
| **SHICop** | 19.6 | 23.3 | **25.4** | 22.9 | 22.6 | <u>15.5</u> | 20.9 | 21.1 | 22.8 | 20.0 | 21.4 |
| **ApprCLF** | 20.3 | 21.7 | **26.4** | 22.2 | 20.5 | <u>14.5</u> | 20.9 | 22.6 | 22.7 | 19.6 | 21.1 |
| **BUNeg** | 19.8 | 24.2 | **25.9** | 23.1 | 21.6 | <u>14.7</u> | 19.6 | 20.1 | 23.2 | 19.0 | 21.1 |
| **Cmpr** | 19.5 | 22.6 | **25.7** | 23.0 | 22.8 | <u>14.5</u> | 19.9 | 20.4 | 22.7 | 19.4 | 21.0 |
| **VerbLE** | 19.7 | 23.3 | **23.6** | 22.3 | 22.1 | <u>15.3</u> | 20.5 | 20.5 | 22.2 | 19.6 | 20.9 |
| **PtcpPP** | 19.7 | 23.3 | **25.7** | 22.4 | 21.1 | <u>14.3</u> | 19.2 | 20.8 | 22.8 | 18.9 | 20.8 |
| **CausCpl** | 19.1 | 21.3 | **25.3** | 20.3 | 23.8 | <u>12.4</u> | 21.3 | 22.5 | 21.5 | 20.0 | 20.8 |
| **Refl** | 18.7 | 20.5 | **22.4** | 21.6 | 21.0 | <u>14.4</u> | 19.8 | 21.5 | 20.8 | 19.2 | 20.0 |
| **AdjP** | 17.1 | 20.7 | **22.5** | 20.1 | 19.3 | <u>12.6</u> | 17.1 | 18.4 | 20.1 | 16.9 | 18.5 |
| **SentIPP** | 17.1 | **19.4** | 19.2 | 17.9 | 17.5 | <u>11.5</u> | 16.1 | 16.9 | 18.4 | 15.5 | 16.9 |
| **ConcCpl** | 14.5 | 17.0 | **18.4** | 16.6 | 16.0 | <u>10.1</u> | 16.6 | 15.9 | 16.6 | 14.7 | 15.6 |

Table 17: Performance in BLEU scores of each system on different grammatical features.

|  | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| LocPP | 60.6 | 66.1 | **66.7** | 62.5 | 63.0 | <u>51.9</u> | 61.0 | 62.8 | 64.0 | 59.7 | 61.8 |
| CondPP | 60.8 | 65.6 | **67.1** | 63.6 | 63.4 | <u>49.7</u> | 60.9 | 62.7 | 64.3 | 59.2 | 61.7 |
| SrcPP | 60.2 | 62.6 | **63.8** | 59.8 | 62.7 | <u>52.8</u> | 61.8 | 63.0 | 61.6 | 60.1 | 60.9 |
| LexNeg | 59.5 | 61.6 | **63.6** | 61.0 | 63.1 | <u>51.8</u> | 60.9 | 62.0 | 61.4 | 59.5 | 60.4 |
| SpcPP | 57.8 | 61.7 | **63.5** | 61.3 | 60.8 | <u>51.5</u> | 59.8 | 61.0 | 61.1 | 58.3 | 59.7 |
| IndCLF | 58.1 | 60.1 | **62.5** | 60.3 | 60.0 | <u>50.6</u> | 58.8 | 59.7 | 60.2 | 57.3 | 58.8 |
| SpanPP | 57.7 | 60.3 | **61.4** | 59.4 | 60.7 | <u>48.7</u> | 59.1 | 59.7 | 59.7 | 57.0 | 58.4 |
| TmpSCpl | 57.1 | 59.9 | **60.3** | 60.0 | 59.4 | <u>48.9</u> | 57.8 | 59.2 | 59.3 | 56.3 | 57.8 |
| RefPP | 57.9 | 60.3 | **61.0** | 59.0 | 59.0 | <u>48.2</u> | 57.7 | 59.1 | 59.5 | 56.0 | 57.8 |
| Rel | 56.8 | 60.8 | **62.6** | 59.9 | 59.1 | <u>47.2</u> | 56.2 | 59.6 | 60.0 | 55.5 | 57.8 |
| Recp | 57.3 | 59.7 | **61.2** | 59.1 | 59.0 | <u>47.7</u> | 58.1 | 59.5 | 59.3 | 56.1 | 57.7 |
| BA | 56.8 | 60.5 | **62.1** | 59.3 | 58.9 | <u>47.0</u> | 57.8 | 58.2 | 59.7 | 55.5 | 57.6 |
| EvCLF | 57.2 | 58.5 | **60.3** | 59.8 | 59.2 | <u>48.6</u> | 57.5 | 59.0 | 59.0 | 56.1 | 57.5 |
| TopPP | 56.7 | 59.7 | **61.1** | 60.2 | 58.8 | <u>46.5</u> | 58.1 | 58.9 | 59.4 | 55.6 | 57.5 |
| NP | 56.6 | 59.7 | **61.5** | 59.0 | 59.1 | <u>48.2</u> | 57.2 | 58.6 | 59.2 | 55.8 | 57.5 |
| GUOPrf | 56.7 | 58.9 | **59.6** | 58.3 | 59.0 | <u>49.6</u> | 58.7 | 58.8 | 58.4 | 56.5 | 57.4 |
| ImpNeg | 56.5 | 59.3 | **62.7** | 58.4 | 58.3 | <u>47.0</u> | 56.8 | 58.5 | 59.2 | 55.1 | 57.2 |
| StdCLF | 55.8 | 60.0 | **60.2** | 57.7 | 60.1 | <u>47.9</u> | 56.1 | 56.4 | 58.4 | 55.1 | 56.8 |
| ZAIProg | 56.2 | 58.4 | **60.8** | 57.7 | 58.2 | <u>47.2</u> | 57.1 | 57.2 | 58.3 | 54.9 | 56.6 |
| Pass | 55.7 | 59.4 | **59.9** | 57.8 | 57.8 | <u>47.8</u> | 56.0 | 57.6 | 58.2 | 54.8 | 56.5 |
| PstVPP | 55.2 | 58.3 | **58.8** | 58.8 | 57.0 | <u>47.5</u> | 57.2 | 58.5 | 57.8 | 55.0 | 56.4 |
| PtcpPP | 55.6 | 58.7 | **59.9** | 57.6 | 57.3 | <u>47.4</u> | 55.7 | 57.5 | 58.0 | 54.5 | 56.2 |
| KindCLF | 55.3 | 58.6 | **59.3** | 57.0 | 57.8 | <u>47.1</u> | 56.5 | 57.4 | 57.5 | 54.7 | 56.1 |
| PreVPP | 55.5 | 57.6 | **59.6** | 57.0 | 57.6 | <u>46.6</u> | 56.8 | 57.9 | 57.4 | 54.7 | 56.1 |
| GoalPP | 54.7 | 57.2 | **58.4** | 58.3 | 56.4 | <u>47.7</u> | 57.0 | 58.2 | 57.2 | 54.8 | 56.0 |
| NPI | 54.7 | 58.0 | **60.5** | 57.8 | 57.1 | <u>43.8</u> | 55.5 | 57.8 | 57.8 | 53.5 | 55.7 |
| PathPP | 55.4 | 58.6 | 58.8 | 56.7 | **58.9** | <u>45.0</u> | 54.9 | 56.3 | 57.4 | 53.8 | 55.6 |
| AdvP | 54.5 | 57.0 | **58.7** | 56.6 | 56.9 | <u>44.3</u> | 55.6 | 56.8 | 56.7 | 53.4 | 55.1 |
| Deixis | 54.1 | 56.7 | **58.3** | 56.5 | 56.6 | <u>46.0</u> | 54.7 | 55.7 | 56.4 | 53.2 | 54.8 |
| VP | 53.9 | 56.9 | **58.3** | 55.6 | 56.1 | <u>45.9</u> | 55.2 | 56.2 | 56.2 | 53.3 | 54.8 |
| DirPP | 54.0 | 57.7 | **58.4** | 55.0 | 55.6 | <u>46.6</u> | 54.6 | 55.9 | 56.3 | 53.2 | 54.7 |
| VerbLE | 54.4 | 56.6 | **57.0** | 56.2 | 56.2 | <u>44.6</u> | 55.4 | 55.6 | 56.0 | 53.0 | 54.5 |
| MEINeg | 53.3 | 56.4 | **58.0** | 56.5 | 56.0 | <u>44.4</u> | 55.5 | 55.8 | 56.0 | 52.9 | 54.5 |
| ZHIQtf | 54.0 | 55.9 | **57.2** | 56.5 | 55.3 | <u>44.9</u> | 55.7 | 55.8 | 55.9 | 52.9 | 54.4 |
| CausCpl | 53.5 | 55.5 | **57.9** | 54.9 | 56.5 | <u>43.7</u> | 54.6 | 56.0 | 55.5 | 52.7 | 54.1 |
| SHICop | 53.3 | 56.2 | **57.2** | 55.3 | 55.3 | <u>45.1</u> | 54.3 | 54.3 | 55.5 | 52.2 | 53.9 |
| Cmpr | 53.1 | 55.4 | **57.4** | 55.8 | 55.7 | <u>43.1</u> | 53.9 | 54.5 | 55.4 | 51.8 | 53.6 |
| BUNeg | 53.0 | 56.1 | **57.3** | 54.9 | 55.1 | <u>44.4</u> | 53.6 | 54.4 | 55.3 | 51.9 | 53.6 |
| AdjP | 52.7 | 55.8 | **57.0** | 55.0 | 55.0 | <u>42.7</u> | 53.3 | 54.6 | 55.1 | 51.4 | 53.3 |
| SentIPP | 53.5 | 55.0 | **55.5** | 54.2 | 53.1 | <u>41.7</u> | 52.9 | 52.9 | 54.5 | 50.2 | 52.4 |
| Refl | 51.2 | 53.3 | **54.7** | 53.2 | 52.7 | <u>42.4</u> | 52.4 | 53.4 | 53.1 | 50.2 | 51.7 |
| ApprCLF | 51.9 | 52.5 | **55.4** | 53.1 | 51.9 | <u>41.8</u> | 51.9 | 53.8 | 53.2 | 49.8 | 51.5 |
| ConcCpl | 49.8 | 51.3 | **52.7** | 49.8 | 51.0 | <u>39.2</u> | 51.4 | 51.2 | 50.9 | 48.2 | 49.6 |

Table 18: Performance in CHRF scores of each system on different grammatical features.

| | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **ZAIProg** | 83.1 | 83.3 | **84.5** | 82.8 | 84.0 | <u>79.0</u> | 83.3 | 83.1 | 83.4 | 82.3 | 82.9 |
| **SpanPP** | 82.9 | 83.0 | **84.0** | 82.7 | 83.7 | <u>78.8</u> | 83.3 | 83.5 | 83.2 | 82.3 | 82.7 |
| **ImpNeg** | 83.2 | 82.7 | **85.1** | 82.8 | 83.3 | <u>77.8</u> | 83.1 | 83.1 | 83.5 | 81.8 | 82.6 |
| **GUOPrf** | 82.2 | 82.3 | **83.2** | 81.8 | **83.2** | <u>78.5</u> | 83.1 | 83.0 | 82.4 | 81.9 | 82.2 |
| **StdCLF** | 82.1 | 82.4 | **83.7** | 81.8 | 83.6 | <u>78.1</u> | 82.5 | 82.4 | 82.5 | 81.7 | 82.1 |
| **Pass** | 82.3 | 82.5 | **83.9** | 82.3 | 82.4 | <u>77.1</u> | 82.4 | 83.1 | 82.8 | 81.2 | 82.0 |
| **BA** | 82.2 | 82.9 | **83.8** | 82.2 | 82.5 | <u>77.2</u> | 82.6 | 82.2 | 82.8 | 81.1 | 82.0 |
| **CondPP** | 82.0 | 83.0 | **84.2** | 81.6 | 82.8 | <u>77.4</u> | 82.0 | 82.4 | 82.7 | 81.2 | 81.9 |
| **TmpSCpl** | 82.1 | 82.3 | **83.1** | 82.0 | 82.6 | <u>77.8</u> | 82.5 | 82.4 | 82.4 | 81.3 | 81.9 |
| **VerbLE** | 81.9 | 82.1 | **82.9** | 82.2 | 82.5 | <u>77.8</u> | 82.4 | 82.2 | 82.3 | 81.2 | 81.8 |
| **MEINeg** | 81.6 | 81.5 | **83.4** | 81.8 | 82.9 | <u>77.0</u> | 82.5 | 82.6 | 82.1 | 81.2 | 81.7 |
| **NPI** | 81.9 | 82.2 | **83.2** | 81.8 | 82.3 | <u>76.9</u> | 81.8 | 82.3 | 82.3 | 80.8 | 81.5 |
| **SrcPP** | 81.9 | 81.9 | **83.0** | 80.9 | 82.3 | <u>77.2</u> | 82.3 | 82.6 | 81.9 | 81.1 | 81.5 |
| **CausCpl** | 81.5 | 81.1 | **83.1** | 80.8 | 82.6 | <u>77.0</u> | 82.9 | 82.9 | 81.6 | 81.3 | 81.5 |
| **Rel** | 81.4 | 81.8 | **83.4** | 81.5 | 82.3 | <u>77.3</u> | 81.6 | 82.4 | 82.0 | 80.9 | 81.5 |
| **ZHIQtf** | 81.6 | 81.6 | **82.8** | 81.6 | 82.6 | <u>76.5</u> | 82.6 | 82.3 | 81.9 | 81.0 | 81.4 |
| **RefPP** | 81.8 | 82.0 | **82.4** | 81.2 | 82.0 | <u>78.1</u> | 81.7 | 82.0 | 81.9 | 81.0 | 81.4 |
| **SpcPP** | 81.3 | 82.0 | **83.0** | 81.4 | 82.0 | <u>77.9</u> | 81.8 | 81.7 | 81.9 | 80.8 | 81.4 |
| **IndCLF** | 81.3 | 81.2 | **82.9** | 81.3 | 82.1 | <u>78.0</u> | 82.0 | 81.9 | 81.7 | 81.0 | 81.3 |
| **VP** | 81.3 | 81.6 | **82.9** | 81.0 | 82.2 | <u>77.2</u> | 81.9 | 82.2 | 81.7 | 80.9 | 81.3 |
| **AdvP** | 81.3 | 81.4 | **82.7** | 81.0 | 82.0 | <u>76.5</u> | 81.8 | 81.9 | 81.6 | 80.6 | 81.1 |
| **SHICop** | 80.8 | 81.1 | **82.6** | 80.5 | 81.9 | <u>77.1</u> | 81.9 | 81.6 | 81.2 | 80.6 | 80.9 |
| **Deixis** | 81.2 | 81.2 | **82.5** | 80.6 | 81.7 | <u>76.7</u> | 81.6 | 81.4 | 81.4 | 80.3 | 80.9 |
| **LocPP** | 80.8 | 82.0 | **82.9** | 81.1 | 81.6 | <u>76.2</u> | 81.2 | 81.2 | 81.7 | 80.0 | 80.9 |
| **Refl** | 80.8 | 81.0 | **82.5** | 80.9 | 81.6 | <u>76.4</u> | 81.6 | 81.9 | 81.3 | 80.4 | 80.8 |
| **PstVPP** | 81.2 | 80.7 | **82.0** | 80.6 | 81.4 | <u>76.6</u> | **82.0** | 81.9 | 81.1 | 80.5 | 80.8 |
| **NP** | 80.8 | 81.1 | **82.4** | 80.6 | 81.7 | <u>76.7</u> | 81.2 | 81.4 | 81.2 | 80.2 | 80.7 |
| **EvCLF** | 80.6 | 80.8 | **82.1** | 81.0 | 81.6 | <u>76.3</u> | 80.9 | 81.4 | 81.1 | 80.0 | 80.6 |
| **GoalPP** | 80.8 | 80.0 | **81.6** | 80.1 | 81.4 | <u>76.3</u> | **81.6** | 81.5 | 80.6 | 80.2 | 80.4 |
| **AdjP** | 80.2 | 80.6 | **81.7** | 80.3 | 81.2 | <u>75.4</u> | 81.0 | 81.2 | 80.7 | 79.7 | 80.2 |
| **Cmpr** | 80.3 | 80.4 | **81.6** | 80.2 | 81.4 | <u>75.6</u> | 80.9 | 81.0 | 80.6 | 79.7 | 80.2 |
| **PtcpPP** | 80.4 | 81.0 | **81.7** | 79.9 | 80.7 | <u>76.5</u> | 80.6 | 80.7 | 80.8 | 79.6 | 80.2 |
| **ConcCpl** | 80.3 | 80.5 | **81.7** | 80.1 | 81.1 | <u>75.1</u> | 81.4 | 80.9 | 80.7 | 79.6 | 80.1 |
| **SentIPP** | 80.5 | 80.5 | **81.3** | 79.9 | 80.6 | <u>75.4</u> | 80.6 | 80.3 | 80.6 | 79.2 | 79.9 |
| **DirPP** | 80.2 | 80.5 | 81.1 | 79.3 | **81.2** | <u>75.1</u> | 80.9 | 80.7 | 80.3 | 79.5 | 79.9 |
| **TopPP** | 80.1 | 80.2 | **81.2** | 79.9 | 80.6 | <u>74.9</u> | 80.6 | 80.5 | 80.3 | 79.2 | 79.8 |
| **Recp** | 80.0 | 80.0 | **81.3** | 79.8 | 80.6 | <u>74.8</u> | 80.3 | 80.5 | 80.3 | 79.0 | 79.7 |
| **LexNeg** | 79.3 | 79.5 | **80.8** | 80.0 | **80.8** | <u>76.4</u> | 79.9 | 80.4 | 79.9 | 79.4 | 79.6 |
| **ApprCLF** | 79.7 | 80.2 | **81.5** | 78.9 | 80.5 | <u>73.7</u> | 80.5 | 80.5 | 80.1 | 78.8 | 79.5 |
| **BUNeg** | 79.2 | 79.7 | **80.4** | 78.4 | 79.8 | <u>74.4</u> | 79.6 | 79.4 | 79.4 | 78.3 | 78.9 |
| **PreVPP** | 78.8 | 79.5 | **81.0** | 78.4 | 79.6 | <u>73.9</u> | 79.8 | 79.5 | 79.4 | 78.2 | 78.8 |
| **PathPP** | 78.4 | 79.1 | **80.0** | 78.3 | 79.5 | <u>73.7</u> | 78.7 | 79.7 | 79.0 | 77.9 | 78.4 |
| **KindCLF** | 77.8 | 78.2 | **78.9** | 77.3 | 78.3 | <u>73.9</u> | 78.4 | 78.5 | 78.0 | 77.3 | 77.7 |

Table 19: Performance in COMET scores of each system on different grammatical features.

|  | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **ZAIProg** | 92.9 | 92.5 | **93.7** | 92.2 | 93.6 | <u>88.8</u> | 93.5 | 93.1 | 92.8 | 92.2 | 92.5 |
| **MEINeg** | 92.1 | 91.7 | 93.9 | 92.1 | **94.0** | <u>88.2</u> | 93.3 | 93.5 | 92.5 | 92.2 | 92.3 |
| **NPI** | 92.0 | 92.4 | **93.0** | 91.6 | **93.0** | <u>88.3</u> | 92.6 | **93.0** | 92.2 | 91.7 | 92.0 |
| **ImpNeg** | 91.9 | 91.0 | **93.5** | 92.1 | 92.8 | <u>87.4</u> | 93.2 | 92.7 | 92.1 | 91.5 | 91.8 |
| **Recp** | 92.0 | 91.0 | **93.3** | 91.2 | 92.8 | <u>87.7</u> | 92.6 | 92.5 | 91.9 | 91.4 | 91.6 |
| **Pass** | 92.0 | 90.8 | **93.2** | 91.2 | 92.6 | <u>86.9</u> | 92.6 | 92.8 | 91.8 | 91.2 | 91.5 |
| **ZHIQtf** | 91.3 | 90.8 | **93.0** | 91.2 | 92.5 | <u>86.6</u> | 92.4 | 92.2 | 91.6 | 90.9 | 91.2 |
| **SHICop** | 91.6 | 90.6 | 92.2 | 90.1 | **92.4** | <u>87.0</u> | **92.4** | 92.0 | 91.1 | 91.0 | 91.0 |
| **Refl** | 91.1 | 90.0 | 92.7 | 90.3 | **92.8** | <u>85.7</u> | 92.5 | 92.7 | 91.0 | 90.9 | 91.0 |
| **GUOPrf** | 91.0 | 90.2 | 92.1 | 90.9 | 91.9 | <u>86.9</u> | **92.4** | 91.6 | 91.0 | 90.7 | 90.9 |
| **StdCLF** | 91.4 | 90.9 | **92.3** | 89.8 | **92.3** | <u>86.0</u> | **92.3** | 91.8 | 91.1 | 90.6 | 90.8 |
| **Deixis** | 91.3 | 90.7 | **92.7** | 89.4 | 92.0 | <u>86.3</u> | 92.1 | 91.8 | 91.0 | 90.5 | 90.8 |
| **SpanPP** | 91.1 | 90.3 | 92.0 | 90.1 | **92.6** | <u>86.2</u> | 92.2 | 91.5 | 90.9 | 90.6 | 90.7 |
| **VP** | 90.8 | 90.3 | **92.3** | 90.7 | 91.6 | <u>86.5</u> | 91.7 | 91.6 | 91.0 | 90.3 | 90.7 |
| **CausCpl** | 91.4 | 89.4 | 92.6 | 88.6 | 92.4 | <u>84.4</u> | 92.7 | **92.8** | 90.5 | 90.6 | 90.5 |
| **CondPP** | 91.0 | 91.1 | **92.8** | 89.0 | 91.7 | <u>86.1</u> | 91.3 | 91.2 | 91.0 | 90.1 | 90.5 |
| **SrcPP** | 90.4 | 90.7 | **92.4** | 90.0 | 91.8 | <u>84.9</u> | 91.9 | 91.9 | 90.9 | 90.1 | 90.5 |
| **GoalPP** | 91.1 | 90.6 | 91.8 | 89.1 | 91.4 | <u>86.1</u> | **92.0** | 91.6 | 90.7 | 90.3 | 90.5 |
| **VerbLE** | 91.0 | 90.0 | 91.6 | 90.4 | 91.4 | <u>85.4</u> | **91.8** | 91.5 | 90.8 | 90.0 | 90.4 |
| **PstVPP** | 91.0 | 90.5 | 91.6 | 89.5 | 91.3 | <u>85.6</u> | **91.7** | **91.7** | 90.7 | 90.1 | 90.4 |
| **AdvP** | 90.7 | 89.9 | **92.0** | 89.6 | 91.5 | <u>85.4</u> | 91.6 | 91.6 | 90.6 | 90.0 | 90.3 |
| **IndCLF** | 90.3 | 90.0 | **91.6** | 89.2 | 91.3 | <u>86.8</u> | 91.1 | 90.6 | 90.3 | 89.9 | 90.1 |
| **TopPP** | 90.4 | 90.0 | 91.5 | 89.8 | **91.6** | <u>85.1</u> | 91.3 | 91.0 | 90.4 | 89.8 | 90.1 |
| **RefPP** | 90.6 | 90.0 | **91.5** | 89.6 | 91.2 | <u>86.2</u> | 90.9 | 90.8 | 90.4 | 89.8 | 90.1 |
| **ConcCpl** | 90.5 | 89.6 | **92.1** | 89.5 | 91.6 | <u>83.6</u> | 91.9 | 91.7 | 90.4 | 89.7 | 90.1 |
| **BA** | 90.4 | 90.4 | **91.5** | 89.9 | 90.7 | <u>84.3</u> | 90.6 | 90.6 | 90.6 | 89.1 | 89.8 |
| **BUNeg** | 89.9 | 89.1 | 90.9 | 89.2 | **91.4** | <u>84.5</u> | 90.9 | 90.6 | 89.8 | 89.3 | 89.6 |
| **DirPP** | 89.9 | 89.7 | 90.8 | 87.5 | **90.9** | <u>85.7</u> | 90.8 | 90.6 | 89.5 | 89.5 | 89.5 |
| **SpcPP** | 89.4 | 89.6 | **91.0** | 89.3 | 90.2 | <u>85.6</u> | 90.2 | 90.4 | 89.8 | 89.1 | 89.5 |
| **LocPP** | 89.5 | 89.8 | **91.3** | 89.1 | 90.9 | <u>83.4</u> | 90.8 | 90.9 | 89.9 | 89.0 | 89.4 |
| **EvCLF** | 89.5 | 89.2 | **90.8** | 89.6 | 90.3 | <u>85.0</u> | 89.9 | **90.8** | 89.8 | 89.0 | 89.4 |
| **NP** | 89.7 | 89.3 | **90.9** | 88.6 | 90.6 | <u>84.8</u> | 90.4 | 90.2 | 89.6 | 89.0 | 89.3 |
| **Cmpr** | 89.4 | 89.1 | 90.9 | 88.8 | 90.9 | <u>83.7</u> | 90.6 | **91.0** | 89.5 | 89.1 | 89.3 |
| **PreVPP** | 89.8 | 89.6 | **91.4** | 88.5 | 89.9 | <u>84.6</u> | 90.2 | 89.7 | 89.8 | 88.6 | 89.2 |
| **SentIPP** | 90.1 | 89.5 | **90.6** | 88.6 | 90.5 | <u>83.7</u> | 90.1 | 89.9 | 89.7 | 88.5 | 89.1 |
| **AdjP** | 89.4 | 88.9 | **90.8** | 88.4 | 90.7 | <u>83.5</u> | 90.5 | 90.2 | 89.4 | 88.7 | 89.0 |
| **TmpSCpl** | 89.7 | 89.3 | **90.8** | 88.5 | 90.3 | <u>83.0</u> | 90.0 | 90.2 | 89.6 | 88.4 | 89.0 |
| **Rel** | 89.5 | 88.4 | **91.1** | 88.0 | 90.4 | <u>82.7</u> | 90.2 | 90.4 | 89.2 | 88.4 | 88.9 |
| **PtcpPP** | 89.1 | 89.0 | **90.7** | 88.3 | 90.4 | <u>83.6</u> | 89.9 | 89.6 | 89.3 | 88.4 | 88.8 |
| **LexNeg** | 88.7 | 88.4 | **90.2** | 88.4 | 90.0 | <u>84.1</u> | 89.7 | 90.0 | 88.9 | 88.5 | 88.7 |
| **ApprCLF** | 89.1 | 88.2 | 90.5 | 87.8 | **90.6** | <u>82.8</u> | 90.1 | 90.0 | 88.9 | 88.4 | 88.6 |
| **KindCLF** | 89.0 | 88.4 | **89.9** | 88.0 | 89.7 | <u>83.8</u> | 89.4 | 89.0 | 88.8 | 88.0 | 88.4 |
| **PathPP** | 84.9 | 84.5 | **86.4** | 83.6 | 85.9 | <u>79.4</u> | 85.1 | 85.9 | 84.8 | 84.1 | 84.4 |

Table 20: Performance in XCOMET scores of each system on different grammatical features.

| | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **ZAIProg** | 74.3 | 73.8 | 74.7 | 73.8 | 75.3 | <u>68.8</u> | **75.4** | 74.7 | 74.2 | 73.5 | 73.9 |
| **SpanPP** | 74.8 | 73.8 | 74.9 | 73.6 | 75.1 | <u>68.1</u> | **75.2** | 74.8 | 74.3 | 73.3 | 73.8 |
| **CausCpl** | 74.6 | 73.3 | 75.6 | 73.6 | 75.0 | <u>65.4</u> | 75.5 | **76.0** | 74.3 | 73.0 | 73.6 |
| **TmpSCpl** | 74.7 | 73.8 | 74.3 | 73.1 | 74.7 | <u>67.5</u> | 74.8 | **74.9** | 74.0 | 73.0 | 73.5 |
| **MEINeg** | 74.5 | 73.2 | 75.1 | 74.1 | 75.1 | <u>65.5</u> | **75.2** | 74.9 | 74.2 | 72.7 | 73.4 |
| **VerbLE** | 74.3 | 73.2 | 74.4 | 73.7 | 74.6 | <u>66.7</u> | **75.1** | 74.8 | 73.9 | 72.8 | 73.4 |
| **BA** | 74.2 | 73.6 | 73.9 | 73.0 | 74.5 | <u>65.8</u> | **75.0** | 74.4 | 73.7 | 72.4 | 73.0 |
| **SpcPP** | 73.7 | 73.2 | **74.3** | 73.1 | 74.0 | <u>66.8</u> | **74.3** | 73.1 | 73.6 | 72.1 | 72.8 |
| **Recp** | 73.4 | 72.6 | 74.2 | 73.1 | 73.7 | <u>66.4</u> | 74.5 | 74.1 | 73.3 | 72.2 | 72.8 |
| **NPI** | 73.8 | 73.3 | 73.3 | 72.7 | 74.0 | <u>66.6</u> | 74.5 | 73.8 | 73.3 | 72.2 | 72.7 |
| **RefPP** | 73.7 | 72.4 | 73.5 | 72.3 | 73.8 | <u>67.3</u> | 74.3 | 73.8 | 73.0 | 72.3 | 72.7 |
| **AdvP** | 73.6 | 72.4 | 73.9 | 72.6 | 74.0 | <u>66.2</u> | 74.4 | 74.1 | 73.1 | 72.2 | 72.6 |
| **ConcCpl** | 73.4 | 72.2 | 74.4 | 72.7 | 74.5 | <u>64.2</u> | 74.8 | 74.6 | 73.2 | 72.0 | 72.6 |
| **StdCLF** | 74.1 | 72.2 | 73.4 | 72.7 | 74.1 | <u>65.5</u> | 74.5 | 73.6 | 73.1 | 71.9 | 72.5 |
| **VP** | 73.4 | 72.5 | 73.8 | 72.4 | 73.8 | <u>65.8</u> | 74.0 | 73.6 | 73.0 | 71.8 | 72.4 |
| **Refl** | 73.2 | 71.6 | 73.6 | 71.4 | 73.5 | <u>65.7</u> | 73.6 | **74.1** | 72.5 | 71.7 | 72.1 |
| **AdjP** | 73.2 | 72.4 | 73.2 | 72.3 | 73.6 | <u>64.2</u> | 74.1 | 73.5 | 72.8 | 71.3 | 72.1 |
| **PtcpPP** | 73.4 | 71.6 | 73.3 | 72.0 | 73.5 | <u>66.0</u> | 73.7 | 72.9 | 72.6 | 71.5 | 72.0 |
| **SHICop** | 73.5 | 71.6 | 73.2 | 71.6 | 73.9 | <u>64.8</u> | 74.2 | 72.9 | 72.5 | 71.4 | 72.0 |
| **ZHIQtf** | 72.7 | 71.7 | 73.1 | 72.8 | 73.5 | <u>64.7</u> | 73.9 | 73.3 | 72.6 | 71.3 | 71.9 |
| **SrcPP** | 73.1 | 71.9 | 73.0 | 71.3 | 73.3 | <u>65.6</u> | 73.8 | 73.5 | 72.3 | 71.5 | 71.9 |
| **Rel** | 73.4 | 71.3 | 72.9 | 71.5 | 73.5 | <u>65.0</u> | 73.9 | 73.2 | 72.3 | 71.4 | 71.8 |
| **Cmpr** | 72.7 | 71.6 | 73.1 | 72.0 | 73.4 | <u>63.6</u> | 73.8 | **73.9** | 72.3 | 71.2 | 71.8 |
| **LexNeg** | 72.4 | 71.0 | 72.6 | 72.0 | 73.1 | <u>66.3</u> | 73.2 | 73.0 | 72.0 | 71.4 | 71.7 |
| **GUOPrf** | 72.8 | 71.6 | 73.0 | 70.9 | 73.1 | <u>66.1</u> | 73.3 | 72.8 | 72.1 | 71.3 | 71.7 |
| **SentIPP** | 72.8 | 72.2 | 73.1 | 71.3 | 73.2 | <u>63.7</u> | 73.5 | 73.0 | 72.3 | 70.8 | 71.6 |
| **Pass** | 73.1 | 71.0 | 73.0 | 72.0 | 72.9 | <u>64.8</u> | 73.2 | 72.9 | 72.3 | 70.9 | 71.6 |
| **EvCLF** | 72.5 | 71.7 | 72.9 | 72.5 | 72.4 | <u>64.1</u> | 72.8 | **73.0** | 72.4 | 70.6 | 71.5 |
| **LocPP** | 72.7 | 71.8 | 72.2 | 72.1 | 72.8 | <u>64.6</u> | **73.0** | 72.7 | 72.2 | 70.8 | 71.5 |
| **NP** | 72.6 | 71.4 | 72.5 | 71.3 | 72.7 | <u>64.9</u> | 73.2 | 72.7 | 72.0 | 70.9 | 71.4 |
| **CondPP** | **73.2** | 71.3 | 71.8 | 70.7 | **73.2** | <u>65.6</u> | 73.0 | 72.4 | 71.8 | 71.1 | 71.4 |
| **ImpNeg** | 72.7 | 70.8 | 72.8 | 72.2 | 72.7 | <u>64.4</u> | **73.5** | 71.9 | 72.1 | 70.6 | 71.4 |
| **GoalPP** | 72.5 | 71.9 | 73.1 | 70.6 | **73.2** | <u>64.1</u> | 72.5 | 72.6 | 72.0 | 70.6 | 71.3 |
| **IndCLF** | 72.1 | 71.6 | 72.3 | 71.0 | 72.3 | <u>65.6</u> | **73.0** | 72.7 | 71.8 | 70.9 | 71.3 |
| **PstVPP** | 72.6 | 71.6 | 72.6 | 70.6 | **73.3** | <u>64.0</u> | 72.8 | 72.9 | 71.8 | 70.8 | 71.3 |
| **DirPP** | 71.6 | 71.4 | 72.1 | 69.8 | 73.3 | <u>63.5</u> | **73.5** | **73.5** | 71.2 | 71.0 | 71.1 |
| **Deixis** | 71.9 | 70.7 | 71.9 | 70.3 | 72.5 | <u>65.3</u> | 72.7 | 72.5 | 71.2 | 70.8 | 71.0 |
| **BUNeg** | 71.3 | 71.0 | 72.7 | 71.2 | 72.6 | <u>63.3</u> | **73.0** | 72.0 | 71.5 | 70.2 | 70.9 |
| **PreVPP** | 71.9 | 71.0 | 72.0 | 70.2 | 72.0 | <u>63.7</u> | **72.1** | 71.6 | 71.3 | 69.8 | 70.6 |
| **TopPP** | 70.6 | 69.5 | 70.9 | 70.3 | 71.1 | <u>62.1</u> | **71.3** | 70.8 | 70.3 | 68.8 | 69.6 |
| **KindCLF** | 70.4 | 69.0 | 70.6 | 69.6 | 70.9 | <u>63.5</u> | **71.3** | 70.7 | 69.9 | 69.1 | 69.5 |
| **ApprCLF** | 70.1 | 69.9 | 69.7 | 68.9 | **70.2** | <u>60.6</u> | 70.1 | 70.1 | 69.7 | 67.8 | 68.7 |
| **PathPP** | 70.1 | 68.2 | 68.9 | 68.4 | 69.3 | <u>62.9</u> | 69.9 | **70.2** | 68.9 | 68.1 | 68.5 |

Table 21: Performance in COMETKIWI-QE scores of each system on different grammatical features.

| | Baidu | Niu | Google | DeepL | Ernie | Qwen | GPT | Claude | NMT-AVG | LLM-AVG | All-AVG |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ZAIProg | 95.3 | 94.5 | 94.9 | 94.4 | 95.6 | <u>91.6</u> | **95.7** | 95.3 | 94.8 | 94.5 | 94.7 |
| MEINeg | 94.9 | 94.2 | 95.3 | 94.0 | **95.5** | <u>90.8</u> | 95.4 | **95.5** | 94.6 | 94.3 | 94.4 |
| NPI | 95.1 | 94.6 | 94.6 | 93.5 | **95.2** | <u>91.3</u> | 95.2 | 95.0 | 94.4 | 94.2 | 94.3 |
| ZHIQtf | 95.0 | 94.3 | 94.6 | 94.2 | 95.0 | <u>90.7</u> | 95.1 | 94.5 | 94.5 | 93.8 | 94.2 |
| Recp | 94.5 | 93.9 | 94.9 | 94.0 | 94.5 | <u>91.7</u> | 95.2 | 94.5 | 94.3 | 94.0 | 94.1 |
| VP | 94.6 | 93.6 | 94.6 | 93.4 | 94.5 | <u>90.4</u> | **94.7** | 94.5 | 94.0 | 93.5 | 93.8 |
| ConcCpl | 94.6 | 93.4 | 94.7 | 93.7 | 94.3 | <u>88.3</u> | 95.2 | 94.5 | 94.1 | 93.1 | 93.6 |
| Refl | 94.4 | 92.8 | 94.5 | 92.7 | 94.1 | <u>90.0</u> | 94.5 | **94.6** | 93.6 | 93.3 | 93.5 |
| SHICop | 94.6 | 93.0 | 94.0 | 92.5 | 94.2 | <u>90.2</u> | 95.0 | 94.0 | 93.5 | 93.3 | 93.4 |
| AdvP | 94.2 | 92.9 | 94.0 | 92.5 | 93.9 | <u>89.6</u> | 94.6 | 94.1 | 93.4 | 93.1 | 93.2 |
| StdCLF | **94.4** | 93.2 | 93.5 | 92.7 | 93.9 | <u>89.6</u> | 94.3 | 93.7 | 93.5 | 92.9 | 93.2 |
| GUOPrf | **94.0** | 92.9 | 93.9 | 92.8 | **94.0** | <u>90.0</u> | 94.0 | 93.7 | 93.4 | 92.9 | 93.2 |
| Deixis | 94.1 | 92.9 | 93.8 | 92.0 | 94.0 | <u>89.8</u> | **94.5** | 94.0 | 93.2 | 93.1 | 93.1 |
| CausCpl | 94.3 | 92.0 | 94.4 | 91.5 | 94.0 | <u>88.0</u> | 94.8 | 94.8 | 93.1 | 92.9 | 93.0 |
| VerbLE | 93.5 | 92.7 | 93.7 | 92.9 | 93.3 | <u>89.5</u> | 94.3 | 93.7 | 93.2 | 92.7 | 93.0 |
| SentIPP | 93.8 | 92.8 | **94.2** | 92.1 | 93.8 | <u>89.2</u> | 93.8 | 93.3 | 93.2 | 92.5 | 92.9 |
| Pass | **93.8** | 92.0 | 93.7 | 92.4 | 93.4 | <u>90.1</u> | 93.3 | 93.7 | 93.0 | 92.6 | 92.8 |
| RefPP | 93.6 | 92.4 | 93.6 | 92.2 | 93.6 | <u>90.0</u> | 93.8 | 93.0 | 93.0 | 92.6 | 92.8 |
| BUNeg | 93.3 | 92.6 | **93.7** | 92.6 | 93.5 | <u>89.2</u> | 93.7 | 93.5 | 93.0 | 92.5 | 92.8 |
| ImpNeg | 93.9 | 92.2 | 93.1 | 92.9 | 93.5 | <u>89.0</u> | 94.1 | 93.2 | 93.0 | 92.5 | 92.7 |
| SpanPP | 93.4 | 92.1 | 93.7 | 92.0 | 93.7 | <u>89.7</u> | 93.8 | 93.2 | 92.8 | 92.6 | 92.7 |
| SpcPP | 93.2 | 92.8 | **93.4** | 92.0 | 93.1 | <u>89.2</u> | 93.4 | 93.2 | 92.8 | 92.2 | 92.5 |
| BA | 93.4 | 92.7 | 92.8 | 92.4 | 93.4 | <u>88.1</u> | 93.8 | 93.2 | 92.8 | 92.1 | 92.5 |
| EvCLF | **93.3** | 92.4 | 92.9 | 92.9 | 92.7 | <u>88.8</u> | 93.3 | 93.0 | 92.9 | 92.0 | 92.4 |
| PstVPP | **94.0** | 92.6 | 92.9 | 91.5 | 93.5 | <u>87.8</u> | 93.1 | 93.1 | 92.8 | 91.9 | 92.3 |
| IndCLF | 93.0 | 92.1 | 92.7 | 91.5 | 93.1 | <u>89.6</u> | 93.5 | 92.9 | 92.3 | 92.3 | 92.3 |
| AdjP | 93.3 | 91.7 | 93.2 | 91.6 | 93.2 | <u>88.2</u> | 93.8 | 93.0 | 92.4 | 92.0 | 92.3 |
| GoalPP | **93.9** | 92.6 | 93.1 | 90.9 | 93.4 | <u>87.7</u> | 93.0 | 92.8 | 92.6 | 91.7 | 92.2 |
| DirPP | 93.1 | 92.3 | 92.8 | 91.0 | 93.0 | <u>88.4</u> | 93.7 | 92.7 | 92.3 | 92.0 | 92.1 |
| TopPP | 92.7 | 92.1 | 92.9 | 92.0 | 92.8 | <u>89.0</u> | 93.0 | 92.5 | 92.4 | 91.8 | 92.1 |
| TmpSCpl | **93.3** | 92.1 | 92.7 | 91.1 | 92.9 | <u>86.9</u> | 93.3 | 93.0 | 92.3 | 91.5 | 91.9 |
| Cmpr | 92.6 | 91.3 | 92.9 | 91.2 | 93.2 | <u>87.5</u> | 93.4 | 93.1 | 92.0 | 91.8 | 91.9 |
| NP | 92.9 | 91.7 | 92.4 | 91.1 | 92.7 | <u>88.4</u> | 93.1 | 92.4 | 92.0 | 91.7 | 91.8 |
| CondPP | 93.1 | 91.5 | 92.4 | 90.4 | 92.8 | <u>88.5</u> | 93.3 | 92.2 | 91.8 | 91.7 | 91.8 |
| SrcPP | 92.2 | 91.4 | 92.3 | 91.2 | **92.5** | <u>88.1</u> | 92.5 | 92.0 | 91.8 | 91.3 | 91.5 |
| PtcpPP | 92.5 | 91.0 | 92.3 | 90.9 | 92.6 | <u>87.5</u> | **93.0** | 92.1 | 91.7 | 91.3 | 91.5 |
| LexNeg | 92.3 | 91.1 | 92.3 | 91.4 | 91.7 | <u>88.1</u> | 92.4 | 92.1 | 91.8 | 91.1 | 91.4 |
| Rel | 92.5 | 90.4 | 92.1 | 90.3 | 92.2 | <u>87.2</u> | 93.3 | 92.0 | 91.3 | 91.2 | 91.3 |
| KindCLF | 92.2 | 90.9 | 91.9 | 90.9 | **92.2** | <u>87.8</u> | 92.2 | 91.7 | 91.5 | 91.0 | 91.2 |
| PreVPP | 92.0 | 91.8 | **92.2** | 90.3 | 92.1 | <u>87.5</u> | 92.1 | 91.6 | 91.6 | 90.8 | 91.2 |
| LocPP | 91.9 | 91.4 | 91.5 | 91.4 | 91.8 | <u>86.6</u> | 92.4 | 92.1 | 91.6 | 90.7 | 91.1 |
| ApprCLF | 92.0 | 91.7 | 91.6 | 90.3 | 91.9 | <u>87.0</u> | 92.2 | 91.3 | 91.4 | 90.6 | 91.0 |
| PathPP | **89.1** | 86.4 | 87.2 | 86.6 | 87.5 | <u>84.5</u> | 88.3 | 88.8 | 87.3 | 87.3 | 87.3 |

Table 22: Performance in XCOMET-QE scores of each system on different grammatical features.

# Improving Statistical Significance in Human Evaluation of Automatic Metrics via Soft Pairwise Accuracy

**Brian Thompson**[†]
Amazon

**Nitika Mathur**
Oracle

**Daniel Deutsch**
Google

**Huda Khayrallah**
Microsoft

## Abstract

Selecting an automatic metric that best emulates human annotators is often non-trivial, because there is no clear definition of "best emulates." A meta-metric is required to compare the human judgments to the automatic metric scores, and metric rankings depend on the choice of meta-metric. We propose Soft Pairwise Accuracy (SPA), a new meta-metric that builds on Pairwise Accuracy (PA) but incorporates the statistical significance of both the human judgments and the metric scores. We show that SPA is more stable than PA with respect to changes in the number of systems/segments used for evaluation. We also show that PA can only assign a small set of distinct output values to metrics, and this results in many metrics being artificially assigned the exact same PA score. We demonstrate that SPA fixes this issue. Finally, we show that SPA is more discriminative than PA, producing more statistically significant comparisons between metrics. SPA was selected as the official system-level metric for the 2024 WMT Metrics Shared Task.

## 1 Introduction

Automatic metrics are crucial because researchers and practitioners in NLP typically can't afford the high cost and latency of high-quality human evaluations. Despite their shortcomings, metrics like word error rate and BLEU (Papineni et al., 2002)—in conjunction with carefully curated test sets—have been crucial for the field of NLP, as they have provided a yardstick to make continual progress over many decades in automatic speech recognition and machine translation (MT), respectively.

Reliance on automatic metrics makes selecting a good automatic metric of paramount importance. Conceptually, an automatic metric should emulate human judgments. Selecting an automatic metric

typically entails generating a set of human judgments for a wide variety of outputs from a large number of different systems, and selecting the automatic metric that produces scores most similar to the human judgments. But how do we quantify similarity? To select the metric which produces output most similar to human judgements, we need a meta-metric to compare metric scores and human judgments. Despite nearly two decades of research on MT meta-evaluation, the community has not reached a consensus on the choice of a meta-metric. Various meta-metrics have been introduced over the years to address problems with prior meta-metrics, while sometimes creating new problems or re-introducing old ones (see § 5).

Recent works (Mathur et al., 2020b; Kocmi et al., 2021) have argued that the primary application of a metric is to choose between two competing systems, therefore the best metric is the one which produces pairwise system rankings most similar to the pairwise system rankings produced by human judgements. This led to Pairwise Accuracy (PA) being adopted by the WMT Metrics shared task for the past several years (Freitag et al., 2021, 2022, 2023). However, this argument omits a key detail: standard best practice when comparing two systems with an automatic metric is to consider not only which system the metric prefers, but also *whether or not that preference is statistically significant* (Koehn, 2004). Thus we argue that metrics should emulate not only the accuracy of human pairwise ranking, but also the confidence or statistical significance of the human pairwise ranking.

To this end, we propose Soft Pairwise Accuracy (SPA), a new meta-metric which takes into account statistical significance of both the metric scores and the human judgments when evaluating the extent to which the metric in question agrees with the human judgments. We show that soft pairwise accuracy, as its name implies, can be viewed as a soft (i.e. non-binarized) version of PA, and present

---

Figure 1: Illustration of the individual components used to calculate both SPA and PA for the Prism metric (Thompson and Post, 2020a,b) on the WMT 2023 English-German language pair. Each box represents a comparison between two systems, systems $i$ and $j$. MT systems are sorted by average human judgment score for easier interpretation. The right column is one minus the absolute difference between the human preference for systems $i$ over system $j$ (left column) and the metric preference for system $i$ over system $j$ (middle column). In PA (top row), human and metric preferences are binarized to 0 and 1, and PA is thus an average of binary terms. In SPA (bottom row), human and metric preferences range from 0 to 1, and as a result SPA is an average of values ranging from 0 to 1. SPA can be viewed as a "soft" extension to pairwise accuracy that incorporates both human judgment and metric uncertainty, allowing for partial credit.

analysis that demonstrates SPA has several distinct advantages over PA. First, we find SPA is more stable with respect to the exact choice of MT systems and segments used. Second, we show that due to the binarization in its formulation, PA can only assign a small set of distinct output values to metrics, and in practice this results in many metrics being artificially assigned the exact same PA score. We demonstrate that SPA fixes this issue. Finally, we argue that PA is effectively equivalent to SPA with added noise due to binarization. We show that removing this noise (i.e. switching to SPA) results in substantially more statistically significant comparisons between metrics, making SPA a more discriminative and therefore more useful meta-metric. Our findings resulted in SPA being selected as the official system-level meta-metric for the 2024 WMT Metrics Shared Task (Freitag et al., 2024).

## 2 Method

We propose a simple meta-metric for evaluating automatic metrics given human judgments, which we denote Soft Pairwise Accuracy:

$$SPA = \binom{N}{2}^{-1} \sum_{i=0}^{N-1} \sum_{j=i+1}^{N-1} 1 - |p_{ij}^h - p_{ij}^m| \quad (1)$$

where $N$ is the number of systems for which we have human judgements and metric scores, $p_{ij}^h$ is the $p$-value for hypothesis that system $i$ is better than system $j$ given the human judgments, and $p_{ij}^m$ is the $p$-value for hypothesis that system $i$ is better than system $j$ given the metric scores. The term $\binom{N}{2}^{-1} = \frac{2}{N(N-1)}$ normalizes the summation by the total number of pairs of systems being compared.

For each pairwise system comparison, we use a permutation test (Fisher, 1935) to estimate statistical significance of the difference in the means of the segment-level scores from a particular metric (or the human judgements) for the two systems. We first randomly split the segment-level scores (ignoring the labels, i.e. which MT system produced each segment) into two parts and compute the difference in metric score mean. Repeating this process many times provides a set of mean differences we can reasonably expect under the null hypothesis that the

two systems are of the same quality. We compute a one-tailed $p$-value by calculating the fraction of the time that the random splits produce differences greater than or equal to the mean difference we observe for the two systems.

Permutation tests are appealing because they don't require any assumptions about the underlying distribution of the data. This fits our use case well because we cannot assume anything about the distribution of segment-level scores of a metric.[1] Permutation tests instead have the assumption of exchangeability (Pitman, 1937; Draper et al., 1993; Good, 2002)—that is, under the null hypothesis (in our case, that the two MT systems are of equal quality) the joint distribution of the observations is invariant under permutations of the data labels. To help ensure exchangeability, we perform permutations such that each split has exactly one translation of each test set sentence, commonly referred to as a paired permutation test (Good, 2013).

Here we present some concrete examples for the sake of intuition. Suppose a metric reports a $+10$ point difference between system $i$ and system $j$, and that the random permutations only produce a metric difference $\geq 10$ points $1\%$ of the time. Thus $p_{ij} = 0.01$ and we conclude that the metric has high confidence that system $i$ is better than system $j$. Likewise, if the metric reports the systems have a $-10$ point difference, we might find that the random permutations produce a metric difference $\geq -10$ points $99\%$ of the time. Thus $p_{ij} = 0.99$ and we conclude the metric has high confidence that system $i$ is worse than system $j$. If the systems have the same metric score, we would expect about half of the random permutations to produce a metric difference $\geq 0$ and thus $p_{ij} = 0.5$, indicating the metric finds the two systems indistinguishable from each other.

## 2.1 Relationship to Pairwise Accuracy

PA is defined as

$$PA = \binom{N}{2}^{-1} \sum_{i=0}^{N-1} \sum_{j=i+1}^{N-1} a_{ij}^m \qquad (2)$$

where $a_{ij}^m$ is 1 when the metric scores and human judgments prefer the same system and 0 otherwise. PA is equivalent to the Kendall rank correlation coefficient (Kendall, 1938), modulo a linear scaling and shifting (see § 5.1).

---

[1]Metric and human annotation distributions are both highly variable (Lo et al., 2023b; Knowles and Lo, 2024).

A $p$-value $p_{ij}$ will be less than 0.5 when the human raters (or automatic metric) prefer system $i$ over system $j$, and greater than 0.5 when the human raters (or automatic metric) prefer system $j$ over system $i$. This allows us to define PA in terms of binarized $p$-values:

$$PA = \binom{N}{2}^{-1} \sum_{i=0}^{N-1} \sum_{j=i+1}^{N-1} 1 - \left| \lfloor p_{ij}^h \rceil - \lfloor p_{ij}^m \rceil \right|$$
(3)

Where binarization is denoted as:

$$\lfloor x \rceil = \begin{cases} 1 & x \geq 0.5 \\ 0 & x < 0.5 \end{cases}$$

Comparing Equation 1 and Equation 3 illustrates that SPA can be viewed as a 'soft' extension to PA that incorporates uncertainty in both the human and metric scores. A visualization of this is provided in Figure 1.

In cases where both the MT metric and the human evaluation both have high statistical significance (regardless of whether the metric agrees with the human judgments or not, i.e. $p_{ij}^m \approx 0$ or $p_{ij}^m \approx 1$), the contribution of that system pair to SPA and PA is approximately identical. However, there are two important cases where our meta-metric differs from PA:

1. The human evaluation has high statistical significance (i.e. $p_{ij}^h \approx 0$ or $p_{ij}^h \approx 1$), but the metric has low statistical significance (i.e. $p_{ij}^m \approx 0.5$): Even if the metric happens to choose the correct winner, we partially penalize the metric for not having high statistical significance.

2. The human evaluation finds the systems are approximately tied (i.e. $p_{ij}^h \approx 0.5$): In this case, we partially penalize the metric if has high statistical significance (i.e. $p_{ij}^m \approx 0$ or $p_{ij}^m \approx 1$) even if it happens to pick the same winner as the human evaluation, and to get full credit the metric must match the human evaluation statistical significance (i.e. $p_{ij}^m \approx p_{ij}^m \approx 0.5$)

## 2.2 Addressing Metric Ties in PA

The fact that PA considers only binary wins/losses (i.e. the binarization in Equation 3) results in an interesting shortcoming in PA. There are $\binom{N}{2}$ pairs of $N$ systems, and thus $\binom{N}{2} + 1$ distinct values that

PA can take on $(0/\binom{N}{2}, 1/\binom{N}{2}, ..., \binom{N}{2}/\binom{N}{2})$. For example, in WMT 2022 En-De, there are $N = 14$ MT systems and thus $\binom{N}{2} + 1 = 92$.

However, metrics tend to perform better than a random baseline, so only the upper half of the range is actually useful (e.g. this leaves 46 distinct values for $N = 14$ systems). We find that this results in PA reporting the same scores for several sets of metrics (see § 4.3). By removing this binarization, SPA has no such issues.

## 3 Experimental Setup

### 3.1 Data

We conduct experiments on the data from the 2022 and 2023 WMT Metrics Shared Tasks (Freitag et al., 2022, 2023). In particular, we use the primary language pairs where MQM judgments were collected. We use the MT Metrics Eval V2 toolkit[2] to retrieve official shared task scores.

We make the somewhat arbitrary decision to compare all metrics, including non-primary metrics but excluding QE metrics (i.e. reference-free metrics) which provide segment-level scores.

In order to compute the statistical significance of comparisons between metrics, we make the simplifying assumption that all system-level metrics are the average of their segment-level metric. This is not true for some metrics, including BLEU (Papineni et al., 2002) and chrF (Popović, 2015). While it would be possible to re-compute BLEU and chrF for each subset, we average the sentence-level versions of these metrics for simplicity. To the best of our knowledge, this approach is also taken in recent WMT metrics shared tasks.

### 3.2 $p$-value Speed Optimization

We estimate each $p$-value from 1000 random permutations.[3] A naive implementation of the paired permutation test is not computationally prohibitive when computing $p$-values for all systems/metrics a single time, but it becomes problematic when we want to compute these values many times in order to estimate statistical significance of metric comparisons.

Experimentally, we find the main speed bottleneck to be generating the random permutations, so when estimating statistical significance of metric

comparisons we cache a batch of permutations and use it for each pair of systems, on a per test-set basis. Additionally, by sharing permutations across system pairs, this allows us to pre-compute the contribution of each system to means of the random permutations, allowing computations to be linear instead of square in the number of systems. See our code[4] for full implementation details. This results in a speedup of over 1000x compared to the implementation in Scipy (Virtanen et al., 2020).[5] Our speed optimization does not change the computation of a $p$-value for a single system-level comparison, but it does mean that the $p$-value for one pair of systems is no longer computed independent from the $p$-value for any other pair of systems. Given that we are using these $p$-values as an approximate level of confidence for the system-level comparisons in the SPA meta-metric formulation, as opposed to making any claims about the actual statistical significance of the system-level comparisons, we believe this lack of independence should be inconsequential.

## 4 Analysis

Meta-metric evaluation is challenging because there is no ground truth (i.e., we don't know the true ranking of the metrics). Instead, we conduct analysis to compare SPA and PA. First, we study how sensitive the meta-metric results are when ablating the number of MT systems and number of segments per MT system, with the assumption that lower sensitivity to the exact systems/segments used indicates a better meta-metric. Second, we examine whether PA indeed has the problem of ties that we hypothesized in § 2.2, and whether SPA fixes this issue. Finally, we test our hypothesis that the binarization in PA is effectively acting as additive random noise, and that SPA is effectively the same underlying meta-metric with the noise term removed.

### 4.1 Ablation: Number of Systems

Each year, WMT and the associated metrics task collect and score many online and submitted MT systems. For an ideal meta-metric, the exact choice of MT systems would have minimal impact on the

---

Figure 2: Final metric ranking stability when ablating the number of MT systems (and thus the number of total MQM judgments), measured as change in Pearson correlation coefficient (Pearson $r$) from the ranking computed on all MT systems. Values are averaged over 1000 random trials. We find SPA to be more stable than PA in all cases.

metric rankings. We perform an ablation on the number of MT systems being scored, keeping the number of annotations per system fixed. We then compute the correlation (as measured by Pearson's $r$) between the meta-metric's ranking of the ablations compared to that same meta-metric's full ranking. This allows us to evaluate how sensitive the metric is to the exact selection of MT systems.

When ablating the number of MT systems (and keeping the number of annotations per system fixed), we find (see Figure 2) that SPA is more stable than PA across all MQM language pairs in the last two years of WMT Metrics Shared Tasks.

### 4.2 Ablation: Sample Size

Since SPA relies on the pairwise $p$-values between MT systems, it is also natural to ask how SPA behaves when the number of available segments used for evaluating systems is small since it is harder to find statistical differences between systems with a smaller sample size. To answer this question, we calculate 95% confidence intervals for both PA and SPA values of two highly performant metrics—in particular, we considered xCOMET (Guerreiro et al., 2023) and MetricX-23 (Juraska et al., 2023)—on WMT 2023 using bootstrapping for various numbers of segments, thereby simulating scenarios with less human annotations but a fixed number of

MT systems.

When ablating the number of segments per MT system (and keeping the number of MT systems fixed), we find (see Figure 3) that SPA has tighter 95% confidence intervals than PA (shown on Metric-X and xCOMET), and that the confidence interval converges to its final value with smaller sample sizes than PA.

### 4.3 Ties

As discussed in § 2.2, the binarization in PA limits the number of distinct values it can assign to metrics. to $\binom{N}{2} + 1$. In practice, we find it tends to take on far fewer values. For example for WMT 2022 En→De, PA could theoretically take on 92 distinct values, but because the metrics fall in a fairly narrow range (PA is $0.626$ for the worst metric and $0.813$ for the best), the 21 metrics have only 11 distinct PA scores, with one 5-way PA tie and several 2- and 3-way PA ties (see Figure 4). Since SPA does not binarize each system comparison, it is able to assign any value to each metric, and is therefore potentially better able to distinguish between metrics.

Results for all language pairs are in Table 1. We find that on average, PA produces about half as many distinct values as there are metrics while SPA produces one unique value per unique metric.

Figure 3: The 95% confidence intervals for SPA (blue) and PA (red) on Metric-X (top) and XCOMET (bottom) when varying the number of annotations per system. We find that SPA has a tighter confidence interval, and that the confidence interval shrinks to its full value with smaller sample sizes than PA.

## 4.4 Statistical Significance of Metric Comparisons

We hypothesize that the binarization in PA is essentially acting as additive random noise on top of the underlying SPA meta-metric. If this is true (and the magnitude of the noise does not dominate the underlying signal), we would expect SPA to produce a similar metric ranking to PA, but with increased statistical significance. To test this, we compute statistical significance of the comparisons between each metric using the PERM-INPUTS (Deutsch et al., 2021) method. We follow recent shared tasks in greedily computing significance clusters, by starting with the highest scoring metric and assigning rank 1 to all metrics until we encounter the first metric that is statistically significantly different from *any* previous metric so far. That metric is assigned rank 2, and the process repeats until all metrics have been assigned a rank. We echo the shared task organizers' warning that this method can place two metrics that are statistically indistinguishable in different significance clusters (and in the case of PA, we observe this multiple times).

On average, SPA increases the number of sta-

tistically significant pairwise comparisons by 31% and the number of significance clusters by 40% compared to PA, while producing similar scores for each metric (see Figure 4 for a visualization for WMT 2022 En→De results and Table 1 for results summary). This is consistent with our hypothesis that PA is effectively SPA with added noise due to binarization. This means that SPA is a more discriminative, and therefore more useful, meta-metric than PA.

## 5 Historical Context and Related Work

WMT has run a machine translation evaluation since 2006 (Koehn and Monz, 2006). Since 2007 (Callison-Burch et al., 2007), there has also been meta-evaluation of automatic metrics on the submitted translation systems. Here we summarize the rich 17 year history of system-level meta-evaluation at the WMT Metrics Shared Tasks[6] and work related to and directly impacting the shared tasks, in order to demonstrate how our work fits into the historical context.

---

[6] The WMT Shared Tasks have typically evaluated at both the system- and segment-level, but we focus on system-level meta-evaluation as it is most relevant to our work.

Figure 4: Metric Comparison Significance, WMT 2022 En→De. Note that PA only assigns 11 distinct values to the 21 metrics (ties are shown in alternating Purple and Yellow text), whereas SPA produces a distinct value for each of the 21 metrics. SPA produces more statistically significant (*p*-value <= 0.05, shown in green) comparisons between metrics (163 vs 108). As a result, SPA divides the metrics into 8 significance clusters (delineated with blue lines) compared to only 5 for PA. Results for other language pairs (not shown) are similar.

1228

| Testset | Language Pairs | # MT Systems | # MT Metrics | Distinct Metric Values (↑) | | | Significant Comparisons (↑) | | | Significant Clusters (↑) | | |
|---------|---------------|--------------|--------------|-----|------|-----|-----|------|-----|-----|------|-----|
| | | | | PA | SPA | Max | PA | SPA | Max | PA | SPA | Max |
| wmt22 | En→De | 14 | 21 | 11 | **21** | 21 | 108 | **163** | 210 | 5 | **8** | 21 |
| wmt22 | Zh→En | 15 | 21 | 12 | **21** | 21 | 150 | **177** | 210 | 6 | **9** | 21 |
| wmt22 | En→Ru | 15 | 20 | 10 | **20** | 20 | 88 | **133** | 190 | 4 | **6** | 20 |
| wmt23 | En→De | 12 | 25$^{\dagger}$ | 12$^{\dagger}$ | **24**$^{\dagger}$ | 24$^{\dagger}$ | 171 | **206** | 276$^{\dagger}$ | 5 | **6** | 24$^{\dagger}$ |
| wmt23 | He→En | 13 | 25 | 11 | **25** | 25 | 180 | **224** | 300 | 5 | **8** | 25 |
| wmt23 | Zh→En | 14 | 25 | 12 | **25** | 25 | 186 | **229** | 300 | **7** | 7 | 25 |

Table 1: Number of distinct values produced, number of statistically significant pairwise comparisons (p-value $<= 0.05$), and number of statistical significance clusters for PA and SPA. We provide the best possible value for each category (Max) for comparison, but note that even an ideal meta-metric would likely not achieve this value due to some metrics being highly correlated with each other (e.g. due to training on the same data). $^{\dagger}$: InstructScore and SEScoreX scores as returned by MT Metrics Eval v2 for WMT23 En-De are identical, causing an exact tie in both PA and SPA. We believe this is an error in MT Metrics Eval v2 but for posterity keep them as-is.

In the WMT 2007-2013 metrics evaluations (Callison-Burch et al., 2007, 2008, 2009, 2010, 2011, 2012; Macháček and Bojar, 2013) Spearman's rank correlation coefficient $\rho$ was used for meta-evaluation of metrics. This was motivated by the fact that Spearman's makes fewer assumptions about the data than the Pearson correlation coefficient.

The WMT 2013 Translation Shared Task (Bojar et al., 2013) introduced system clusters (groups of systems that cannot be distinguished given the human judgments), and the 2013 metrics task (Macháček and Bojar, 2013) introduced empirical confidence of Spearman's $\rho$ using bootstrap resampling. Since they were not able to resample on the submitted metrics, they only re-sampled human judgments. This iteration also discussed the fact that Spearman's $\rho$ does not give partial credit. The penalty is equal for all wrong judgments, regardless of if the systems are close or far in quality. To compensate they present additional methods of analysis: Pearson's, and correlation with systems' clusters from the translation task (Bojar et al., 2013). Those clusters were treated as 'ranks with ties,' and then correlation computed against Pearson's and Pearson's correlation against 'fuzzy ranks' (the average over ranks of all systems that are not significantly different in human quality).

In 2014, the metrics task (Macháček and Bojar, 2014) fully switched to Pearson's $r$ from Spearman's $\rho$. They also did bootstrap resampling to get empirical confidence intervals of system level correlations. This change to Pearson's was due to the concerns pointed out in the previous year's shared task, which had explored other meta-metrics.

The 2015 metrics task (Stanojević et al., 2015) continued with Pearson's $r$, and also presented analysis of Pearson's $r$ vs Spearman's $\rho$, and highlighted the instability of Spearman's $\rho$ when MT systems are similar.

The 2016 metrics task (Bojar et al., 2016) stuck with Pearson's $r$, but changed the confidence to be the Williams test (Williams, 1959), as Graham and Baldwin (2014) had noted that this test is appropriate for dependent correlations.

The 2017 metrics task (Bojar et al., 2017) kept Pearson's $r$, and Williams test. They also added a pairwise significance test using Williams test. This continued in 2018 and 2019 (Ma et al., 2018, 2019)

The 2020 metrics task (Mathur et al., 2020b) continued to use Pearson's, but also includes Kendall's Tau for analysis. Kendall's Tau is a closer match for the system ranking use case, since it is evaluating whether the ordering of a pair of systems is the same as the human ordering. However, it does not take into account the magnitude difference.

In 2021, the metrics task (Freitag et al., 2021) adopted pairwise accuracy (Kocmi et al., 2021), motivated in part by the fact that MT system outliers had an outsized impact on Pearson correlation when it is used to rank MT Metrics (Mathur et al., 2020a). Pairwise accuracy produces the same system-level ranking as Kendall's Tau, as they are equivalent modulo a linear scaling and shifting (see § 5.1). The PERM-BOTH hypothesis test of Deutsch et al. (2021) was used to determine significance. 2021 and 2022 (Freitag et al., 2022, 2023) follow.

In summary, the historical context of the WMT metric evaluations demonstrates that meta-evaluation is very challenging due to the numerous issues that must be simultaneously addressed, and underscores the pitfalls of making changes to meta-evaluation without considering the full set of ramifications. Most relevant to our work, it appears that the switch to pairwise accuracy in 2021 reduced the influence of outliers (Mathur et al., 2020a) and (somewhat) aligned meta-evaluation with the standard use of comparing two systems with a metric, but it also reintroduced a problem that was first pointed out by Macháček and Bojar (2013) and more fully addressed by the change to Pearson's $r$ from Spearman $\rho$ by Macháček and Bojar (2014): a disregard for the magnitude of differences. We address this issue by considering empirical confidence, which was first added by Macháček and Bojar (2013), and in the process we also better align meta-evaluation to the (more correct) usage of comparing two systems with a metric while also considering the statistical significance of the results.

## 5.1 Relationship to Kendall's Tau

Our work builds on pairwise accuracy, typically attributed to Kocmi et al. (2021). Pairwise accuracy is equivalent to the the widely used Kendall rank correlation coefficient (Kendall, 1938), modulo a linear scaling and shifting. Kocmi et al. (2021) present pairwise accuracy as simply "accuracy" and make no mention of its relation to Kendall, which was already in use for MT meta-evaluation (Mathur et al., 2020b). The term "pairwise accuracy" appears to have been coined by Freitag et al. (2021) to distinguish it from other types of accuracy.

Kendall's Tau is defined in terms of concordance (equivalent to our previously defined $a_{ij}^m$) and discordance $d_{ij}^m$, defined to be 1 when the metric and human judgments disagree and 0 otherwise:

$$\tau = \binom{N}{2}^{-1} \sum_{i=0}^{N-1} \sum_{j=i+1}^{N-1} (a_{ij}^m - d_{ij}^m) \quad (4)$$

Any system pair which is not concordant is discordant,[7] and thus $d_{ij}^m = 1 - a_{ij}^m$. Given this and the

---

[7]We ignore tie handling, as ties are extremely unlikely in system-level evaluation. Ties in *segment-level* evaluation are an entirely different matter (Deutsch et al., 2023a).

definition of PA from Equation 2, we have:

$$\tau = \binom{N}{2}^{-1} \sum_{i=0}^{N-1} \sum_{j=i+1}^{N-1} a_{ij}^m - (1 - a_{ij}^m)$$
$$= 2\left(\binom{N}{2}^{-1} \sum_{i=0}^{N-1} \sum_{j=i+1}^{N-1} a_{ij}^m\right) - 1 \quad (5)$$
$$= 2\,PA - 1$$

## 5.2 Additional Connections to Prior Work

Graham and Liu (2016) proposed a method of sampling translations from every pair of competing MT systems, creating synthetic systems for scoring. Our work has clear similarities in that we create and score synthetic permutations, but differs in how those synthetic systems are used in the meta-metric formulation.

Mathur et al. (2020a) showed that MT system outliers had an outsized impact on Pearson correlation. In SPA, outliers impact is limited because $p$-values saturate at 0 or 1.

Knowles (2021) highlights that as WMT annotation protocols have shifted the original statistical assumptions, and questions the validity of the resulting protocols. Similarly, we show that shifts over the years have caused problems in meta-evaluation.

Lo et al. (2023a) investigated what magnitude of metric changes tend to be statistically significant. SPA uses statistical significance measures ($p$-values) directly, as opposed to the magnitude of metric differences (e.g. as in Pearson correlation).

Deutsch et al. (2023a) demonstrated that principled tie handling is crucial when comparing MT metrics at the segment level, because some metrics produce quantized scores that often result in ties. SPA is system level (i.e. sentence level scores averaged over the entire test set), so exact ties are very unlikely. However, SPA can be seen as giving full credit for (statistical) ties, which is similar in spirit.

We show that quantization (specifically binarization) is problematic in PA. Quantization in evaluation has proved problematic in other spaces as well—for example, Schaeffer et al. (2024) attributes the widely repeated claim that LLMs have emergent properties to quantization in evaluation.

## 6 Conclusions

We introduce a new meta-metric which we denote soft pairwise accuracy, and show that it improves on pairwise accuracy in a number of ways, most notably that it is more stable than pairwise accuracy

when ablating the number of systems and annotations per system, it fixes an issue of metric ties observed in pairwise accuracy, and it produces more statistically significant comparisons between metrics than pairwise accuracy. We also discuss how soft pairwise accuracy fits into and builds upon the nearly two decade history of meta-evaluation at the WMT Metric Shared Tasks.

## Acknowledgments

## Limitations

When computing $p$-values, we assume that system-level metric scores are the average of segment-level metrics scores. There is a line of recent work that seeks to incorporate contextual information into automatic metrics. Many such works still produce scores at the segment level (e.g. Vernikos et al., 2022; Hu et al., 2023; Agrawal et al., 2024) but others produce one score per window of a few sentences (Raunak et al., 2024) or one score per paragraph (Deutsch et al., 2023b). Our method should still be applicable in such cases, but would require permuting windows or paragraphs instead of segments. Additionally, as previously noted, some metrics—notably BLEU (Papineni et al., 2002) and chrF (Popović, 2015)—compute statistics at the segment level and combine them to create document-level scores. Again, permutations would still work but would require some modification. To the best of our knowledge, this issue is not limited to our work—the same assumption is made in prior work computing statistical significance of metrics, including the WMT shared tasks (Freitag et al., 2021, 2022, 2023) and Deutsch et al. (2021).

It is worth noting that the permutations in this work (as in prior works) are done on a single test set, and do not necessarily reflect variations in performance that could result from using the metrics in another domain. Prior work has shown that trained metrics are sensitive to a shift in domain relative to the data domain they were trained on (Zouhar et al., 2024).

# References

Sweta Agrawal, Amin Farajian, Patrick Fernandes, Ricardo Rei, and André F. T. Martins. 2024. Is context helpful for chat translation evaluation? *Preprint*, arXiv:2403.08314.

Ondřej Bojar, Christian Buck, Chris Callison-Burch, Christian Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, and Lucia Specia. 2013. Findings of the 2013 Workshop on Statistical Machine Translation. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 1–44, Sofia, Bulgaria. Association for Computational Linguistics.

Ondřej Bojar, Yvette Graham, and Amir Kamran. 2017. Results of the WMT17 metrics shared task. In *Proceedings of the Second Conference on Machine Translation*, pages 489–513, Copenhagen, Denmark. Association for Computational Linguistics.

Ondřej Bojar, Yvette Graham, Amir Kamran, and Miloš Stanojević. 2016. Results of the WMT16 metrics shared task. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 199–231, Berlin, Germany. Association for Computational Linguistics.

Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2007. (meta-) evaluation of machine translation. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 136–158, Prague, Czech Republic. Association for Computational Linguistics.

Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2008. Further meta-evaluation of machine translation. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 70–106, Columbus, Ohio. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, Kay Peterson, Mark Przybocki, and Omar Zaidan. 2010. Findings of the 2010 joint workshop on statistical machine translation and metrics for machine translation. In *Proceedings of the Joint Fifth Workshop on Statistical Machine Translation and MetricsMATR*, pages 17–53, Uppsala, Sweden. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, and Lucia Specia. 2012. Findings of the 2012 workshop on statistical machine translation. In *Proceedings of the Seventh Workshop on Statistical Machine Translation*, pages 10–51, Montréal, Canada. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, and Josh Schroeder. 2009. Findings of the 2009 Workshop on Statistical Machine Translation. In *Proceedings of the Fourth Workshop on Statistical Machine Translation*, pages 1–28, Athens, Greece. Association for Computational Linguistics.

Chris Callison-Burch, Philipp Koehn, Christof Monz, and Omar Zaidan. 2011. Findings of the 2011 workshop on statistical machine translation. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 22–64, Edinburgh, Scotland. Association for Computational Linguistics.

Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. A statistical analysis of summarization evaluation metrics using resampling methods. *Transactions of the Association for Computational Linguistics*, 9:1132–1146.

Daniel Deutsch, George Foster, and Markus Freitag. 2023a. Ties matter: Meta-evaluating modern metrics with pairwise accuracy and tie calibration. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12914–12929, Singapore. Association for Computational Linguistics.

Daniel Deutsch, Juraj Juraska, Mara Finkelstein, and Markus Freitag. 2023b. Training and meta-evaluating machine translation evaluation metrics at the paragraph level. In *Proceedings of the Eighth Conference on Machine Translation*, pages 996–1013, Singapore. Association for Computational Linguistics.

David Draper, James S Hodges, Colin L Mallows, and Daryl Pregibon. 1993. Exchangeability and data analysis. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 156(1):9–28.

RA Fisher. 1935. *The design of experiments.* Oliver & Boyd.

Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are LLMs breaking MT metrics? Results of the WMT24 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, Miami, USA. Association for Computational Linguistics.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774, Online. Association for Computational Linguistics.

Phillip Good. 2002. Extensions of the concept of exchangeability and their applications. *Journal of Modern Applied Statistical Methods*, 1:243–247.

Phillip Good. 2013. *Permutation tests: a practical guide to resampling methods for testing hypotheses*. Springer Science & Business Media.

Yvette Graham and Timothy Baldwin. 2014. Testing for significance of increased correlation with human judgment. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 172–176, Doha, Qatar. Association for Computational Linguistics.

Yvette Graham and Qun Liu. 2016. Achieving accurate conclusions in evaluation of automatic machine translation metrics. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1–10, San Diego, California. Association for Computational Linguistics.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. XCOMET: Transparent Machine Translation Evaluation through Fine-grained Error Detection. *Preprint*, arXiv:2310.10482.

Xinyu Hu, Xunjian Yin, and Xiaojun Wan. 2023. Exploring context-aware evaluation metrics for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15291–15298, Singapore. Association for Computational Linguistics.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

M. G. Kendall. 1938. A NEW MEASURE OF RANK CORRELATION. *Biometrika*, 30(1-2):81–93.

Rebecca Knowles. 2021. On the stability of system rankings at WMT. In *Proceedings of the Sixth Conference on Machine Translation*, pages 464–477, Online. Association for Computational Linguistics.

Rebecca Knowles and Chi-kiu Lo. 2024. Calibration and context in human evaluation of machine translation. *Natural Language Processing*, page 1–25.

Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 478–494, Online. Association for Computational Linguistics.

Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.

Philipp Koehn and Christof Monz. 2006. Manual and automatic evaluation of machine translation between European languages. In *Proceedings on the Workshop on Statistical Machine Translation*, pages 102–121, New York City. Association for Computational Linguistics.

Chi-kiu Lo, Rebecca Knowles, and Cyril Goutte. 2023a. Beyond correlation: Making sense of the score differences of new MT evaluation metrics. In *Proceedings of Machine Translation Summit XIX, Vol. 1: Research Track*, pages 186–199, Macau SAR, China. Asia-Pacific Association for Machine Translation.

Chi-kiu Lo, Samuel Larkin, and Rebecca Knowles. 2023b. Metric score landscape challenge (MSLC23): Understanding metrics' performance on a wider landscape of translation quality. In *Proceedings of the Eighth Conference on Machine Translation*, pages 776–799, Singapore. Association for Computational Linguistics.

Qingsong Ma, Ondřej Bojar, and Yvette Graham. 2018. Results of the WMT18 metrics shared task: Both characters and embeddings achieve good performance. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 671–688, Belgium, Brussels. Association for Computational Linguistics.

Qingsong Ma, Johnny Wei, Ondřej Bojar, and Yvette Graham. 2019. Results of the WMT19 metrics shared task: Segment-level and strong MT systems pose big challenges. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 62–90, Florence, Italy. Association for Computational Linguistics.

Matouš Macháček and Ondřej Bojar. 2013. Results of the WMT13 metrics shared task. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 45–51, Sofia, Bulgaria. Association for Computational Linguistics.

Matouš Macháček and Ondřej Bojar. 2014. Results of the WMT14 metrics shared task. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 293–301, Baltimore, Maryland, USA. Association for Computational Linguistics.

Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020a. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4984–4997, Online. Association for Computational Linguistics.

Nitika Mathur, Johnny Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020b. Results of the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 688–725, Online. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Edwin James George Pitman. 1937. Significance tests which may be applied to samples from any populations. ii. the correlation coefficient test. *Supplement to the Journal of the Royal Statistical Society*, 4(2):225–232.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Vikas Raunak, Tom Kocmi, and Matt Post. 2024. SLIDE: Reference-free evaluation for machine translation using a sliding document window. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 205–211, Mexico City, Mexico. Association for Computational Linguistics.

Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. 2024. Are emergent abilities of large language models a mirage? In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA. Curran Associates Inc.

Miloš Stanojević, Amir Kamran, Philipp Koehn, and Ondřej Bojar. 2015. Results of the WMT15 metrics shared task. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 256–273, Lisbon, Portugal. Association for Computational Linguistics.

Brian Thompson and Matt Post. 2020a. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 90–121, Online. Association for Computational Linguistics.

Brian Thompson and Matt Post. 2020b. Paraphrase generation as zero-shot multilingual translation: Disentangling semantic similarity from lexical and syntactic diversity. In *Proceedings of the Fifth Conference on Machine Translation*, pages 561–570, Online. Association for Computational Linguistics.

Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 118–128, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. 2020. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272.

Evan James Williams. 1959. *Regression Analysis*, volume 14 of *WILEY SERIES in PROBABILITY and STATISTICS: APPLIED PROBABILITY and STATISTICS SECTION Series*. Wiley.

Vilém Zouhar, Shuoyang Ding, Anna Currey, Tatyana Badeka, Jenyuan Wang, and Brian Thompson. 2024. Fine-tuned machine translation metrics struggle in unseen domains. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 488–500, Bangkok, Thailand. Association for Computational Linguistics.

# Speech is More Than Words:
# Do Speech-to-Text Translation Systems Leverage Prosody?

**Ioannis Tsiamas**◇*    **Matthias Sperber**†    **Andrew Finch**†    **Sarthak Garg**†

◇Universitat Politècnica de Catalunya    †Apple

ioannis.tsiamas@upc.edu, sperber@apple.com

## Abstract

The prosody of a spoken utterance, including features like stress, intonation and rhythm, can significantly affect the underlying semantics, and as a consequence can also affect its textual translation. Nevertheless, prosody is rarely studied within the context of speech-to-text translation (S2TT) systems. In particular, end-to-end (E2E) systems have been proposed as well-suited for prosody-aware translation because they have direct access to the speech signal when making translation decisions, but the understanding of whether this is successful in practice is still limited. A main challenge is the difficulty of evaluating prosody awareness in translation. To address this challenge, we introduce an evaluation methodology and a focused benchmark (named CONTRAPROST) aimed at capturing a wide range of prosodic phenomena. Our methodology uses large language models and controllable text-to-speech (TTS) to generate contrastive examples. Through experiments in translating English speech into German, Spanish, and Japanese, we find that (a) S2TT models possess some internal representation of prosody, but the prosody signal is often not strong enough to affect the translations, (b) E2E systems outperform cascades of speech recognition and text translation systems, confirming their theoretical advantage in this regard, and (c) certain cascaded systems also capture prosodic information in the translation, but only to a lesser extent that depends on the particulars of the transcript's surface form.[1]

## 1   Introduction

Prosody, which includes features like stress, intonation, and rhythm, is crucial for conveying meaning in spoken language beyond the literal words used (Ladd, 1980; Bolinger, 1989). Among others, prosody can direct focus and clarify meaning (Bolinger, 1961; Halliday, 1967), disambiguate

| Example: *These are German teachers.* | | |
|---|---|---|
| A | Prosody | These are GERMAN teachers. |
|   | Explanation | Teachers from Germany |
|   | Translation | Dies sind Deutschlehrer. |
| B | Prosody | These are German TEACHERS. |
|   | Explanation | Teachers that teach German |
|   | Translation | Dies sind deutsche Lehrer. |
| Example: *John laughed at the Party.* | | |
| A | Prosody | John LAUGHED (pause) at the Party. |
|   | Explanation | Laughed while at the party (literal) |
|   | Translation | John lachte während der Party. |
| B | Prosody | John LAUGHED AT (pause) the Party. |
|   | Explanation | Ridiculed the party (idiomatic) |
|   | Translation | John lachte über die Party. |

Table 1: Examples of prosody-aware Speech Translation from English to German.

syntax and sentence structure (Bolinger, 1989), convey the emotional state of the speaker (Banse and Scherer, 1996), and provide useful cues that make communication more effective (Shriberg et al., 1998). For example, the phrase "*Really?*" can express surprise, genuine interest or disbelief, depending on the intonation with which is spoken.

Table 1 illustrates the importance of considering prosody when generating translations in S2TT. Sperber and Paulik (2020) suggest that E2E S2TT systems may have an inherent advantage over cascaded systems in this regard, because only the former have access to the speech signal when making translation decisions. However, our understanding of whether prosody informs translation choices in practice is currently still limited, as prior research on this topic either shows only anecdotal evidence (Huang et al., 2023b), focuses on only a small subset of prosodic phenomena (Zhou et al., 2024; Chen et al., 2024), or considers how prosody informs target-side speech with regards to generated prosody but not lexical choice (§6).

---

* Work done during an internship at Apple.
[1] github.com/apple/ml-speech-is-more-than-words

In this paper, we take steps toward a reliable and comprehensive evaluation methodology, which is one of the most important prerequisites for achieving prosody-aware S2TT. We identify three central challenges that must be addressed: (1) Existing S2TT benchmarks often do not include prosody-rich spontaneous speech and/or do not include translations that are informed by the audio, limiting the extent to which reference translations are influenced by source-side prosody. (2) General-purpose evaluation methods like BLEU (Papineni et al., 2002) and COMET (Guerreiro et al., 2023) are insensitive to the often subtle changes in translation caused by input prosody. (3) Existing prosody-centric benchmarks are difficult to scale to broader coverage of languages and prosodic phenomena, which hinders comprehensive analysis.

To address these challenges, we take inspiration from prior work on behavioral testing (Ribeiro et al., 2020; Ferrando et al., 2023) and contrastive evaluation (Sennrich, 2017). We address the first challenge by synthesizing prosody-rich data that covers a wide range of prosodic phenomena through the use of large language models (LLMs) and controllable TTS (cTTS). We tackle the second challenge by developing a double-contrastive evaluation approach, i.e. a directional behavioral test that relies on minimal pairs (differing only in prosody) to evaluate prosody-awareness in S2TT in isolation. The resulting benchmark, CONTRAPROST (Contrastive Prosody ST), covers a variety of language pairs and prosodic phenomena. Since it is mostly automated, it can be further extended, thus addressing also the third challenge.

To investigate how well current state-of-the-art models understand and leverage prosody, we evaluate S2TT models of various sizes and types, including both E2E and cascaded systems. We find indications that S2TT models represent prosody internally, but this knowledge is often not manifested in the translations. We observe that while tested cascaded systems perform better on traditional evaluation (COMET), E2E models outperform cascaded models on CONTRAPROST. We also find indications that some amount of prosody is carried through transcripts in cascaded setups, but this depends on the particulars of the transcriptions. The most important implication of our findings is the need for exploring improvements of S2TT regarding prosody-awareness, e.g. through auxiliary losses or finetuning on prosody-rich data.

## 2 The CONTRAPROST Benchmark

CONTRAPROST is composed of double-contrastive examples (see Table 1), where each example is composed of a sentence in English that could be semantically ambiguous, along with two different pairs of <speech, translation> that capture contrastive cases of prosody.

As it would be expensive and practically difficult to collect such test data manually, we employ an automatic data generation process, illustrated in Fig. 1. First, we identify several relevant categories where prosody influences sentence semantics in important ways, and construct illustrative examples that reflect the respective phenomena of each category, while highlighting differences in prosody-induced meaning (§2.1). We then prompt GPT-4[2] (OpenAI, 2024) to generate sentences similar to the examples for each subcategory using in-context learning, grounding the generation on different text domains to increase diversity (§2.2). Next, GPT-4 is prompted to translate each prosodic case, while also being given access to the prosodies, meanings and general information of the category, thus acting as a prosody- and context-aware oracle translator (§2.3). Finally, we use the OpenAI TTS API[3] to synthesize the prosodic speech of each case (§2.4). Each generation stage is coupled with filtering and quality assessment to ensure the data are of high quality.

### 2.1 Categorization of Prosodic Phenomena

Below, we summarize the examined prosodic categories. Details and examples are available in the Appendices A and B.

**(1) Sentence Stress.** This is usually manifested through increased loudness, vowel length or higher pitch (Fry, 1955), invoking emphasis on certain words within a sentence, potentially changing the semantics by shifting focus (Wagner, 2020). We further categorize prosodic stress in four subcategories according to the purpose of the stress or its use in disambiguation of linguistic phenomena (see Appendix A.1).

**(2) Prosodic Breaks.** Here we consider the existence or placement of longer breaks in the flow of speech, primarily associated with tempo, that create different phrasal boundaries and help disambiguate syntax and sentence structure (Bolinger, 1989). We follow Hirschberg (2017) and use the

---

[2]GPT-4O-2024-05-13
[3]TTS-1-HD, platform.openai.com/docs/models/tts

Figure 1: The Data Generation process for CONTRAPROST.

subcategories outlined in Appendix A.2.

**(3) Intonation Patterns.** This concerns the modality of the sentence, specifically whether it is a statement (falling tone), or a declarative question (rising tone) (Gunlogson, 2002).

**(4) Emotional Prosody.** A different emotional tone can indicate a speaker's emotional state and thus affect the semantics of the utterance (Banse and Scherer, 1996). Emotional tone is usually manifested through changes in pitch, tempo, and loudness. For example, happiness is associated with higher values in pitch and tempo, while sadness exhibits lower values for pitch, tempo, and loudness (Larrouy-Maestri et al., 2024). Here, we focus on the seven *basic* emotions: happy, sad, angry, disgust, surprisal, fear, and neutral (Ekman and Friesen, 1971; Ekman, 1992), based on which we construct all possible pairs, thus having 21 subcategories.

**(5) Politeness.** The level of politeness can be conveyed by non-verbal cues, and influences the pragmatic context of a conversation. A polite tone is associated with a higher pitch and a smooth rhythm, while an impolite tone is manifested through low pitch, irregular rhythm and very high or low loudness levels (Culpeper et al., 2003; Culpeper, 2011).

## 2.2 Prosodic Example Generation

For each category, we prompt GPT-4 to generate sentences based on hand-crafted category-specific examples. More specifically, we have the LLM generate English sentences, each with two different textual prosodic annotations and respective meanings/interpretations to guide subsequent translation (§2.3). The generated annotations include rich text that indicates different levels of emphasis, pause tags, and special punctuation such as ellipsis, ex-

---

> **Prompt 1: Prosodic Example Generation**
>
> You are a helpful assistant with expert knowledge in linguistics, speech, and prosody. Your task is to come up with examples of English sentences where different prosody would change the meaning of the sentence significantly.[1]
> {Details for Category & Subcategory}[2]
> Here are some examples to guide you:
> {List of Examples}[3]
> Strictly follow these rules:
> {List of Rules}[4]
> Provide a rating of how significant is the difference between the two meanings.[5]
> Generate {n} such examples, with rating as high as possible,[6] in the domain of {domain}.[7]

clamation, or interobang (!?). The sentence itself is generated to be as simple as possible, ending with a full stop or question mark.

The general prompt template is displayed in Prompt 1. It starts with some general information about the task, see superscript (1). The prompt then continues with details describing the current category/subcategory (2). The next part refers to in-context learning (Brown et al., 2020), where we provide a list of illustrative, hand-crafted examples for the LLM to follow (3). In certain subcategories, due to repeated mistakes observed in preliminary explorations, we also provide examples to avoid. In (4) we provide a list of rules for the LLM to adhere to, indicating the desired structure of the sentence and how to use prosodic notation, which might not be obvious from the examples (3). Examples of such rules are "do not include prosodic annotations in the sentence," or "stress different noun-phrases in each prosodic case." We further-

1237

more use *self-criticism* (Huang et al., 2023a) by instructing the model to rate its own generations, according to how different the two prosodic interpretations are (5). Then we instruct the LLM to generate examples that have high scores after self-reflection (6). These scores are also used later during filtering. Finally, to avoid repetitive examples and enhance diversity, we condition the generation on specific text domains (7) (Chung et al., 2023). The list of domains is also generated by GPT-4 based on the context that its subcategory would naturally occur (e.g. *legal testimonies*). For each text domain in the subcategory the LLM then generates $n$ candidate examples. We use several hand-crafted text-based filtering steps to ensure that the examples generated by the LLM at this stage comply with the instructions specified in (4).

## 2.3 Oracle Translation

Recent research on the emerging capabilities of LLM-based MT (Vilar et al., 2023; Alves et al., 2023; Zhang et al., 2023) has shown that LLMs can attain very high translation quality, especially for high-resource languages (Robinson et al., 2023) and including translation factors such as emotions (Brazier and Rouas, 2024), suggesting the possibility that LLMs can be leveraged for prosodic translation synthesis. To obtain the translations of the prosodic cases, we thus utilize GPT-4 as a prosody- and context-aware oracle translator. The LLM is prompted to translate, while having access to the sentence, the textual prosodic annotations (prosody-awareness), and the semantic interpretations (context-awareness). The template prompt is shown in Prompt 2. We provide a list of contraints to the LLM with several goals in mind: (i) avoid generating prosodic annotations in the translations; (ii) avoid translating the interpretations rather than the sentences; (iii) encourage the model to generate different translations for each case; (iv) ensure that differences in the translations are only due to the difference in the prosodies.

Although prosody variants substantially influence sentence semantics, this does not always imply that the ideal translations must differ. In particular, sometimes a translation that leaves semantics ambiguous may be preferred as the most natural translation.[4] As a consequence, constraint (iii) is sometimes overly strict and even in conflict with constraint (iv), leading to changes in the translations

---

[4]This is essentially an instance of the fluency-accuracy trade-off (Lim et al., 2024).

---

> **Prompt 2: Oracle Translation**
>
> You are a helpful assistant with expert knowledge in speech, prosody, linguistics and translation, particularly in English and {Target Lang}. You will be provided with a sentence in English (S) and two different prosodic variations ($S_A$, $S_B$), focused on {Category}, which correspond to two different semantic interpretations.
> Your task is to translate S, $S_A$ and $S_B$ into {Target Lang}, as T, $T_A$, and $T_B$.
> Carry out the translation in these steps:
> (1) Translate S into T.
> (2) Translate $S_A$ to $T_A$ and $S_B$ to $T_B$, by focusing on how T should change in order to reflect the additional information from the prosodies.
> The following constraints should be applied: {List of Constraints}
> The sentence S is: {sentence}
> The two different prosodic variations are:
> $S_A$. {prosody$_A$} ({meaning$_A$})
> $S_B$. {prosody$_B$} ({meaning$_B$})

that do not stem from the prosodies, that are not idiomatic. To account for that, we include a post-editing step, where GPT-4 is instructed to choose the most fitting translation among $\{T, T_A, T_B\}$ for each prosodic case, independently from the other prosodic cases, while having access only the prosody information (Prompt 3). We prompt the LLM to first provide an explanation, before selecting the most appropriate translation, in order to induce *chain-of-thought* reasoning effect (Kojima et al., 2024).

> **Prompt 3: Translation Post-editing**
>
> You are a helpful assistant and an expert translator. You will be provided with a sentence in English and different possible translations in {Target Lang}. The English sentence can contain rich prosodic text with {Category-specific information}, that affects the meaning of the sentence. Your task is to select the most appropriate and prosody-aware translation. First provide a brief explanation of your reasoning and then the index of the selected translation.
> The sentence S to be translated is {sentence} and the candidate translations are: $[T, T_A, T_B]$

After post-editing we remove all examples where the prosodic cases have identical translations, i.e. $(T_A = T_B)$. As an extra measure, we also remove examples where the word length-ratio of the non-

prosodic translation $T$ and one of the prosodic translations $T_A, T_B$ is not within $(0.75, 1.25)$[5]. This aims to remove translations that are overly explanatory, including new bits of information that can be due to the prosody, but are making the translation unnatural (see Table 10 in App. D.2 for examples.).

## 2.4 Controllable Speech Synthesis

We use the OpenAI TTS which can synthesize very natural speech with high-quality audio, offering six different voice profiles. While there are no clear guidelines[6] on how to control prosody, we identified some effective prompting strategies to control the TTS output through trial-and-error (Table 2).

| Effect | TTS Prompting |
| --- | --- |
| Strong Emphasis | *WORD* |
| Normal Emphasis | *word* |
| Slight Emphasis | _word_ |
| Pause | <pause> |
| Statement Intonation | Prepend <statement> |
| Question Intonation | Prepend <question> & Append ???? |
| Emotional/Polite Tone | Prepend & Append Emojis |

Table 2: OpenAI TTS prompting strategies.

To ensure that the generated audio follows the correct wording and exhibits the intended prosodic characteristics we use the following process: First, we generate six candidates (one per voice) for each prosody, discarding invalid candidates (WER $\neq 0$) using an ASR model. Then we estimate prosody quality using category-specific tests in order to rank or filter examples. These tests employ techniques such as forced alignment (Kürzinger et al., 2020), signal processing, punctuation probability, and speech emotion classification. They are explained in detail in Appendix C.

## 3 Contrastive Evaluation

General-purpose MT metrics like BLEU and COMET may be insensitive to subtle changes caused by prosody, and do not allow disentangling prosody awareness from overall translation quality. Thus, to assess how well an S2TT model can handle prosody specifically, we develop a contrastive evaluation framework (Sennrich, 2017). Note that previous work on contrastive evaluation uses a single source and two or more targets (Sennrich, 2017; Vamvas and Sennrich, 2021; Zhou et al., 2024) of which only one is correct. The model likelihood

is then estimated for each target, and models are preferred that assign a better score to the correct example than to the foil(s). Here, we generalize this approach to leverage CONTRAPROST's *double-contrastive* pairs, i.e. two sources and two targets (Fig. 1).

Formally, each double-contrastive pair has two cases $\{X^a, Z, Y^a\}$ and $\{X^b, Z, Y^b\}$, where $X^a, X^b$ are the two different prosodic speech signals, $Z$ is the source text (same for both cases), and $Y^a, Y^b$ are the different translated texts for each case. Thus, each example has two correct pairs $(X^a, Y^a)$, $(X^b, Y^b)$ and two incorrect ones $(X^a, Y^b)$, $(X^b, Y^a)$. We propose the following conditions to assess whether the S2TT model can correctly solve the contrastive example, and to what degree:

$$C_{\mathcal{G}} = \mathbf{1}\Big[ f(Y^a \mid X^a; \theta) - f(Y^b \mid X^a; \theta) > 0$$
$$\text{and } f(Y^b \mid X^b; \theta) - f(Y^a \mid X^b; \theta) > 0 \Big]$$
$$C_{\mathcal{D}} = \mathbf{1}\Big[ f(Y^a \mid X^a; \theta) - f(Y^b \mid X^a; \theta)$$
$$+ f(Y^b \mid X^b; \theta) - f(Y^a \mid X^b; \theta) > 0 \Big]$$

Here, $\mathbf{1}[\cdot]$ is the indicator function, and $f(\cdot) > 0$ is a function that measures the agreement between audio input $X$ and target translation $Y$ under the S2TT model with parameters $\theta$. $C_{\mathcal{G}}$ is a *global* condition, requiring the model to prefer both of the correct pairs versus the incorrect ones according to $f$. $C_{\mathcal{D}}$ is a *directional* condition (Ribeiro et al., 2020) where we require a net positive directional movement for the two comparisons. We expect a model to have a strong internal representation of prosody if it can solve the global condition, and weak representation if it can only solve the directional one.[7]

We consider two different functions $f$ to measure the agreement of $X$ and $Y$.

### 3.1 Contrastive Likelihood

Similar to prior work on contrastive evaluation (Sennrich, 2017; Vamvas and Sennrich, 2021; Zhou et al., 2024) we use the model likelihood to measure the level of agreement between input audio and target text. We obtain the model likelihood $\mathcal{L} \in \mathbb{R}^+$ for a reference $Y = (y_1, \ldots, y_{|Y|})$, given a speech signal $X \in \mathbb{R}^k$ and an E2E S2TT model with parameters $\theta_{\text{E2E}}$. It is defined as the product

---

[5]We use character-based length-ratio for Japanese.

[6]platform.openai.com/docs/guides/text-to-speech

[7]Note that $C_{\mathcal{G}}$ is a sufficient condition for $C_{\mathcal{D}}$.

of the conditional probabilities, normalized by the length of the reference. Formally:

$$\mathcal{L}(Y \mid X; \theta_{\text{E2E}}) = \frac{1}{|Y|} \prod_{i=1}^{|Y|} p_{\theta_{\text{E2E}}}(y_i \mid X, y_{<i})$$

For a cascaded S2TT model we approximate the true likelihood by considering the top-n ASR hypotheses $\mathcal{Z} = \{Z^{(1)}, \ldots, Z^{(n)}\}$. Assuming the lengths of the $\mathcal{Z}$ are generally similar, we get:

$$\mathcal{L}(Y \mid X; \theta_{\text{casc}}) \approx \mathcal{L}(Y \mid \mathcal{Z}; \theta_{\text{MT}})\mathcal{L}(\mathcal{Z} \mid X; \theta_{\text{ASR}})$$

$$\approx \frac{\sum_{j=1}^{n} \left[ \mathcal{L}(Y \mid Z^{(i)}; \theta_{\text{MT}}) \cdot \mathcal{L}(Z^{(i)} \mid X; \theta_{\text{ASR}}) \right]}{\sum_{j=1}^{n} \mathcal{L}(Z^{(i)} \mid X; \theta_{\text{ASR}})}$$

Furthermore, to remove a potential bias of the model against rare translations, we normalize by the unconditioned decoder likelihood of the reference:[8]

$$f_{\overline{\mathcal{L}}}(Y \mid X; \theta) = \frac{\mathcal{L}(Y \mid X; \theta)}{\mathcal{L}(Y \mid \theta)} \quad (1)$$

## 3.2 Contrastive Translation Quality

A common criticism of using model likelihoods is that they do not assess whether the correct output is actually generated in practice, due to teacher forcing. To address this, we propose another function that leverages translation quality estimation (QE) to compare unconstrained autoregressively generated model outputs. We obtain the hypothesis $\hat{Y}$ of input $X$ by generating with the S2TT model $\mathcal{M}_\theta$, and use xCOMET (Guerreiro et al., 2023) to measure the quality of the translation. Thus:

$$f_{\mathcal{Q}}(Y \mid X; \theta) = \mathcal{Q}\big(Y, \mathcal{M}_\theta(X)\big) = \mathcal{Q}(Y, \hat{Y}) \quad (2)$$

The contrastive metrics using $f_{\mathcal{Q}}$ are expected to give us a better insight into how influential prosody is when translating with S2TT models, as compared to using $f_{\mathcal{L}}$ (Eq. 1), since they consider autoregressive generation and beam search.

## 4 Experimental Setup

### 4.1 Data Generation

For prosodic example generation with GPT-4 (§2.2) we used a temperature of 1, and 20 text domains per subcategory. The model was prompted to generate 10 examples[9] for each pair of (subcategory,

domain). The total number of subcategories is 27 (more details in App. A), amounting to 5.5k examples of English sentences with pairs of prosodies and meanings created initially. Then we generated the candidates for the six voices with the TTS (5.5k×6×2 = 66k) and choose the 11k best candidates as described in §2.4. After quality assessment we end up with 2.8k examples with good prosody quality in the generated audio. Then we separately translated each one to the three target languages German (De), Spanish (Es), and Japanese (Ja). After post-editing and filtering we obtained 1.3k–1.4k full examples for each language pair (Table 3).

| Category | En-De | En-Es | En-Ja |
|---|---|---|---|
| Emotional prosody | 373 | 379 | 376 |
| Sentence stress | 277 | 279 | 342 |
| Prosodic breaks | 276 | 252 | 289 |
| Politeness | 212 | 193 | 206 |
| Intonation patterns | 173 | 173 | 173 |
| **Total** | 1,311 | 1,294 | 1,386 |

Table 3: Number of examples for each language pair in CONTRAPROST. More details are in Appendix D.1.

### 4.2 Speech-to-text Translation Models

We evaluated S2TT models that fall under these three categories:

- E2E, where inference is done without an intermediate transcription step. The decoder of this model has full access to the prosody of the input.

- AED-based cascade, which is composed of an attentional encoder-decoder (AED) (Vaswani et al., 2017) ASR model and an MT model. We expect the decoder of the MT model to have limited access to prosody, unless the ASR model is able to encode it in the transcription. This is possible mainly though punctuation, but also when the ASR model is acting more interpretative (i.e. generating synonyms that better fit the prosody rather than the spoken words).

- CTC-based cascade, which uses a CTC encoder (Graves et al., 2006) for the ASR part. The decoder of the MT model is expected to have almost no access to prosody since CTC model outputs are not punctuated and cannot be interpretative.

---

[8]Estimated by using an empty audio for E2E case and empty source text in the MT model for the cascade.
[9]We generated 15/20 examples for intonation patterns/politeness, respectively.

We are evaluating the following S2TT models:

- SEAMLESSM4T (Seamless Communication, 2023b) is a multilingual and multimodal encoder-decoder. It is trained with multi-task learning on ASR, MT, S2TT and also on speech-to-speech translation (S2ST), and can thus be used in either E2E or cascaded (AED) mode.
- XLS-R (Babu et al., 2021) is a multilingual E2E model, of which the encoder is based on WAV2VEC2.0 and its decoder on MBART50 (Tang et al., 2020).
- ZEROSWOT (Tsiamas et al., 2024) is a zero-shot E2E model that connects a WAV2VEC 2.0 CTC encoder and NLLB (NLLB Team, 2022).
- SALMONN (Tang et al., 2024) is an audio LLM that connects WHISPER (Radford et al., 2022) and BEATs (Chen et al., 2023) to the Vicuna LLM (Peng et al., 2023), and can be used as an E2E S2TT model.
- WHISPER & NLLB (AED-based cascade).
- CTC & NLLB (CTC-based cascade) with WAV2VEC 2.0 or HUBERT (Hsu et al., 2021).

We considered different versions of these 6 models, thus evaluating in total 31 S2TT model variants of different sizes and capabilities (App. E).

### 4.3 Metrics

We used beam search with beam size 5 to generate hypotheses. For estimating the conditional likelihood of the cascade (§3.1) we used the top-5 ASR hypotheses. For the contrastive translation quality (§3.2) we used xCOMET-XL[10] (Guerreiro et al., 2023), which is a state-of-the-art neural quality estimation metric based on XLM-R (Conneau et al., 2020). For all evaluated models we present their *contrastive likelihood* and *contrastive translation quality* scores, both *global* and *directional* versions, as a percentage of solved examples. We also evaluate them on standard QE using xCOMET-XL, by using the 2 correct pairs of each example (2.6k samples). For statistical significance testing we used bootstrap resampling (Efron, 1979) with 10k resamples and a 95% confidence interval.

## 5 Experimental Results

In Table 4 we present the results of evaluating a selection of large and recent model versions all three

language pairs. We find that most S2TT models have at least some internal representation of prosody, enabling them to outperform the random baseline of 50% for the directional contrastive likelihood. On the other hand, when we consider autoregressive generation, we observe that the scores for the directional contrastive quality are relatively low[11], indicating that prosody is often not prominent enough in the internal representations of the models for it to be manifested in the generated translations. Furthermore, we find that the task of correctly solving both sub-cases of each example (global agreement) is very challenging for all models, with scores ranging around 10% for both contrastive metrics. We observe that even though the best performing model according to standard evaluation (xCOMET) is a cascade system, it falls behind the best E2E models when considering the contrastive evaluation on CONTRAPROST. This finding illustrates why it is beneficial to separate prosody evaluation from general accuracy evaluation to study the phenomenon, which is further supported by our observation that the prosody and general accuracy metrics are only moderate correlated (see Fig. 5 in App. F).

**Are model type and model size important for prosody-awareness?** We evaluate all 31 S2TT models using *global contrastive quality*, and run a regression analysis with the model type (E2E/AED-cascade/CTC-cascade) and model size as inputs. We use a mixed effects model (Pinheiro and Bates, 2006) to group together each model family, and thus account for random effects, such as the training data and hyperparameters. Specifically:

$$y_{ij} = \beta_0 + \beta_1 S_{ij} + \beta_2 AED_{ij} + \beta_3 CTC_{ij} + u_j + \epsilon_{ij},$$

where $y_{ij}$ is the score of $i$-th model variant of the $j$-th model family, $\beta_0$ is the intercept, $S$ is the log of the model size, AED and CTC are binary variables, $u_j$ is the random effect for $j$-th model family, and $\epsilon_{ij}$ is a residual error term. All scores are available in Table 11 in App. F. In Figure 2 we confirm with statistical significance that the E2E models outperform the cascades in all three language directions.[12] There is also a statistically significant negative impact on prosody-awareness when the cascade is based on a CTC ASR model that may be explained by the absence of punctuation in CTC

---

[10]hf.co/Unbabel/XCOMET-XL

[11]Assuming xCOMET is 0 for randomly generated text, the baseline scores are also 0.

[12]Note that results are borderline non-significant for En-Ja against the AED-cascade.

| Model Name | Contrastive Likelihood | | Contrastive Quality | | xCOMET |
|---|---|---|---|---|---|
| | **Directional** | **Global** | **Directional** | **Global** | |
| *English → German* | | | | | |
| SEAMLESSM4T-v2-LARGE | 61.2 | **13.5** | 37.4 | 14.5 | 0.988 |
| XLS-R 2B | 59.3 | 4.6 | 31.1 | 7.3 | 0.980 |
| ZEROSWOT-LARGE | 60.6 | 9.7 | 29.2 | 8.7 | 0.990 |
| SALMONN-13B | **62.8** | 7.2 | **43.2** | **15.9** | 0.975 |
| SEAMLESSM4T-v2-LARGE | 60.2 | 12.9 | 31.1 | 10.4 | 0.991 |
| WHISPER-v3-LARGE & NLLB-3.3B | 60.7 | 5.8 | 23.1 | 5.5 | **0.992** |
| HUBERT-XL & NLLB-3.3B | 39.4 | 0.5 | 20.5 | 2.6 | 0.979 |
| *English → Spanish* | | | | | |
| SEAMLESSM4T-v2-LARGE | **64.9** | **13.4** | 37.9 | 11.0 | 0.982 |
| XLS-R 2B | 57.6 | 5.6 | 32.0 | 8.4 | 0.930 |
| ZEROSWOT-LARGE | 57.5 | 9.2 | 31.1 | 5.6 | 0.948 |
| SALMONN-13B | 61.3 | 3.6 | **39.6** | **12.3** | 0.967 |
| SEAMLESSM4T-v2-LARGE | 61.3 | 11.7 | 29.5 | 7.6 | 0.984 |
| WHISPER-v3-LARGE & NLLB-3.3B | 63.2 | 2.9 | 25.4 | 4.8 | **0.987** |
| HUBERT-XL & NLLB-3.3B | 41.8 | 0.2 | 20.8 | 2.4 | 0.968 |
| *English → Japanese* | | | | | |
| SEAMLESSM4T-v2-LARGE | 59.4 | **12.4** | 40.3 | 13.8 | 0.956 |
| XLS-R 2B | 60.0 | 4.6 | 27.4 | 7.0 | 0.950 |
| ZEROSWOT-LARGE | 58.8 | 7.9 | 23.6 | 7.9 | **0.970** |
| SALMONN-13B | **60.4** | 10.8 | **46.1** | **16.1** | 0.859 |
| SEAMLESSM4T-v2-LARGE | 59.4 | 9.1 | 31.0 | 8.7 | 0.961 |
| WHISPER-v3-LARGE & NLLB-3.3B | 59.8 | 4.9 | 21.5 | 5.3 | 0.960 |
| HUBERT-XL & NLLB-3.3B | 40.4 | 0.8 | 15.7 | 2.5 | 0.922 |
| *Average* | | | | | |
| SEAMLESSM4T-v2-LARGE | **61.8** | **13.1** | 38.5 | 13.1 | 0.975 |
| XLS-R 2B | 59.0 | 4.9 | 30.2 | 7.6 | 0.953 |
| ZEROSWOT-LARGE | 59.0 | 8.9 | 28.0 | 8.1 | 0.969 |
| SALMONN-13B | 61.5 | 7.2 | **42.9** | **14.8** | 0.933 |
| SEAMLESSM4T-v2-LARGE | 60.3 | 11.2 | 30.5 | 8.9 | 0.979 |
| WHISPER-v3-LARGE & NLLB-3.3B | 61.2 | 4.5 | 23.3 | 5.2 | **0.980** |
| HUBERT-XL & NLLB-3.3B | 40.5 | 0.5 | 19.0 | 2.5 | 0.956 |

Table 4: Contrastive Evaluation of S2TT models on CONTRAPROST. Grey background indicates a cascaded system.

transcripts, which if present can at least approximately signal some prosodic phenomena. Finally, although there is some evidence that larger models are more prosody-aware, results are not statistically significant. We speculate that larger models have more capacity to encode prosody in the weights, but since prosody is perhaps not sufficiently represented in the training data, this effect is limited.

**How do results compare across categories and models?** In Figure 3 we present results across individual prosodic categories for four different English-German models, and perform pairwise model comparisons via bootstrap resampling[13]. The only category models are able to solve consistently is *intonation patterns*, which can also be solved by cascaded models due to the presence of

punctuation in the transcription. The comparably lower scores in the other four categories further demonstrate the inability of current state-of-the-art models to use prosody, with *sentence stress* being the most challenging. Through the pairwise comparisons, we find that an LLM-based model (SALMONN) is not statistically different from a more standard S2TT model, like SEAMLESSM4T. Next, comparing the SEAMLESSM4T model in both E2E and cascade allows us to control for parameters such as training data and architecture, in order to observe the effect of model type, giving more clarity of our results on the theoretical advantage of E2E models. Finally, we observe a clear performance gain by using the SEAMLESSM4T cascade over the WHISPER & NLLB one. We hypothesize this advantage is due to the multitasking nature of SEAMLESSM4T, which makes its ASR

---

[13]English-Spanish/Japanese are available at Figures 6, 7 in App. F.

Figure 2: Regression Analysis of model types and model sizes per language pair.

mode more interpretative than standard ASR models. This allows the ASR part of the cascade to escape the word-by-word paradigm, and use more fitting words in the transcription (such as synonyms) that fit better the prosody of the audio. Supporting this hypothesis. we observe a worse WER score for SEAMLESSM4T (11%) compared to WHISPER (4%).

**Is the level of prosody-awareness language-dependant?** In Figure 4 we carry out a similar regression analysis as in Figure 2, but with the language pair as an independent categorical variable. Interestingly, we observe that there are differences between the three language pairs, and also significant for Spanish vs. German, which indicates that prosody-awareness in S2TT could be language-dependant. We hypothesize that the expressivity of the target language might be a relevant factor, since more expressive languages might be able to easier encode the prosody of the source speech into text.

## 6 Related Work

Prosody has traditionally been an important topic for TTS research (Kohler, 1991), either for transferring (Skerry-Ryan et al., 2018) or encoding it (Pamisetty and Sri Rama Murty, 2022) in the synthesized speech. Furthermore, Torresquintero et al. (2021) created a dataset for evaluating prosody transfer in TTS models, which contains several categories, similar to our study here. Naturally prosody has also been the focus of S2ST systems, in order to translate in a more expressive way (Aguero et al., 2006; Do et al., 2017; Communication et al., 2023). The topic has received less

attention in the context of S2TT. Chen et al. (2024) present a dataset for emotional prosody based on speech and translations from TV series, and show that finetuning with emotion labels, can improve translation quality. Zhou et al. (2024) studied the prosody-awareness of WHISPER in E2E and cascade mode, in translating Korean *wh-phrases* using contrastive likelihood, and find evidence of the E2E model outperforming the cascade. Here we contribute a broader study of prosody in S2TT, by proposing a double-contrastive benchmark that covers several prosodic categories, the use of more generative-like contrastive evaluation, and evaluating a plethora of S2TT models. Finally, de Seyssel et al. (2023) present a benchmark for evaluating prosody-awareness in self-supervised acoustic representations. Similarly to our study they present evidence of prosody awareness in the representations. Contrary to our results, they conclude that size has a positive effect on prosody awareness.

## 7 Conclusions

We presented CONTRAPROST, a benchmark based on double-contrastive examples for evaluating prosody-awareness in S2TT models, covering several categories and languages. In addition to standard contrastive evaluation based on model likelihoods, we proposed a generative contrastive metric based on quality estimation. We evaluated a plethora of models, and found that they exhibit some signs of prosody-awareness, but the effect is often not strong enough to influence the translations. We also confirmed the previously hypothesized inherent advantage of E2E models com-

Figure 3: Upper: Model performance per category (En-De). Lower Model performance comparisons (En-De), (a): SALMONN-13B vs. SEAMLESSM4T-V2-LARGE, (b) SEAMLESSM4T-V2-LARGE(E2E) vs. SEAMLESSM4T-V2-LARGE(cascade), (c) SEAMLESSM4T-V2-LARGE(cascade) vs. WHISPER-V3-LARGE/NLLB-3.3B.



Figure 4: Regression Analysis of language pairs.

pared to cascaded models. We hope that our benchmark and findings will motivate more research into prosody-aware S2TT in the future, enabling us to better understand it and improve it.

## Limitations

For creating CONTRAPROST we relied on an almost entirely automated data generation process. This allowed us to create a comprehensive dataset covering several prosodic phenomena and three language pairs, in a fast and cost-effective way. It would also enable expanding the coverage of lan-

guages and prosodic phenomena relatively easy in the future. Nevertheless, despite our best efforts regarding filtering and quality assessment (§2 and App. C), the data is not perfect and includes a certain amount of noise. We observed the following sources of noise in order of decreasing importance: (1) prosody not prominent in the generated speech; (2) translations overly explanatory or not encoding prosody; (3) semantic interpretations of the two cases rather similar. We do not expect these issues to be so frequent as to alter the findings of this work in a systematic way, but additional human annotation or verification would be a valuable step for future work. Furthermore, as the landscape of available generative models, in particular controllable TTS, is changing quickly, the quality of results using our data generation process would expectantly become less of a concern in future iterations.

Our study follows a contrastive evaluation methodology in order to isolate prosody-related behavior. As a consequence, our study does not allow drawing conclusions on how much prosody matters in real life data, and in what domains it is especially important. In addition, we hypothesize that some prosodic phenomena could be correctly translated by having access to the broader context of the conversation (context-aware S2TT), which we leave for future research.

# References

P.D. Aguero, J. Adell, and A. Bonafonte. 2006. Prosody Generation for Speech-to-Speech Translation. In *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*, volume 1, pages I–I.

Duarte Alves, Nuno Guerreiro, João Alves, José Pombal, Ricardo Rei, José de Souza, Pierre Colombo, and Andre Martins. 2023. Steering Large Language Models for Machine Translation with Finetuning and In-Context Learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11127–11148, Singapore. Association for Computational Linguistics.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. Common Voice: A Massively-Multilingual Speech Corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4218–4222, Marseille, France. European Language Resources Association.

Arun Babu, Changhan Wang, Andros Tjandra, Kushal Lakhotia, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, Alexei Baevski, Alexis Conneau, and Michael Auli. 2021. XLS-R: Self-supervised Cross-lingual Speech Representation Learning at Scale. *Preprint*, arXiv:2111.09296.

Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

Rainer Banse and Klaus R Scherer. 1996. Acoustic Profiles in Vocal Emotion Expression. *Journal of personality and social psychology*, 70(3):614.

Dwight Bolinger. 1989. *Intonation and Its Uses*. Stanford University Press, Redwood City.

Dwight L. Bolinger. 1961. Contrastive Accent and Contrastive Stress. *Language*, 37(1):83–96.

Charles Brazier and Jean-Luc Rouas. 2024. Conditioning LLMs with Emotion in Neural Machine Translation. In *Proceedings of the 21st International Conference on Spoken Language Translation (IWSLT 2024)*, pages 33–38, Bangkok, Thailand (in-person and online). Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, Wanxiang Che, Xiangzhan Yu, and Furu Wei. 2023. BEATs: Audio Pre-training with Acoustic Tokenizers. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Sirou Chen, Sakiko Yahata, Shuichiro Shimizu, Zhengdong Yang, Yihang Li, Chenhui Chu, and Sadao Kurohashi. 2024. MELD-ST: An emotion-aware speech translation dataset. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 10118–10126, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

John Chung, Ece Kamar, and Saleema Amershi. 2023. Increasing Diversity While Maintaining Accuracy: Text Data Generation with Large Language Models and Human Interventions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 575–593, Toronto, Canada. Association for Computational Linguistics.

Yu-An Chung, Yu Zhang, Wei Han, Chung-Cheng Chiu, James Qin, Ruoming Pang, and Yonghui Wu. 2021. w2v-BERT: Combining Contrastive Learning and Masked Language Modeling for Self-Supervised Speech Pre-Training. In *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 244–250.

Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, John Hoffman, Min-Jae Hwang, Hirofumi Inaguma, Christopher Klaiber, Ilia Kulikov, Pengwei Li, Daniel Licht, Jean Maillard, Ruslan Mavlyutov, Alice Rakotoarison, Kaushik Ram Sadagopan, Abinesh Ramakrishnan, Tuan Tran, Guillaume Wenzek, Yilin Yang, Ethan Ye, Ivan Evtimov, Pierre Fernandez, Cynthia Gao, Prangthip Hansanti, Elahe Kalbassi, Amanda Kallet, Artyom Kozhevnikov, Gabriel Mejia Gonzalez, Robin San Roman, Christophe Touret, Corinne Wong, Carleigh Wood, Bokai Yu, Pierre Andrews, Can Balioglu, Peng-Jen Chen, Marta R. Costa-jussà, Maha Elbayad, Hongyu Gong, Francisco Guzmán, Kevin Heffernan, Somya Jain, Justine Kao, Ann Lee, Xutai Ma, Alex Mourachko, Benjamin Peloquin, Juan Pino, Sravya Popuri, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Anna Sun, Paden Tomasello, Changhan Wang, Jeff Wang, Skyler Wang, and Mary Williamson. 2023. Seamless: Multilingual Expres-

sive and Streaming Speech Translation. *Preprint*, arXiv:2312.05187.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Jonathan Culpeper. 2011. *"It's not what you said, it's how you said it!" Prosody and Impoliteness*, pages 57–84. De Gruyter Mouton, Berlin, New York.

Jonathan Culpeper, Derek Bousfield, and Anne Wichmann. 2003. Impoliteness Revisited: With Special Reference to Dynamic and Prosodic Aspects. *Journal of Pragmatics*, 35:1545–1579.

Maureen de Seyssel, Marvin Lavechin, Hadrien Titeux, Arthur Thomas, Gwendal Virlet, Andrea Santos Revilla, Guillaume Wisniewski, Bogdan Ludusan, and Emmanuel Dupoux. 2023. ProsAudit, a prosodic benchmark for self-supervised speech models. In *Proc. INTERSPEECH 2023*, pages 2963–2967.

Nicole Dehé. 2014. *Parentheticals in Spoken English : The Syntax-Prosody Relation*. Cambridge [u.a.] : Cambridge University Press.

Quoc Truong Do, Sakriani Sakti, and Satoshi Nakamura. 2017. Toward Expressive Speech Translation: A Unified Sequence-to-Sequence LSTMs Approach for Translating Words and Emphasis. In *Proc. Interspeech 2017*, pages 2640–2644.

B. Efron. 1979. Bootstrap Methods: Another Look at the Jackknife. *The Annals of Statistics*, 7(1):1 – 26.

Paul Ekman. 1992. Facial Expressions of Emotion: New Findings, New Questions. *Psychological Science*, 3(1):34–38.

Paul Ekman and Wallace V Friesen. 1971. Constants across cultures in the face and emotion. *Journal of personality and social psychology*, 17(2):124.

Javier Ferrando, Matthias Sperber, Hendra Setiawan, Dominic Telaar, and Saša Hasan. 2023. Automating Behavioral Testing in Machine Translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1014–1030, Singapore. Association for Computational Linguistics.

D. B. Fry. 1955. Duration and Intensity as Physical Correlates of Linguistic Stress. *The Journal of the Acoustical Society of America*, 27(4):765–768.

Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, page 369–376, New York, NY, USA. Association for Computing Machinery.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xCOMET: Transparent Machine Translation Evaluation through Fine-grained Error Detection. *Preprint*, arXiv:2310.10482.

Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. 2020. Conformer: Convolution-augmented Transformer for Speech Recognition. In *Proc. Interspeech 2020*, pages 5036–5040.

Christine Gunlogson. 2002. Declarative questions. In *Proceedings of Semantics and Linguistic Theory (SALT) XII*, pages 124–143, Ithaca, NY. CLC Publications.

M. A. K. Halliday. 1967. Notes on transitivity and theme in English. Part 1 and 2. *Journal of Linguistics*, 3:199–244.

Julia Hirschberg. 2017. Pragmatics and Prosody (Chapter 28). In *The Oxford Handbook of Pragmatics*. Oxford University Press.

Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 29:3451–3460.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *International Conference on Learning Representations*.

Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2023a. Large Language Models Can Self-Improve. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1051–1068, Singapore. Association for Computational Linguistics.

Zhichao Huang, Rong Ye, Tom Ko, Qianqian Dong, Shanbo Cheng, Mingxuan Wang, and Hang Li. 2023b. Speech Translation with Large Language Models: An Industrial Practice. *arXiv preprint arXiv:2312.13585*.

Hirofumi Inaguma, Sravya Popuri, Ilia Kulikov, Peng-Jen Chen, Changhan Wang, Yu-An Chung, Yun Tang, Ann Lee, Shinji Watanabe, and Juan Pino. 2023. UnitY: Two-pass direct speech-to-speech translation with discrete units. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15655–15680, Toronto, Canada. Association for Computational Linguistics.

Ray Jackendoff. 1972. *Semantic Interpretation in Generative Grammar*. MIT Press, Cambridge, Massachusetts.

Wouter Jansen, Michelle L. Gregory, and Jason M. Brenier. 2001. Prosodic correlates of directly reported speech: Evidence from conversational speech. In *Proc. ITRW on Prosody in Speech Recognition and Understanding*, page paper 14.

J. Kahn, M. Riviere, W. Zheng, E. Kharitonov, Q. Xu, P.E. Mazare, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen, T. Likhomanenko, G. Synnaeve, A. Joulin, A. Mohamed, and E. Dupoux. 2020. Libri-Light: A Benchmark for ASR with Limited or No Supervision. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE.

G. Klewitz and E. Couper-Kuhlen. 1999. Quote-unquote? The role of prosody in the contextualization of reported speech sequences. *Pragmatics*, 9(4):459–485.

K.J. Kohler. 1991. Prosody in speech synthesis: the interplay between basic research and TTS application. *Journal of Phonetics*, 19(1):121–138. Speech Synthesis and Phonetics.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2024. Large Language Models are Zero-shot Reasoners. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.

Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. HiFi-GAN: Generative Adversarial Networks for Efficient and High Fidelity Speech Synthesis. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

Ludwig Kürzinger, Dominik Winkelbauer, Lujun Li, Tobias Watzel, and Gerhard Rigoll. 2020. CTC-Segmentation of Large Corpora for German End-to-End Speech Recognition. In *Speech and Computer*, pages 267–278, Cham. Springer International Publishing.

J.D.R. Ladd. 1980. *The Structure of Intonational Meaning: Evidence from English*. Indiana University Press, Bloomington.

Pauline Larrouy-Maestri, David Poeppel, and Marc D. Pell. 2024. The Sound of Emotional Prosody: Nearly 3 Decades of Research and Future Directions. *Perspectives on Psychological Science*, 0(0):17456916231217722. PMID: 38232303.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Mark Y. Liberman and Richard Sproat. 1992. The Stress and Structure of Modified Noun Phrases in English. In *Lexical Matters*.

Zheng Wei Lim, Ekaterina Vylomova, Trevor Cohn, and Charles Kemp. 2024. Simpson's paradox and the accuracy-fluency tradeoff in translation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 92–103, Bangkok, Thailand. Association for Computational Linguistics.

Steven R. Livingstone and Frank A. Russo. 2018. The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLOS ONE*, 13(5):1–35.

Marina Nespor and Irene Vogel. 1986. Prosodic Phonology. *Phonology*, 5(1):161–168.

Giridhar Pamisetty and K. Sri Rama Murty. 2022. Prosody-TTS: An End-to-End Speech Synthesis System with Prosody Control. *Circuits Syst. Signal Process.*, 42(1):361–384.

Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. LibriSpeech: An ASR Corpus Based on Public Domain Audio Books. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5206–5210.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction Tuning with GPT-4. *Preprint*, arXiv:2304.03277.

Gabriel Peyré and Marco Cuturi. 2019. *Computational Optimal Transport: With Applications to Data Science*. Now Foundations and Trends.

Jose Pinheiro and Douglas M. Bates. 2006. *Mixed-effects Models in S and S-PLUS*. Statistics and Computing. Springer Science & Business Media, New York.

Patti Price, Mari Ostendorf, Stefanie Shattuck-Hufnagel, and Cynthia Fong. 1991. The Use of Prosody in Syntactic Disambiguation. In *Proceedings of the Workshop on Speech and Natural Language*, HLT '91, page 372–377, USA. Association for Computational Linguistics.

Joel Pynte. 1996. Prosodic Breaks and Attachment Decisions in Sentence Parsing. *Language and Cognitive Processes*, 11(1-2):165–192.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust Speech Recognition via Large-Scale Weak Supervision. *Preprint*, arXiv:2212.04356.

Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.

Nathaniel Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. ChatGPT MT: Competitive for High- (but Not Low-) Resource Languages. In *Proceedings of the Eighth Conference on Machine Translation*, pages 392–418, Singapore. Association for Computational Linguistics.

Rico Sennrich. 2017. How Grammatical is Character-level Neural Machine Translation? Assessing MT Quality with Contrastive Translation Pairs. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 376–382, Valencia, Spain. Association for Computational Linguistics.

Elizabeth Shriberg, Andreas Stolcke, Daniel Jurafsky, Noah Coccaro, Marie Meteer, Rebecca Bates, Paul Taylor, Klaus Ries, Rachel Martin, and Carol van Ess-Dykema. 1998. Can Prosody Aid the Automatic Classification of Dialog Acts in Conversational Speech? *Language and Speech*, 41(3-4):443–492. PMID: 10746366.

RJ Skerry-Ryan, Eric Battenberg, Ying Xiao, Yuxuan Wang, Daisy Stanton, Joel Shor, Ron Weiss, Rob Clark, and Rif A. Saurous. 2018. Towards End-to-End Prosody Transfer for Expressive Speech Synthesis with Tacotron. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4693–4702. PMLR.

Matthias Sperber and Matthias Paulik. 2020. Speech Translation and the End-to-End Promise: Taking Stock of Where We Are. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7409–7421, Online. Association for Computational Linguistics.

Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. 2024. SALMONN: Towards Generic Hearing Abilities for Large Language Models. In *The Twelfth International Conference on Learning Representations*.

Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. Multilingual Translation with Extensible Multilingual Pretraining and Finetuning. *Preprint*, arXiv:2008.00401.

Alexandra Torresquintero, Tian Huey Teh, Christopher G.R. Wallis, Marlene Staib, Devang S. Ram Mohan, Vivian Hu, Lorenzo Foglianti, Jiameng Gao, and Simon King. 2021. ADEPT: A Dataset for Evaluating Prosody Transfer. In *Proc. Interspeech 2021*, pages 3880–3884.

Ioannis Tsiamas, Gerard Gállego, José Fonollosa, and Marta Costa-jussà. 2024. Pushing the Limits of Zero-shot End-to-End Speech Translation. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 14245–14267, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Jannis Vamvas and Rico Sennrich. 2021. On the Limits of Minimal Pairs in Contrastive Evaluation. In *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 58–68, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for Translation: Assessing Strategies and Performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15406–15427, Toronto, Canada. Association for Computational Linguistics.

Michael Wagner. 2020. Prosodic Focus. In Daniel Gutzmann, Lisa Matthewson, Ceclia Meier, Hotze Rullmann, and Thomas E. Zimmermann, editors, *The Wiley Blackwell Companion to Semantics*. Wiley–Blackwell.

Changhan Wang, Anne Wu, Jiatao Gu, and Juan Pino. 2021. CoVoST 2 and Massively Multilingual Speech Translation. In *Proc. Interspeech 2021*, pages 2247–2251.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Qiantong Xu, Alexei Baevski, Tatiana Likhomanenko, Paden Tomasello, Alexis Conneau, Ronan Collobert,

Gabriel Synnaeve, and Michael Auli. 2020. Self-training and Pre-training are Complementary for Speech Recognition. *Preprint*, arXiv:2010.11430.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting Large Language Model for Machine Translation: A Case Study. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Giulio Zhou, Tsz Kin Lam, Alexandra Birch, and Barry Haddow. 2024. Prosody in Cascade and Direct Speech-to-Text Translation: a case study on Korean Wh-Phrases. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 674–683, St. Julian's, Malta. Association for Computational Linguistics.

## A  Prosodic Subcategories

Here we expand the categorization of §2.1, and discuss the identified subcategories for sentence stress and prosodic breaks, which are 4 and 6 respectively. Intonation patterns and Politeness do not have subcategories. For emotional prosody we have 15 emotion pairs[14], thus having 15 subcategories. Examples are available at Tables 5 and 6.

### A.1  Sentence Stress Subcategories

(1.1) *Contrastive Stress*, which highlights differences or corrects previous statements, emphasizing contrasts between elements (Bolinger, 1961).

(1.2) *New vs. Given Information*, which differentiates between new and given information, emphasizing what is considered new (Halliday, 1967).

(1.3) *Relational vs. Descriptive Adjectives*, where stressing the adjective or the noun can differentiate between the relational and descriptive uses of attributive adjectives (Liberman and Sproat, 1992).

(1.4) *Focus-Sensitive Operators*, where stress indicates the focus of adverbs of quantification (*only*, *just*, etc), shifting the meaning of the sentence accordingly (Halliday, 1967; Jackendoff, 1972).

### A.2  Prosodic Break Subcategories

(2.1) *Direct vs. Indirect Statements*, where a prosodic break can indicate whether a phrase is a direct or an indirect quote (Klewitz and Couper-Kuhlen, 1999; Jansen et al., 2001).

(2.2) *Restrictive vs. Non-Restrictive Clauses*, which involves the use of prosodic breaks to differentiate between essential and non-essential information, impacting the specificity of the noun being described (Nespor and Vogel, 1986).

(2.3) *VP vs. NP Attachment*, where a trailing phrase can be attached either to the verb-phrase or the noun-phrase, depending on the existence of a prominent prosodic break (Pynte, 1996).

(2.4) *Particle vs. Preposition*, where a prosodic break can disambiguate between the literal and idiomatic meaning of phrasal verbs, by grouping the preposition with or without it (Price et al., 1991).

(2.5) *Broad vs. Narrow Scope*, where the existence of a prosodic break can signal that a modifier (adjective) has narrow scope, and refers only to one of two nouns that follow it (Hirschberg, 2017).

(2.6) *Complementizer vs. Parenthetical*, where the location of a prosodic break indicates whether an intermediate phrase acts as a complementizer or simply parenthetical to the main one (Dehé, 2014).

## B  Examples for In-context Learning

In Tables 5, 6 and 7 we present some of the examples used for in-context learning when generating new examples with GPT-4 (§2.2).

## C  Quality Assessment for TTS candidates

Here we present the objectives we defined for assessing the quality of the generated speech candidates for each contrastive example. The objective is applied only to candidates that had WER = 0 using WHISPER. If all candidates are invalid for a prosodic case, the whole example is removed. We also defined some threshold levels for the objectives after trial-and-error, in order to remove examples where the best candidate was below it.

**Sentence Stress.** We use forced-alignment with WAV2VEC 2.0 (Baevski et al., 2020) to obtain the segment for each word in the signal, and extract their loudness, pitch and duration features. Then we define the stress level *stress* for a word $w$ as the weighted sum of these three features. Finally we select the best candidate according to a simple objective $obj_{stress}$ that has three goals: (1) maximize the stress of the target word ($stress_{tgt}$), (2) minimize the stress of the target word of the contrastive case ($stress_{foil}$), and (3) minimize the average stress of the rest.

$$stress_w = \lambda_1 loud_w + \lambda_2 pitch_w + \lambda_3 dur_w$$
$$obj_{stress} = 2 \cdot stress_{tgt} - stress_{foil}$$
$$- \frac{1}{n-1} \sum_{w \neq tgt} stress_w,$$

where we used $\lambda_1 = 0.5$, $\lambda_2 = 0.3$, and $\lambda_3 = 0.2$. Note that in the sentence stress examples, there is

---

[14]Removed *fearful* emotion due to issues with the TTS.

| 1.1 Contrastive Stress | |
|---|---|
| Sentence | She didn't give the book to John. |
| Prosody$_A$ | She didn't give the *BOOK* to John. |
| Meaning$_A$ | Something else was given to John. |
| Prosody$_B$ | She didn't give the book to *JOHN*. |
| Meaning$_B$ | The book was given to someone else. |
| **1.2 New vs. Given Information** | |
| Sentence | The committee decided to postpone the meeting. |
| Prosody$_A$ | The *COMMITTEE* decided to postpone the meeting. |
| Meaning$_A$ | Given: Someone decided to postpone the meeting; New: It was the committee who decided. |
| Prosody$_B$ | The committee decided to *POSTPONE* the meeting. |
| Meaning$_B$ | Given: The committee decided something; New: The decision was to postpone it. |
| **1.3 Relational vs. Descriptive Adjectives** | |
| Sentence | They are German teachers. |
| Prosody$_A$ | They are *GERMAN* teachers. |
| Meaning$_A$ | Teachers who teach the German language. (Relational) |
| Prosody$_B$ | They are German *TEACHERS*. |
| Meaning$_B$ | Teachers who are German. (Descriptive) |
| **1.4 Focus-Sensitive Operators** | |
| Sentence | I only introduced John to Maria at yesterday's party. |
| Prosody$_A$ | I only introduced *JOHN* to Maria at yesterday's party. |
| Meaning$_A$ | John was the only person I introduced to Maria. |
| Prosody$_B$ | I only introduced John to *MARIA* at yesterday's party. |
| Meaning$_B$ | Maria was the only person I introduced John to. |

Table 5: Examples in the category *Sentence Stress* that were used for in-context learning.

always exactly 1 target word in each contrastive prosodic case.

**Prosodic Breaks.** Likewise, after forced-alignment, we measure the duration *dur* of each gap $l$ between the words in the utterance, and define a similar objective $obj_{break}$ as:

$$obj_{break} = 2\frac{1}{|tgt|}\sum_{l \in tgt} dur_l - \frac{1}{|foil|}\sum_{l \in foil} dur_l$$
$$- \frac{1}{n - |tgt|}\sum_{l \notin tgt} dur_l$$

In this category, there can be 0 to 2 breaks in each prosodic case, which could be shared between the two prosodic cases. In the objective we consider only the ones that are not common in the two cases.

**Intonation Patterns.** We use teacher-forcing with WHISPER to extract the punctuation probabilities given the transcription text without the ending punctuation. The probability of the sentence to be a statement is the sum of the probabilities of the tokens "." and "!", while the probability of a question is the probability of the token "?". Thus

the objective $obj_{inton}$ for a statement is defined as:

$$obj_{inton} = p(. \mid X, Z_{<n}) + p(! \mid X, Z_{<n})$$
$$- p(? \mid X, Z_{<n}),$$

where $X$ is the speech signal and $Z_{<n}$ are the tokens of the transcription, excluding the final one, which corresponds in all cases of this category. to the punctuation. The negative objective $-obj_{inton}$ is used for a case that is a question.

**Emotional Prosody.** We employ an emotion classifier[15] which is a based on a finetuned WAV2VEC 2.0 on the RAVDESS dataset (Livingstone and Russo, 2018), and define the objective as:

$$obj_{emo} = p(e_{tgt} \mid X) - p(e_{foil} \mid X),$$

where $\theta$ are the parameters of the classifier, $e_{tgt}$ is the target emotion label and $e_{foil}$ is the emotion label of the other prosodic case.

**Pragmatic Prosody.** To the best of our knowledge there is no open-sourced audio classifier to detect politeness levels, thus we re-purpose the emotion classifier and define the probabilities of politeness

---

[15]hf.co/ehcalabres/wav2vec2-lg-XLS-R-en-speech-emotion-recognition

| | **2.1 Direct vs. Indirect Statements** |
|---|---|
| Sentence | Alex announced Jamie will meet the manager. |
| Prosody$_A$ | Alex *ANNOUNCED* \| Jamie will meet the manager. |
| Meaning$_A$ | (Direct Statement) |
| Prosody$_B$ | Alex announced Jamie will meet the manager. |
| Meaning$_B$ | (Indirect Statement) |
| | **2.2 Restrictive vs. Non-Restrictive Phrases** |
| Sentence | The students who were talking were sent out. |
| Prosody$_A$ | The students who were *TALKING* \| were sent out. |
| Meaning$_A$ | Only the students who were talking were actually sent out. (Restrictive) |
| Prosody$_B$ | The *STUDENTS* \| who were talking \| were sent out. |
| Meaning$_B$ | All students were sent out, and the fact they were talking is additional information. (Non-restrictive) |
| | **2.3 Verb-phrase vs. Noun-phrase Attachment** |
| Sentence | Paula phoned her friend from Alabama. |
| Prosody$_A$ | Paula phoned her friend \| from *ALABAMA*. |
| Meaning$_A$ | Paula called her friend while she was in Alabama. (VP Attachment) |
| Prosody$_B$ | Paula phoned \| her *FRIEND* from Alabama. |
| Meaning$_B$ | Paula phoned her friend who is from Alabama. (NP Attachment) |
| | **2.4 Phrasal Verbs** |
| Sentence | John laughed at the party. |
| Prosody$_A$ | John *LAUGHED* \| at the party. |
| Meaning$_A$ | John laughed while he was at the party. (Literal) |
| Prosody$_B$ | John *LAUGHED AT* \| the party. |
| Meaning$_B$ | John made fun of the party. (Idiomatic) |
| | **2.5 Complementizer vs. Parenthetical** |
| Sentence | We only suspected they all knew that a burglary had been committed. |
| Prosody$_A$ | We only *SUSPECTED* \| they all knew that a burglary had been committed. |
| Meaning$_A$ | The suspicion was that they all knew about the burglary. (Complementizer) |
| Prosody$_B$ | We only suspected \| they all *KNEW* \| that a burglary had been committed. |
| Meaning$_B$ | They all knew that we only suspected that a burglary had been committed. (Parenthetical) |
| | **2.6 Modifier Scope** |
| Sentence | This collar is dangerous to younger dogs and cats. |
| Prosody$_A$ | This collar is dangerous to *YOUNGER* dogs and cats. |
| Meaning$_A$ | Younger refers to both dogs and cats. (Broad Scope) |
| Prosody$_B$ | This collar is dangerous to *YOUNGER* dogs \| and *CATS*. |
| Meaning$_B$ | Younger refers only to dogs. (Narrow Scope) |

Table 6: Examples in the category *Prosodic Breaks* that were used for in-context learning.

and impoliteness as a weighted sum of the 8 available emotion classes.

$$p(polite) = \frac{\sum_e w_e p(e \mid X)}{\sum_e w_e},$$

and similarly for impolite. We used the weighted scheme displayed in Table 8, which was obtained by prompting GPT-4.

## D   Data

### D.1   Data Statistics

In Table 9 we provide the analytic data statistics for each category/subcategory, throughout the generation process stages. The poor quality of the cTTS, where prosody was not always encoded in the speech, led us to remove a large percentage of the examples before translating them. Also many examples where removed because the oracle translations for both cases were the same.

### D.2   Overly Explanatory Examples

In Table 10 we present two examples where GPT-4 acting as an oracle translator (§2.3) proposed overly explanatory translations in the emotional prosody category. Both are inline with the emotion of the speaker, but they contain new bits of information, not initially there. These were removed in filtering due to excessive word-length ratio between the two cases.

| 3. Intonation Patterns | |
|---|---|
| Sentence | You can solve this problem |
| Prosody$_A$ | You *CAN* solve this problem. |
| Meaning$_A$ | Encouraging or asserting the person's ability to solve this problem. |
| Prosody$_B$ | You _can_ solve this problem? |
| Meaning$_B$ | Questioning the person's ability to solve this problem. |
| **4. Emotional Prosody (Happy/Sad)** | |
| Sentence | The surgery went as expected. |
| Prosody$_A$ | <happy> The surgery went *AS EXPECTED*! |
| Meaning$_A$ | The surgery's successful outcome aligns with hopes and predictions, leading to joy and relief. |
| Prosody$_B$ | <sad> The surgery went _as expected_ ... |
| Meaning$_B$ | The expected outcome was not favorable, leading to a somber tone. |
| **4. Emotion Prosody (Fearful/Angry)** | |
| Sentence | Can we talk about this later? |
| Prosody$_A$ | <fearful> Can we... talk about this... later? |
| Meaning$_A$ | Indicates hesitation or fear about the topic, or the situation in general. |
| Prosody$_B$ | <angry> Can we *TALK* about this later!? |
| Meaning$_B$ | Implies urgency or frustration, and a demand for immediate attention. |
| **5. Politeness** | |
| Sentence | Can you move your car? |
| Prosody$_A$ | <polite> Can you _move_ your car? |
| Meaning$_A$ | A polite request to move the car. |
| Prosody$_B$ | <impolite> Can you *MOVE* your *CAR*?! |
| Meaning$_B$ | A rude demand to move the car, with an aggressive tone. |

Table 7: Examples in the categories *Intonation Patterns*, *Emotional Prosody*, *Politeness*, and that were used for in-context learning.

| Emotion | Politeness | Impoliteness |
|---|---|---|
| Happy | 0.3 | -0.1 |
| Calm | 0.3 | -0.2 |
| Neutral | 0.2 | 0.1 |
| Surprised | 0.1 | 0.1 |
| Sad | 0.0 | 0.2 |
| Disgust | -0.1 | 0.3 |
| Angry | -0.2 | 0.4 |
| Fearful | -0.1 | 0.0 |

Table 8: Weighting scheme for Politeness and Impoliteness labels based on the emotion classifier.

# E  Evaluated Speech Translation Models

Here we describe in more detail the model families and the specific versions used. We evaluated in total 31 S2TT model variants. All models are available in the Transformers Huggingface Library (Wolf et al., 2020). For inference we used the default generation parameters and a beam search of 5.

1. SEAMLESSM4T (Seamless Communication, 2023a) and its updated version v2 (Seamless Communication, 2023b) is a recently proposed family of unified encoder-decoder models that are both multilingual (many-to-many, 100 languages) and multimodal (speech/text input or output), meaning they can carry out the tasks of ASR, TTS, MT, S2TT, and also S2ST. The architecture is composed of a text encoder, text decoder, speech encoder, and speech decoder, and different parts are active depending on the input/output modalities. The text encoder-decoder is based on NLLB (NLLB Team, 2022), the speech encoder on a newly proposed conformer (Gulati et al., 2020) w2v-BERT (Chung et al., 2021), and the speech decoder on a unit decoder (Inaguma et al., 2023) and a HiFi-GAN vocoder (Kong et al., 2020). The original version has a medium (1.2B)[16] and a large (2.3B)[17] variant, while the updated v2 has a large variant (2.3B)[18]. For cascade S2TT we first use the model in ASR mode, and then the

---

[16] hf.co/facebook/seamless-m4t-medium
[17] hf.co/facebook/seamless-m4t-large
[18] hf.co/facebook/seamless-m4t-v2-large

| Category / Subcategory | Initial | Generated | Synthesised | Translated | | |
|---|---|---|---|---|---|---|
| | | | | De | Es | Ja |
| Contrastive Stress (General) | 200 | 199 | 183 | 87 | 76 | 97 |
| Relational/Descriptive Adjectives | 200 | 199 | 147 | 42 | 33 | 51 |
| Contrastive Stress (Noun-Phrase) | 200 | 199 | 124 | 37 | 36 | 39 |
| New/Given Information | 200 | 197 | 146 | 51 | 65 | 91 |
| Focus-sensitive Operators | 200 | 181 | 118 | 60 | 42 | 64 |
| **Sentence Stress** | 1000 | 975 | 718 | 277 | 252 | 342 |
| Complementizer/Parenthetical | 200 | 200 | 171 | 59 | 46 | 73 |
| VP/NP Attachment | 200 | 200 | 66 | 23 | 18 | 20 |
| Modifier Scope | 200 | 200 | 200 | 83 | 107 | 81 |
| Restrictive/Nonrestrictive | 200 | 199 | 177 | 65 | 82 | 40 |
| Direct/Indirect | 200 | 198 | 154 | 41 | 25 | 70 |
| Phrasal Verbs | 42 | 42 | 17 | 5 | 1 | 5 |
| **Prosodic Breaks** | 1042 | 1039 | 785 | 276 | 279 | 289 |
| **Intonation Patterns** | 300 | 263 | 174 | 173 | 173 | 173 |
| Sad-Happy | 200 | 200 | 1 | 1 | 1 | 1 |
| Neutral-Angry | 200 | 199 | 185 | 123 | 111 | 119 |
| Neutral-Happy | 200 | 198 | 161 | 81 | 97 | 81 |
| Disgust-Angry | 200 | 198 | 18 | 4 | 5 | 3 |
| Disgust-Sad | 200 | 198 | - | - | - | - |
| Neutral-Surprised | 200 | 198 | 43 | 33 | 35 | 30 |
| Disgust-Neutral | 200 | 197 | 7 | 2 | 5 | 5 |
| Happy-Angry | 200 | 197 | 138 | 50 | 65 | 72 |
| Sad-Surprised | 200 | 197 | 3 | 2 | 2 | 2 |
| Sad-Neutral | 200 | 196 | 4 | 3 | 2 | 2 |
| Sad-Angry | 200 | 196 | 5 | 1 | 4 | 4 |
| Disgust-Surprised | 200 | 196 | 4 | 2 | 2 | 1 |
| Disgust-Happy | 200 | 195 | 10 | 5 | 7 | 6 |
| Happy-Surprised | 200 | 195 | 52 | 34 | 27 | 21 |
| Angry-Surprised | 200 | 193 | 68 | 32 | 34 | 30 |
| **Emotional Prosody** | 3000 | 2953 | 699 | 433 | 418 | 377 |
| **Politeness** | 400 | 375 | 387 | 212 | 193 | 206 |
| **Total** | 5742 | 5605 | 2763 | 1311 | 1294 | 1386 |

Table 9: Number of Examples by Category and Subcategory

same model is MT mode.

2. XLS-R (Babu et al., 2021) is a multilingual E2E S2TT model that is based on a multilingual WAV2VEC 2.0 (Baevski et al., 2020) trained with self-supervised learning on a large speech corpus on 128 languages. For S2TT, the encoder is coupled with the decoder from MBART50 (Tang et al., 2020), and finetuned on paired speech-translation data. We use the following versions that are finetuned on English-to-15 on CoVoST2 (Wang et al., 2021): 300M[19], 1B[20], and 2B[21].

3. ZEROSWOT is a zero-shot E2E S2TT model that softly connects a WAV2VEC 2.0 encoder and an NLLB model, by compressing the speech representation into subword units and Optimal Transport (Peyré and Cuturi, 2019) alignment, using only ASR data. The versions used here are based on NLLB that were finetuned on the text data of CoVoST2, and the ZEROSWOT model was trained on Com-

---

[19]hf.co/facebook/wav2vec2-xls-r-300m-en-to-15
[20]hf.co/facebook/wav2vec2-xls-r-1b-en-to-15
[21]hf.co/facebook/wav2vec2-xls-r-2b-en-to-15

| | | |
|---|---|---|
| | Example 1: *This will only take a minute.* | |
| A | (neutral) | Das dauert nur eine Minute. (This will only take a minute.) |
| B | (angry) | Das dauert nur eine Minute, also machen Sie keinen Aufstand. (This will only take a minute so don't make a fuzz about it.) |
| | Example 2: *Our case was dismissed.* | |
| A | (neutral) | Unser Fall wurde abgewiesen. (Our case was dismissed.) |
| B | (sad) | Unser Fall wurde abgewiesen und das macht mich fassungslos. (Our case was dismissed which is just perplexing.) |

Table 10: Examples of overly explanatory translations proposed by GPT-4.

monVoice (Ardila et al., 2020). The MEDIUM version[22] has 1B parameters and the LARGE version[23] has 1.7B parameters.

4. SALMONN (Tang et al., 2024) is a general-purpose audio LLM that is capable of several speech- and audio-related tasks, including S2TT. It is build on top of the Vicuna LLM (Peng et al., 2023), and uses two encoders, one from WHISPER and one from BEATs (Chen et al., 2023). The concatenated output representations from the two encoders are processed by a Q-former (Li et al., 2023) and fed to the LLM which is finetuned with LoRA (Hu et al., 2022). There is a 7B version[24] and a 13B version[25]. To translate speech into a target language we use the recommended prompt from the paper: "Listen to the speech and translate it into {Target Language}".

5. WHISPER & NLLB is an AED-based cascade. WHISPER (Radford et al., 2022) is an encoder-decoder ASR and many-to-en S2TT model. We use three different versions for this casdade, namely the WHISPER-MEDIUM[26], the WHISPER-LARGE[27], and the latest v3

large version[28]. We primarily present results with the WHISPER-LARGE-V3, but since it was also used for filtering we also discuss v1 in order to avoid biasing our results. NLLB (NLLB Team, 2022) is a massively multilingual many-to-many MT model with access to 200 languages. We used the two distilled versions from the 54B MoE model, namely the distilled-600M[29] and the distilled-1.3B[30], as well as the 3.3B model[31]. We evaluated all possible combinations, thus having 9 cascade variants with these models.

6. CTC & NLLB is a CTC-based cascade. We use three different CTC encoders for the cascades. The first one is the Large version (300M) of WAV2VEC 2.0[32] which is finetuned on Libri-Light (Kahn et al., 2020) and Librispeech (Panayotov et al., 2015), additionally using self-training (Xu et al., 2020). The second is the Large version (300M) of HU-BERT[33] (Hsu et al., 2021), finetuned on Librispeech. The third is also based on HUBERT, more specifically to the XL version[34] with 1B parameters. We use the same three versions of NLLB, as we did for the AED-based cascade, thus having in total 9 variants of the CTC-based cascade.

## F Supplementary Results

In Figure 5 we present the Spearman rank correlation for the four contrastive metrics and the standard evaluation metric XCOMET. They were computed by evaluating all 31 models (§E) for all 3 language pairs, thus having a total of 93 observations.

In Table 11 we present the global contrastive quality scores for all 31 S2TT models for the 3 language pairs, which were used for the analysis of Figure 2 in §5 of the main text.

In Figures 6 and 7 we present the comparisons of the 4 models for Spanish and Japanese, similar to what we did in Figure 3 for German in the main text. In general, the findings and observations here coincide with those for German.

---

[22]hf.co/johntsi/ZeroSwot-Medium-cv-covost2-en-to-15
[23]hf.co/johntsi/ZeroSwot-Large-cv-covost2-en-to-15
[24]hf.co/tsinghua-ee/SALMONN-7B
[25]hf.co/tsinghua-ee/SALMONN
[26]hf.co/openai/whisper-medium
[27]hf.co/openai/whisper-large

[28]hf.co/openai/whisper-large-v3
[29]hf.co/facebook/nllb-200-distilled-600M
[30]hf.co/facebook/nllb-200-distilled-1.3B
[31]hf.co/facebook/nllb-200-3.3B
[32]hf.co/facebook/wav2vec2-large-960h-lv60-self
[33]hf.co/facebook/hubert-large-ls960-ft
[34]hf.co/facebook/hubert-xlarge-ls960-ft

| Model | Model Type | Model Size (B) | Contrastive Quality (Global) | | | |
|---|---|---|---|---|---|---|
| | | | En-De | En-Es | En-Ja | Average |
| SEAMLESSM4T-v1-MEDIUM | E2E | 1.2 | 9.1 | 10.4 | 10.1 | 9.9 |
| SEAMLESSM4T-v1-LARGE | E2E | 2.3 | 7.8 | 8.7 | 8.2 | 8.2 |
| SEAMLESSM4T-v2-LARGE | E2E | 2.3 | 14.5 | 11.0 | 13.9 | 13.1 |
| XLS-R 300M | E2E | 0.3 | 10.2 | 9.3 | 9.5 | 9.6 |
| XLS-R 1B | E2E | 1.0 | 9.3 | 8.3 | 8.1 | 8.6 |
| XLS-R 2B | E2E | 2.0 | 7.3 | 8.4 | 7.0 | 7.6 |
| ZEROSWOT-MEDIUM | E2E | 0.9 | 8.9 | 7.5 | 7.5 | 8.0 |
| ZEROSWOT-LARGE | E2E | 0.9 | 8.7 | 7.7 | 7.9 | 8.1 |
| SALMONN-7B | E2E | 7.0 | 12.7 | 9.4 | 12.9 | 11.7 |
| SALMONN-13B | E2E | 13.0 | **15.9** | **12.3** | **16.1** | **14.8** |
| SEAMLESSM4T-v1-MEDIUM | Cascade-AED | 2.4 | 7.3 | 7.6 | 6.9 | 7.2 |
| SEAMLESSM4T-v1-LARGE | Cascade-AED | 4.6 | 6.5 | 5.8 | 4.3 | 5.5 |
| SEAMLESSM4T-v2-LARGE | Cascade-AED | 4.6 | 10.5 | 7.7 | 8.7 | 8.9 |
| WHISPER-v1-MEDIUM & NLLB-600M | Cascade-AED | 1.4 | 5.7 | 5.0 | 5.4 | 5.4 |
| WHISPER-v1-MEDIUM & NLLB-1.3B | Cascade-AED | 2.1 | 6.2 | 5.0 | 5.8 | 5.6 |
| WHISPER-v1-MEDIUM & NLLB-3.3B | Cascade-AED | 4.1 | 6.6 | 5.1 | 5.4 | 5.7 |
| WHISPER-v1-LARGE & NLLB-600M | Cascade-AED | 2.2 | 5.9 | 4.9 | 6.1 | 5.6 |
| WHISPER-v1-LARGE & NLLB-1.3B | Cascade-AED | 2.9 | 6.0 | 4.7 | 5.8 | 5.5 |
| WHISPER-v1-LARGE & NLLB-3.3B | Cascade-AED | 4.9 | 6.3 | 4.6 | 5.6 | 5.5 |
| WHISPER-v3-LARGE & NLLB-600M | Cascade-AED | 2.2 | 5.3 | 4.6 | 6.5 | 5.5 |
| WHISPER-v3-LARGE & NLLB-1.3B | Cascade-AED | 2.9 | 5.3 | 4.7 | 5.4 | 5.1 |
| WHISPER-v3-LARGE & NLLB-3.3B | Cascade-AED | 4.9 | 5.5 | 4.8 | 5.3 | 5.2 |
| WAV2VEC 2.0 & NLLB-600M | Cascade-CTC | 0.9 | 1.4 | 1.3 | 2.0 | 1.5 |
| WAV2VEC 2.0 & NLLB-1.3B | Cascade-CTC | 1.6 | 1.7 | 1.0 | 1.4 | 1.3 |
| WAV2VEC 2.0 & NLLB-3.3B | Cascade-CTC | 3.6 | 1.6 | 0.9 | 1.7 | 1.4 |
| HUBERT & NLLB-600M | Cascade-CTC | 0.9 | 3.2 | 2.7 | 2.5 | 2.8 |
| HUBERT & NLLB-1.3B | Cascade-CTC | 1.6 | 2.2 | 2.4 | 2.9 | 2.5 |
| HUBERT & NLLB-3.3B | Cascade-CTC | 3.6 | 2.7 | 2.6 | 2.9 | 2.7 |
| HUBERT-XL & NLLB-600M | Cascade-CTC | 1.6 | 2.4 | 2.9 | 3.0 | 2.8 |
| HUBERT-XL & NLLB-1.3B | Cascade-CTC | 2.3 | 3.5 | 1.7 | 2.7 | 2.6 |
| HUBERT-XL & NLLB-3.3B | Cascade-CTC | 4.3 | 2.6 | 2.4 | 2.5 | 2.5 |

Table 11: Contrastive Quality (Global) scores for English-German, English-Spanish, and English-Japanese, including their averages.

Figure 5: Correlation Matrix of the metrics across all language pairs and models.

Figure 6: Upper: Model performance per category (En-Es). Lower: Model performance comparisons (En-Es), (a): SALMONN-13B vs. SEAMLESSM4T-v2-LARGE, (b) SEAMLESSM4T-v2-LARGE(E2E) vs. SEAMLESSM4T-v2-LARGE(cascade), (c) SEAMLESSM4T-v2-LARGE(cascade) vs. WHISPER-v3-LARGE/NLLB-3.3B.



Figure 7: Upper: Model performance per category (En-Ja). Lower: Model performance comparisons (En-Ja), (a): SALMONN-13B vs. SEAMLESSM4T-v2-LARGE, (b) SEAMLESSM4T-v2-LARGE(E2E) vs. SEAMLESSM4T-v2-LARGE(cascade), (c) SEAMLESSM4T-v2-LARGE(cascade) vs. WHISPER-v3-LARGE/NLLB-3.3B.

# Cultural Adaptation of Menus: A Fine-Grained Approach

**Zhonghe Zhang, Xiaoyu He, Vivek Iyer, Alexandra Birch**
University of Edinburgh
zhonghe.zhang@hotmail.com, claire.xiaoyu.he@gmail.com
{vivek.iyer, a.birch}@ed.ac.uk

## Abstract

Machine Translation of Culture-Specific Items (CSIs) poses significant challenges. Recent work on CSI translation has shown some success using Large Language Models (LLMs) to adapt to different languages and cultures; however, a deeper analysis is needed to examine the benefits and pitfalls of each method. In this paper, we introduce the ChineseMenuCSI dataset, the largest for Chinese-English menu corpora, annotated with CSI vs Non-CSI labels and a fine-grained test set. We define three levels of CSI figurativeness for a more nuanced analysis and develop a novel methodology for automatic CSI identification, which outperforms GPT-based prompts in most categories. Importantly, we are the first to integrate human translation theories into LLM-driven translation processes, significantly improving translation accuracy, with COMET scores increasing by up to 7 points. The code and dataset are available at https://github.com/Henry8772/ChineseMenuCSI.

## 1 Introduction

Translating restaurant menus is a challenging, non-literal translation task. Unlike other texts, dish names are not merely lists of ingredients and culinary methods; they are short (Pellatt and Liu, 2010) and culturally rich expressions that require an understanding of cultural traditions (Amenador and Wang, 2022), symbolism (Lam et al., 2018), and local nuances. This complexity is compounded by LLMs and Neural Machine Translation (NMT) systems that often lack the cultural awareness necessary to accurately understand these nuances (Liu et al., 2024; Naous et al., 2023; Tao et al., 2024). This results in mistranslations that can confuse and mislead the target audience (Garcea et al., 2023; Gallo et al., 2021), such as in Figure 1.

A key challenge in menu translation lies in the handling of Culture-Specific Items (CSIs), defined



Figure 1: CSI translation errors by Google Translate and ChatGPT 3.5 in translating Chinese culinary terms.

as "concepts that are specific to a particular language or group" (Aixelá, 1996). For example, the literal translation of the Chinese dish 蚂蚁上树 is "Ants Climbing a Tree" – but this is actually a figurative Chinese expression that should be translated in English as "Sauteed Vermicelli with Minced Pork". The Chinese name creatively expresses the idea that pork resembles ants, while vermicelli represents tree branches. Existing machine translation systems, trained on plain, sentence-level translations, fail to capture these cultural subtleties and generate literal translations (Figure 1).

However, there has been little work in NLP exploring CSI translations in-depth, particularly focusing on how the translation outputs generated by neural models should be improved. There has been foundational work on improving translations of CSIs by LLMs through enhanced prompting strategies (Yao et al., 2024). Simultaneously, there has also been work on adapting CSIs (Peskov et al., 2021; Cao et al., 2024; Singh et al., 2024), but their focus has been on adapting culture-specific named entities. In this work, we seek to go beyond entities, and approach the translation of figurative language imbued with cultural nuance, as exemplified in Figure 1 – which is quite underexplored.

In linguistics, cultural translation theories have been developed and widely adopted by human theorists and translators over decades. We aim to improve MT of CSIs by bringing the wisdom of Translation Theory research to modern NLP mod-

els like LLMs. Our approach improves the identification and translation of figurative and culturally nuanced CSIs. Unlike previous CSI identification methods, our method does not depend on parallel corpora or extensive knowledge graphs (Yao et al., 2024; Han et al., 2023) – but at the same time, we also show how recipes can be *optionally* leveraged as a source of external knowledge to enhance performance even further. We also propose a novel CSI taxonomy for Chinese-English, that allows for a detailed analysis of figurative and culturally nuanced language and the translation challenges therein. We evaluate our proposed methods using a large dataset of Chinese dish names, ensuring robust and reliable results.

Our key contributions are as follows:

1. We introduce *ChineseMenuCSI*, a fine-grained dataset of 4,275 bilingual Chinese-English restaurant menu entries from UK Chinese restaurants. The dataset is categorised into CSI and Non-CSI entities, with 480 entries further bifurcated into specific CSI categories, enabling an in-depth analysis of CSI translation efficacy of LLMs and NMT systems.

2. We propose novel techniques for identifying CSIs, grounded in human translation theory. These techniques match or outperform current GPT-based prompts – all without needing external knowledge graphs or parallel corpora.

3. Lastly, we show how external knowledge, in the form of recipes, can add to the benefits of translation strategies and enhance CSI translation performance further - achieving significant improvements in COMET scores, with gains of +3 to +7 points across CSI categories.

The code and datasets are available at https://github.com/Henry8772/ChineseMenuCSI.

## 2 Related Work

### 2.1 CSIs in Translation Studies

Aixelá (1996) was among the first scholars to introduce the term "culture-specific items" (CSIs) to refer to elements in texts that are unique and significant in a specific culture. CSIs may include objects, classification systems, or measurement tools common in the source culture but foreign to the target culture. Additionally, CSIs can encompass transcriptions of opinions/habits specific to a culture,

which are often reflected in the language structure, style, and content.

Culture is closely related to understanding and translating CSIs, as Aixelá (1996) highlighted. In the 1960s, Nida and Taber (2003) introduced the concepts of formal and dynamic equivalence in translation to distinguish between structurally accurate and fluency-focused approaches to translation. These concepts have laid the foundation for subsequent translation theories, including those related to cultural translation. Expanding on this, Newmark (1988) proposed a set of robust strategies for translating cultural elements, which have been particularly influential in translating Chinese culinary CSIs, as noted by Amenador and Wang (2022).

According to Newmark (1988), adaptation uses a recognized equivalent between two cultures. This strategy has been explored by Pellatt and Liu (2010) on Chinese menu translation. Newmark (1988) proposed three equivalent strategies for translation: cultural, functional, and descriptive – which we introduce later to improve LLM translation performance in §5.2.

Neutralisation is another translation strategy related to CSI translation. As proposed by Chou et al. (2016), on the continuum between foreignisation (focusing on source culture) and domestication (focusing on target culture), there are intermediary approaches, including neutrality and neutralisation. For culture-specific text, neutralisation involves paraphrasing to convey the meaning of a CSI. After analysing the translations of Chinese dish names into English, Amenador and Wang (2022) found that neutralisation, by substituting the source text element with a more or less detailed explanation of its meaning, is the most commonly used translation strategy by human translators for translating Chinese dish names into English.

In this paper, we use these conventional translation strategies employed by human translators as instructions in zero-shot prompts, to enhance CSI translation quality of LLMs (in §5.2).

### 2.2 Culture-Aware NMT

Despite the early successes of NMT (Bahdanau et al., 2015; Sutskever et al., 2014), translation of culture-specific texts has remained a daunting task. In addition to the challenge of translating rarer words and adapting to under-resourced domains (Koehn and Knowles, 2017), CSIs are deeply intertwined with cultures (Hershcovich et al., 2022; Liebling et al., 2022; Yao et al., 2024) – something

even the most capable neural models of today fail to grasp, particularly for non-Western cultures (Masoud et al., 2023; AlKhamissi et al., 2024; Nayak et al., 2024).

While there have been related works on domain-specific translation, including terminology translation (Dinu et al., 2019), disambiguation (Iyer et al., 2023a,b) and named entity translation (Hu et al., 2022), CSIs often lack direct equivalents in other languages, making translation complex and hard to understand cross-culturally (Yao et al., 2024).

Our approach uniquely combines translation studies with modern NLP techniques to identify and translate CSIs, resulting in more culturally sensitive and comprehensible translations.

### 2.3 Cultural Awareness and Adaptation in Large Language Models

In recent times, many works have shown that LLMs contain significant cultural biases against non-Western cultures (Cao et al., 2023; Liu et al., 2024; Masoud et al., 2023; Naous et al., 2023; Tao et al., 2024). In response, there has been a growing focus on improving cultural awareness in LLMs through prompt-engineering techniques (Wang et al., 2024; Tao et al., 2024) and fine-tuning on culture-specific data (Chan et al., 2023; Li et al., 2024a,b). Various tasks have been used to assess LLMs' cultural awareness, including tasks like culturally aware inference (Huang and Yang, 2023; Yao et al., 2024) and common sense reasoning on specific languages (Koto et al., 2024a,b).

Previous works on cultural awareness have primarily focused on understanding cultural norms in different languages rather than accurately translating culture-specific text. While well-explored in translation studies, cultural adaptation is rather understudied in NLP. Initial efforts in this direction have included adaptation of recipes (Cao et al., 2024) and localisation of named entities through adaptation (Peskov et al., 2021) or explicitation (Kementchedjhieva et al., 2020; Garcea et al., 2023; Han et al., 2023). Most similar to our work is that of Yao et al. (2024), who also released a CSI dataset covering 6 languages, on which they benchmark LLMs and NMT systems.

In contrast, our goal is to conduct a more fine-grained evaluation, given multiple CSI types in any given language. So, we leverage translation studies to create a dataset that classifies Chinese-English dishes into fine-grained categories, which we use for downstream evaluation, analysis and a detailed

ablation of our proposed techniques. While we focus on the Chinese-English pair and culinary domain in this work, our framework and proposed techniques are agnostic of language/domain, and are designed to be easily scalable.

## 3 ChineseMenuCSI Dataset

We introduce a new bilingual Chinese-English Restaurant Menu (ChineseMenuCSI) dataset consisting of 4,275 human-verified dish entries collected from restaurants in UK.

### 3.1 Data Collection

We develop a Selenium-based web crawler[1] to gather Chinese menu translations from UK restaurant websites. After manually reviewing 50 restaurants, we selected those with ratings above 3 out of 5 and average meal prices over £20, ensuring higher-quality menus not generated by commercial Machine Translation systems like Google Translate. These restaurants were sourced from TripAdvisor[2].

Additionally, we developed a heuristic menu parser capable of accurately extracting structural content from image-based menus. Details of this parser are provided in Appendix A.2.

### 3.2 CSI Taxonomy

Translating Chinese menu items into English presents unique challenges because the dish names contain non-descriptive, picturesque elements (Pellatt and Liu, 2010). Our initial data inspection revealed that CSIs within these dish names contribute varying degrees of complexity to the translation process, and carry differing levels of figurativeness brought by cultural and linguistic nuances. Inspired by translation theory literature that tends to categorise Chinese dishes into concrete and abstract categories (Lam et al., 2018), we develop an approach to categorise the Chinese dish names in our dataset into three groups based on the degree of figurativeness in each CSI.

**Category 1: Concrete CSIs (With a Low-level/No Figurative Meaning)**

**Definition:** The CSIs in this category have a minimal figurative meaning, often referring to tangible attributes like ingredients, colour, taste, container, processing method, and dish appearance. Readers

---

[1]Selenium: https://www.selenium.dev/
[2]TripAdvisor: https://www.tripadvisor.com

1260

can easily understand these dish names as the information is either shared between the source and target cultures or has widely used translations in the target culture.

**Example:** An example from the corpus is *"咕噜猪肉" (sweet and sour pork)*. The first two Chinese characters "咕噜" denote the Guangdong-style "sweet and sour" method, a culinary translation widely recognised outside the Chinese culture. The last two characters "猪肉" mean "pork", a culturally universal ingredient.

**Category 2: Creative CSIs (With Some Figurative Meaning)**

**Definition:** This category features dish names that blend concrete lexical terms with figurative meanings, creating inventive expressions that extend beyond literal definitions. Understanding these dishes necessitates integration of creative flair with concrete information, presenting challenges.

**Example:** *"水煮鱼"* (*Poached fish fillet with chilli oil and herb* or *Sichuan-style boiled fish*) originates from Sichuan, China. While the literal translation of the characters is "water-boiled fish", "water-boil" carries a creative description, representing the cooking state. This dish involves a Sichuan cooking style that uses hot chilli oil and Chinese herbs. "Poached fish fillet with chilli oil and herb" effectively describes the ingredients and cooking method, while "Sichuan-style boiled fish" adds cultural context by highlighting the dish's regional origin. Both are valid translations but different strategies are used.

**Category 3: Abstract CSIs (With a High-level of Figurative Meaning)**

**Definition:** This category encompasses dish names that exist beyond the realm of literal translation, and require in-depth cultural knowledge to understand. Crafted from metaphors, idioms, allegories, and other figurative language, these names disconnect from straightforward translations to engage in storytelling, aiming to convey broader narratives, evoke emotions, or reflect cultural heritage.

**Example:** *"佛跳墙"* (*Buddha Jumps Over the Wall* or *Steamed Abalone with Fish Maw in Chicken Broth*) metaphorically describes a dish so enticing that even a vegetarian and divine figure like Buddha would leap over a wall to taste it. Popular translations include "Buddha jumps over the wall"

as a direct translation, and "Steamed Abalone with Fish Maw in Chicken Broth" includes ingredients and cooking methods, reflecting the dish's cultural and culinary nuances.

### 3.3 Data Annotation

To annotate our data, we first seek to classify the data into CSI and non-CSI entities, and if it is a CSI, we want to categorise it into one of the above-listed groups. Given the dataset has as many as 4.3K CSIs, we approach the annotation process in two stages: a) in Stage 1, we conduct a broad, albeit rough, annotation of the entire dataset by two volunteer annotators, and b) we uniformly sample from the annotations of Stage 1 to ensure a fair distribution across categories, and conduct a more focused and rigorous annotation process using five volunteer translators – to create our fine-grained test set. We describe these stages in detail below:

**Stage 1 (Broad) Annotation:** Two annotators, who are postgraduate students and professional Chinese-English translators, reviewed and labelled all 4,275 entries. Both are native Chinese speakers proficient in English, ensuring high linguistic and cultural expertise. Firstly, we classify the entries into CSIs and non-CSIs. We use Cohen's kappa (Cohen, 1960) to measure agreement between annotators and obtained a high score of 0.91 - likely because the classification of CSIs and non-CSIs is mostly unambiguous. For the 187 entries without consensus, we invited a third annotator to label and assigned the final label using a majority vote.

For entries with CSI, the annotators further categorised the dish into one of the three CSI categories. The annotation results are reported in Table 1.

| Label | Category | Count |
|-------|--------------|-------|
| 0 | Non-CSIs | 2003 |
| 1 | Concrete CSIs | 1658 |
| 2 | Creative CSIs | 494 |
| 3 | Abstract CSIs | 120 |

Table 1: Distribution of menu items across CSI taxonomy in the ChineseMenuCSI dataset.

**Stage 2 (Focused) Annotation:** In Table 1, we note that Category 3 is the smallest, with only 120 items. To evenly balance our test set, we randomly sample 120 items from each category: 0 (Non-CSIs), 1 (Concrete CSIs), 2 (Creative CSIs), and 3 (Abstract CSIs), totalling 480 items. This subset was annotated by a larger and more diverse group of five annotators, who, like the first-stage

annotators, included professional Chinese-English translators and postgraduate students – all native Chinese speakers proficient in English.

For span annotation, we first segment the dish name into spans and phrases using Jieba[3]. Annotators then label the spans that correspond to CSI. To assess inter-annotator agreement, we use Fleiss' kappa (Fleiss, 1971) across two levels: CSI fine-grained categorisation and span-level CSI identification; the results are summarised in Table 2.

| Annotation | Kappa | Interpretation |
|---|---|---|
| CSI vs. Non-CSI | 0.91 | High |
| CSI Category | 0.63 | Substantial |
| CSI Identification | 0.70 | Substantial |

Table 2: Inter-annotator agreement scores for different annotation tasks.

For the fine-grained CSI categorisation, we exclude items that do not attain majority consensus (at least 3 out of 5 annotators in agreement), resulting in a kappa score of 0.63. This score falls within the range of substantial agreement (0.6-0.8) but is lower than the CSI vs. Non-CSI score due to the subjective nature of the fine-grained categorisation. This level of agreement is comparable to ranges reported in related work (Huang and Yang, 2023; Soderstrom et al., 2021). Lastly, for CSI identification at the span-level, i.e. within a given dish name, the kappa score is 0.70 - which indicates substantial agreement as well.

## 4 CSI Automatic Identification

To accurately translate CSIs, it is essential to first identify which parts of the text comprise CSIs. Previous studies have approached this challenge in different ways. Han et al. (2023) focus on implicit detection and use a relative distance of terms in Wikidata, but it does not include all Chinese dish CSI. Yao et al. (2024) rely on parallel corpora with entity-linking to find CSIs; however, the approach is infeasible for online MT, where we need to identify CSIs beforehand to produce translations from the monolingual source text.

Inspired by these methods, we propose a method called **Combined CSI Identification**, that uses a combination of three checking criteria for CSI identification, and classifies CSI if at least two of the following three checks are met. The checks are Round-trip Translation (RTT), Cultural Uniqueness (CU), and Historical Significance (HS).

---

[3]Jieba: https://github.com/fxsjy/jieba

### 4.1 Round-trip Translation (RTT)

Since CSIs are defined as terms unique to a specific language or culture (Álvarez and Vidal, 1996), based on the assumption that they do not have corresponding translations in the target language, we propose using round-trip translation (RTT) as one of the identifying criteria.

1. **Initial Translation**: Translate the Chinese dish name to English using Google Translate.

2. **RTT Translation**: Translate the English version back to Chinese using DeepL Translate and split it into Jieba-segmented words. Using different translation systems for RTT proved most effective in identifying CSIs.

3. **Identification**: Subtract the segmented words in the RTT from those in the original text. The remaining words are potential CSIs.

$$\text{CSIs} = \text{Original Words} - \text{RTT Words}$$

Using Jieba's cut-for-search module, which returns all words and phrases, a phrase is considered CSI only if all of its words are omitted in the RTT, and not otherwise.

This method has its limitations, for example, it could also: a) return words that are not CSIs and are just difficult to translate, and b) miss CSIs that have literal translations. We find in §6.3 that it performs strongly in identifying CSIs in most cases.

### 4.2 Cultural Uniqueness (CU)

According to Newmark (1988), "unfindable" words are often less frequently seen within a language. Words with cultural and historical references can be deeply embedded in a specific culture or history, making them rare or unfamiliar to outsiders.

We use Jieba to segment words in the Chinese-MenuCSI dataset, then measure each word's frequency and calculate its inverse frequency. A cut-off at the 95th percentile of these inverse frequencies is set based on a manual review of 100 words. Words above this cut-off are marked as potential CSI. Words not previously seen are given an inverse frequency of 1, indicating they are potential CSI. No smoothing techniques are applied, as inverse frequency is used against a fixed threshold rather than for probability calculations.

### 4.3 Historical Significance (HS)

Chinese food names and CSIs often contain historical narratives such as historical events, figures and periods (Lam et al., 2018; Amenador and Wang, 2022). To identify these, we use the Wikipedia API to search for individual words or entire dish names. If a word's Wikipedia page includes a "History" section, it is considered a potential CSI. The words appearing 30 times or more, such as "chicken" or "sauce" are excluded as generic terms.

## 5 CSI Translation

Having identified CSIs, we propose prompting strategies to improve CSI translations. Our strategies fall into two categories: Recipe-based Translation and Translation Studies-inspired Prompting.

### 5.1 Recipe-based Translation

We explore using recipe information to improve the translation of CSI dish names. By incorporating the most relevant recipe as external knowledge, we experimented with two zero-shot prompt strategies: **Default Recipe** prompting and **Recipe + Explain-then-Translation** prompting (Figure 2).

**Recipe Retrieval Pipeline**  Given over 50% of the CSI dish entries in our test set lack detailed publicly available descriptions from sources like Wikipedia, we use the Xiachufang recipe database (Liu et al., 2022) – which contains approximately 1.4 million Chinese monolingual recipes – to retrieve recipes and enhance translation accuracy.

We take inspiration from Translation Studies research, which emphasizes the importance of cooking methods and ingredients in translating Chinese dish names (Amenador and Wang, 2022).

Our retrieval pipeline involves two key stages:

1. **Query and document Construction:** The query is the full Chinese dish name, with the CSI span from previous annotations. We concatenated each recipe name and instructions into a single recipe document.

2. **Filtering and Ranking Recipes:** To filter and rank the recipes, we employ the BM25 algorithm (Robertson et al., 1995), which assigns a score to each word in the Chinese recipe document based on its term frequency and inverse document frequency. For each word in the recipe document, the score is enhanced by applying a weighting factor when the word matches either the dish name (weight = 5) or the CSI span (weight = 3), with an additional multiplier of 3 applied to words within the dish name to prioritize their importance. If there is no exact match for the dish name, the process shifts focus to matching with the CSI span. Additionally, we apply a length penalty to the final score, adjusting it based on the difference between the recipe's length and the average length of all recipes. We select the top-ranked recipe as the final output.

**Prompt Strategy: CSI Recipe**  In this prompting strategy, we use the most relevant recipe returned from the aforementioned search pipeline to aid CSI translation. We provide the name and cooking instructions of the closest-matching recipe while noting that it might not correspond exactly to the given dish but is beneficial as external knowledge since it contains the CSIs to translate.

**Prompt Strategy: CSI Explain-then-Translate**  Inspired by Chain of Thought (CoT) prompting (Wei et al., 2024) and Self-Explanation (Yao et al., 2024), we formulate another prompting strategy that first asks the LLMs to explain the meaning of the CSIs described in the recipe and then generate the translation for the dish. The motivation is to help the LLMs conduct advanced reasoning on the recipe instructions, such as interpreting dish names with CSIs not explicitly defined in the recipe. For example, a recipe might instruct to "cut it first, then stir fry" or note that "it can be very spicy" without explaining the CSI. The LLM's task is to infer the meaning of the CSIs based on the recipe's instructions. The prompt is shown in Figure 2.

### 5.2 Translation Studies-inspired Prompting

Unlike the conventional prompt engineering used in related work, our second set of prompting strategies differs in that they incorporate *human translation strategies*, inspired by the rich literature in Translation Studies, directly into the design of the prompt. We provide the prompt template for both of these strategies in Figure 2 and complete prompt in Appendix 3 and 4.

**Prompt Strategy: Equivalents**  Using the recipe in §5.1 as external knowledge, we ask the LLMs to produce three translations, each based on a translation strategy (i.e. cultural, functional and descriptive) and **select the best translation**.
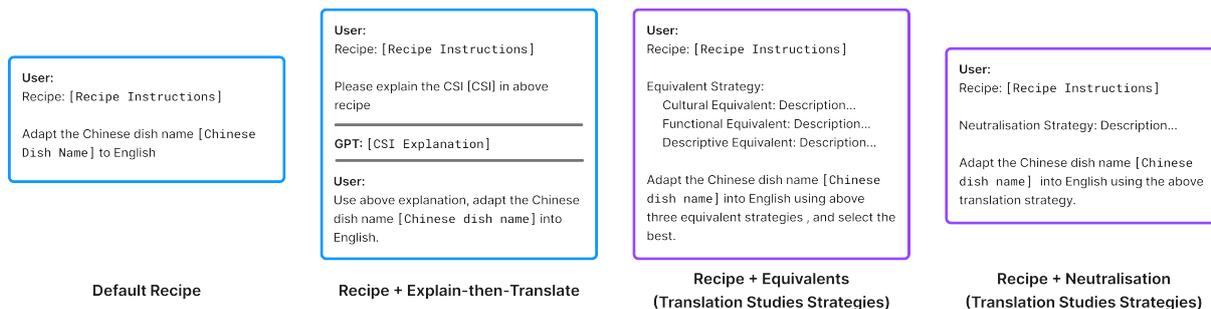
Figure 2: Four adaptation prompt strategies

These translation strategies are inspired by New-mark (1988)'s theories on equivalent translation. We define the equivalent strategies below, providing examples for the reader's understanding:

1. **Cultural Equivalent:** Replacing a CSI in the source text with a term that is culturally relevant and functionally equivalent in the target culture. This strategy aims to evoke the same response in the target audience. (i.e. translating "粽子" as "tamale" in Spanish – given both are traditional wrapped food items made with a starchy substance and fillings, albeit from different cultures.)

2. **Functional Equivalent:** This strategy focuses on the function or purpose of the item (i.e., translating "粽子" as "rice dumpling" to convey the idea of a food made of rice)

3. **Descriptive Equivalent:** Providing a detailed description or explanation of the CSI to convey its meaning and significance. This approach is useful when the CSI is essential for understanding the text but has no equivalent in the target language (i.e. translating "粽子" as "a traditional Chinese sticky rice dumpling wrapped in bamboo leaves")

**Prompt Strategy: Neutralisation** Another human translation strategy we used in the prompting experiments is neutralisation. Again, we provide recipe information as external knowledge as in the previous strategy and incorporate an explanation of the neutralisation strategy to guide the LLM translation of dish names.

**Neutralisation:** Using culturally neutral language to describe or explain a cultural word, phrase, or rhetorical expression from the source text. It answers the question, "What is this?" (Amenador and Wang, 2022) by adding information such as ingredients, culinary methods and key characteristics. Compared with the descriptive equivalent

strategy, the neutralisation strategy we used for prompt design confines the information used in the translations to ingredients, culinary methods and key characteristics (i.e. translating "粽子" as "sticky rice wrapped in bamboo leaves").

## 6 Results and Analysis

Our experiments include five parts: 1. Baseline evaluation of MT performance using three models on the ChineseMenuCSI dataset (§6.2); 2. Assessment of CSI span identification accuracy (§6.3); 3. Exploration of main adaptation strategies (§6.4); 4. Exploration of individual equivalents translation strategies for enhancing CSI translation (§6.5); 5. Human evaluation of translation quality on a subset of the dataset (§6.6).

### 6.1 Experimental Settings

We evaluate the effectiveness of LLM translations for CSIs by comparing various prompting strategies across two SOTA LLMs — GPT-3.5 (gpt-3.5-turbo-0125[4]) and the advanced GPT-4o (gpt-4o-2024-05-13[5]) — against the robust commercial MT system, Google Translate. This approach allows us to assess the strengths of LLM prompting versus a widely used commercial MT.

### 6.2 Evaluation of CSIs vs Non-CSIs

We use the Stage 1 annotated version of the ChineseMenuCSI dataset (§3.3). This version has a high inter-annotator agreement of 0.91 for the CSI vs non-CSI classification task, with conflicts further resolved using a third annotator. We compare the MT performance using the COMET version wmt22-comet-da (Rei et al., 2022).

While GPT-4o yields major improvements in the CSI translation performance, GPT-3.5 and GPT-4o

---

[4]GPT-3.5:https://platform.openai.com/docs/models/gpt-3-5-turbo

[5]GPT-4o:https://platform.openai.com/docs/models/gpt-4o

1264

| Method | Non-CSIs | CSIs |
|---|---|---|
| Google Translate | **74.08** | 64.48 |
| GPT-3.5 | 72.88 | 64.34 |
| GPT-4o | 73.67 | **65.97** |

Table 3: Comparison of COMET scores across different translation systems for CSI and Non-CSI menu items. The prompt used: "Translate the [Chinese dish name] into English".

| | Method | Precision | Recall | F1 Score |
|---|---|---|---|---|
| **CSI-1: Concrete CSIs** | | | | |
| Ours | Combined | 64.9 | 34.0 | 44.7 |
| | RTT | 58.8 | 28.4 | 38.3 |
| | CU | 38.3 | 65.3 | 48.2 |
| | HS | **86.3** | 31.2 | 45.8 |
| GPT | 3.5 | 32.4 | **80.4** | 46.2 |
| | 4o | 37.9 | 67.1 | **48.5** |
| **CSI-2: Creatives CSIs** | | | | |
| Ours | Combined | 66.1 | 53.1 | 58.9 |
| | RTT | 63.4 | 60.1 | **61.7** |
| | CU | 35.4 | 70.4 | 47.1 |
| | HS | **68.6** | 16.4 | 26.5 |
| GPT | 3.5 | 34.1 | **82.0** | 48.2 |
| | 4o | 40.6 | 73.9 | 52.4 |
| **CSI-3: Abstract CSIs** | | | | |
| Ours | Combined | 81.4 | 68.6 | 74.4 |
| | RTT | **81.7** | 73.6 | **77.4** |
| | CU | 43.9 | 88.4 | 58.6 |
| | HS | 80.0 | 9.9 | 17.6 |
| GPT | 3.5 | 50.4 | **94.3** | 65.7 |
| | 4o | 59.1 | 78.9 | 67.6 |

Table 4: Evaluation of CSI span identification accuracy by CSI category: precision, recall, and F1 scores.

show worse scores for non-CSIs than Google Translate. These results suggest that Google Translate is particularly strong at translating straightforward, culturally neutral content, likely due to its extensive and diverse training dataset which prioritizes general language accuracy over cultural nuances.

Interestingly, despite Google Translate's general strength in multilingual tasks (Zhu et al., 2024), GPT-4o shows better handling of CSIs, highlighting the benefits of pretraining at scale on diverse corpora from many cultures, as opposed to NMT systems that are typically trained on narrow-domain sentence-level parallel corpora.

### 6.3 Evaluation of CSI Span Identification

We further assess the capability of different methods to *pinpoint specific CSI spans* – a task we call **CSI Span Identification** – within dish names. This fine-grained analysis is crucial for understanding the elements that contribute to cultural specificity and translation complexity.

Table 4 shows GPT 4o as the best performer in

CSI-1 and RTT as the best in other categories. This is likely due to RTT's strength in identifying CSIs, which are often figurative and lack general interpretations. However, the combined metrics only improved RTT's performance in CSI-1, likely because of low recall in HS. Moreover, CU underperforms in CSI-2, as those CSIs typically involve figurative messages in creative combinations of frequent words, which cannot be captured by frequency.

The combined method does not outperform the individual highest method as it requires majority agreement, where a single correct check is insufficient. HS shows high precision but low recall, likely because CSIs often have historical backgrounds, though the inverse is not always true.

### 6.4 Evaluation of Main Adaptation Strategies

This section examines whether external knowledge, such as recipes, can enhance MT of CSI-rich dish names. Four prompting strategies were tested: a) Default Recipe prompting, b) Recipe + Explain-then-Translate (EtT) prompting, c) Equivalents and d) Neutralisation. For the latter two, which are translation strategies, we try baselines with and without incorporation of recipes, to ablate the dependency on external knowledge. In Table 5, we see that while all our proposed methods yield overall improvements in performance, translation strategy-based methods – that do not involve any external knowledge – yield the largest gains, of up to +4.8 COMET points! Moreover, when recipes are added to these strategies, the gap widens even further, with the maximum gain reaching *as high as +7.87 COMET points* for our best performing *Recipe + Equivalents* strategy. The second trend we note is that the largest gains come from the more complex CSI-2 and CSI-3 categories, indicating the efficacy of our translation theory-inspired methods for translating highly culturally nuanced text. Finally, while the more advanced model GPT-4o naturally yields the largest improvements, we note that we get pretty good results with the far cheaper GPT-3.5 model too, indicating that our methods could be used quite economically.

Revisiting the relatively lower improvements in CSI-1 by examining GPT-generated translations, we find that the LLMs sometimes focus on irrelevant details in the provided recipes. In CSI-1, which involves shorter CSI terms, finding an exact match for dish names is harder, forcing the inclusion of noise in the recipe. For instance, the term "咕嚕" *(Sweet and sour)* applies to various dishes

| | GPT-3.5 | | | | GPT-4o | | | |
|---|---|---|---|---|---|---|---|---|
| | CSI-1 | CSI-2 | CSI-3 | Overall | CSI-1 | CSI-2 | CSI-3 | Overall |
| Baseline | 62.68 | 55.38 | 43.92 | 53.33 | 62.68 | 55.38 | 43.92 | 53.33 |
| *Recipe-based Translation* | | | | | | | | |
| Recipe | +0.16 | -0.90 | +3.44 | +0.50 | -0.08 | -3.02 | +3.49 | +1.93 |
| Recipe + EtT | +1.13 | -1.33 | +4.92 | +1.04 | +1.10 | +1.61 | +4.87 | +2.16 |
| *Translation Studies Prompting* | | | | | | | | |
| Neutralisation | +0.74 | +1.15 | +3.62 | +1.56 | +0.46 | **+4.84** | +4.29 | +3.02 |
| Equivalents | **+1.44** | +3.24 | +2.52 | +2.38 | **+2.34** | +3.89 | +0.94 | +2.62 |
| *Recipe + Translation Studies Prompting* | | | | | | | | |
| Recipe + Neutralisation | -1.29 | +1.15 | **+7.72** | +1.71 | +0.09 | +3.47 | +4.25 | +2.34 |
| Recipe + Equivalents | +0.95 | **+3.85** | +3.01 | **+2.54** | +1.80 | +3.24 | **+7.87** | **+3.74** |

Table 5: COMET score comparisons for GPT-3.5 and GPT-4o using various translation strategies across CSI categories. The overall score is calculated as the average of CSI-1, CSI-2, and CSI-3 scores for each method.

| | GPT-3.5 | | | | GPT-4o | | | |
|---|---|---|---|---|---|---|---|---|
| | CSI-1 | CSI-2 | CSI-3 | Overall | CSI-1 | CSI-2 | CSI-3 | Overall |
| Baseline | 62.68 | 55.38 | 43.92 | 53.33 | 63.43 | 54.50 | 47.50 | 55.14 |
| *Equivalents Strategy Prompting* | | | | | | | | |
| Cultural | -0.06 | -0.86 | +0.99 | +0.02 | **+0.91** | +0.90 | -2.77 | -0.32 |
| Descriptive | -6.73 | -1.90 | +0.93 | -2.57 | -3.83 | +2.62 | +2.10 | +0.96 |
| Functional | +0.54 | +2.69 | +0.78 | +1.34 | +0.06 | +3.47 | +1.09 | +1.54 |
| *Recipe + Equivalents Strategy Prompting* | | | | | | | | |
| Recipe + Cultural | -2.56 | -0.89 | -0.06 | -1.83 | +0.84 | +1.73 | +0.72 | +1.10 |
| Recipe + Descriptive | -8.69 | -1.81 | +2.29 | -2.74 | -4.74 | **+5.27** | +3.86 | +1.46 |
| Recipe + Functional | **+2.27** | **+4.02** | **+2.57** | **+2.95** | -0.96 | +2.80 | **+7.97** | **+3.27** |

Table 6: Ablation study comparing COMET scores for GPT-3.5 and GPT-4o using different equivalent strategies across CSI categories. The overall score is calculated as the average of CSI-1, CSI-2, and CSI-3 scores for each method.

like pork, chicken, or fish, making it difficult to provide the correct ingredient. In contrast, CSI-2 and CSI-3 usually involve longer, more specific phrases like "蚂蚁上树" *(Fried vermicelli with pork)*, making it easier to find an exact recipe match, reduce noise, and majorly improve accuracy.

### 6.5 Evaluation of Individual Equivalent Strategies

We further perform an ablation analysis of the recipe and individual equivalent strategies, including cultural, descriptive and functional, against the baseline results.

Table 6 shows that for GPT-3.5, the functional equivalent strategy outperforms others, especially when combined with the recipe. For GPT-4o, both descriptive and functional strategies yield better results in CSI-2 and CSI-3, with descriptive strategy excelling in CSI-2 when a recipe is included. In CSI-3, "Recipe + Functional" strategy leads to a significant performance boost of +7.97.

Upon reviewing the translations, both descriptive and functional strategies align well with the gold standards for CSI-2 and CSI-3. However, due to its complexity, the descriptive strategy produces

longer translations with trivial details for CSI-3, which is likely to negatively affect COMET scores.

### 6.6 Human Evaluation

We collect ratings from 10 native Chinese speakers fluent in English, based on the concept of cross-cultural adaptation on a scale of 0 to 10, alongside automatic quantitative metrics. We select the top-performing methods with recipes, as evaluated by COMET in Tables 5 and 6. We then randomly sample 15 entries with perfect agreement from each CSI category (1: Concrete, 2: Creative, 3: Abstract), totalling 45 entries.

The human evaluation results reveal a trend of performance improvement from CSI-1 to CSI-3 in GPT-3.5 and 4o. We use green to highlight cells with major improvements, i.e. over 1 point. Interestingly, for more complex CSIs (i.e. CSI-2 and CSI-3) we have larger improvements. We also observe that these trends align well with COMET trends in Table 5, noting that by both metrics, translation theory prompts yield significantly better results than basic prompting across categories.

Interestingly, human evaluators prefer "*Recipe + Neutralisation*" instead of "Recipe + Equivalent",

|  | GPT-3.5 | | | | GPT-4o | | | |
|---|---|---|---|---|---|---|---|---|
|  | CSI-1 | CSI-2 | CSI-3 | Overall | CSI-1 | CSI-2 | CSI-3 | Overall |
| Baseline | 6.33 | 3.88 | 3.18 | 4.47 | 6.22 | 4.23 | 3.65 | 4.70 |
| *Recipe-based Translation* | | | | | | | | |
| Recipe | -0.93 | +0.67 | +1.74 | +0.49 | -0.04 | +0.80 | +2.28 | +1.01 |
| Recipe + Explain-then-Translate | -0.03 | +0.60 | +1.35 | +0.64 | +0.43 | +0.98 | +1.68 | +1.03 |
| *Recipe + Translation Studies Prompting* | | | | | | | | |
| Recipe + Functional | -1.15 | +1.10 | +1.71 | +0.60 | +0.75 | +1.77 | +2.14 | +1.05 |
| Recipe + Neutralisation | +0.37 | +1.21 | +1.03 | +0.62 | +1.32 | +2.83 | +3.20 | +1.95 |
| Recipe + Equivalent | -0.85 | +0.62 | +2.05 | +0.81 | +0.71 | +0.99 | +2.38 | +1.36 |

Table 7: Difference in human evaluation of translation quality compared to baseline for different models and strategies across CSI categories. The overall score is calculated as the average of CSI-1, CSI-2, and CSI-3 scores for each method.

the highest in COMET. This preference may stem from the neutralisation definition used in this study, based on the findings of Amenador and Wang (2022). They note that neutralisation is the most common strategy employed by human translators for Chinese names, suggesting a familiarity that could influence the evaluators' preferences towards human-like translation outputs.

Table 8 illustrates the effectiveness of various translation strategies applied to the Chinese dish "三不沾", known for its non-stick quality when served, featuring osmanthus eggs. The full translation examples are provided in Appendix A.3.

| Strategy | Translation |
|---|---|
| Baseline | Not sticky in three ways |
| Equivalent | Sweet Egg Pastry |
| Neutralisation | Osmanthus Egg Custard |

Table 8: Selected Translations from GPT-4o Using Different Translation Strategies with Recipe.

"Sweet Egg Pastry" generated using the Equivalent strategy by GPT-4o, effectively conveys the essence of the dish by focusing on its key ingredients and flavour profile. "Osmanthus Egg Custard," produced through the Neutralisation strategy, is also an accurate translation as it highlights "Osmanthus egg," the main ingredient, and "Custard," indicating the dish's texture. In contrast, the baseline translation "Not sticky in three ways" fails to provide meaningful information about the dish, making it the weakest.

## 7 Discussion

The CSI categorisation can be applied to wider cultural domains that contain figurative elements. Future research can use this taxonomy to analyze how different translation methods perform on figurativeness and cultural specificity, suggesting a new

framework for evaluating CSI translation. This is similar to the evaluating framework in cultural inference, categorising entailment in different levels to better assess an LLM's ability to understand cultural inference (Huang and Yang, 2023).

CSI automatic identification offers a cost-effective approach that outperforms GPT-based prompting in CSI-2 and 3. This method is versatile and applicable to both general and domain-specific CSIs, as it focuses on preserving meaning in translation. It could enhance the quality of translations in a wide variety of domains where maintaining cultural integrity is essential – like literature, media, marketing and cross-cultural communication.

The findings of this paper also demonstrate the effectiveness of prompt strategies inspired by translation studies in overcoming the challenges of translating CSIs, particularly when direct equivalents are lacking across cultures. This approach shows promise for using LLMs with tailored prompts, integrating human translation insights, and translating diverse cultural elements more effectively.

## 8 Conclusion

In this paper, we introduce the ChineseMenuCSI dataset for CSI-rich dishes and propose a detailed classification in the test set. The results show that LLMs outperformed NMT systems, while NMT is better for Non-CSI translations. Additionally, automatic methods are better than GPT-based prompting at identifying CSIs in most categories.

Incorporating translation studies and recipe details improves LLMs' translation of Chinese dish names. Equivalence strategies, aligned with popular restaurant translations, yield consistently high-quality results, while neutralisation, based on previous analyses, is well-received by evaluators.

## Limitations

We acknowledge a few limitations of our study. Firstly, we use COMET as the primary automatic evaluation metric for CSIs. While COMET provides a robust evaluation, assessing cultural awareness may require an even deeper understanding of cultural backgrounds in both source and target languages, which COMET may not fully cover. Currently, in the absence of a metric that can evaluate text-to-text cultural similarity, we use COMET due to its high correlations with human judgment.

Secondly, while we only test zero-shot prompting for translation studies and recipe information, other research, such as Nayak et al. (2024), has demonstrated promising results using few-shot in-context learning strategies, which should also be explored.

Lastly, we only sample 45 menu entries from the test set, which can be relatively small compared to the studies with a larger test set. To achieve more robust and reliable results, increasing the number of human evaluators and the sample size of evaluation entries would be beneficial.

## Ethical Considerations

In conducting this research, we adhere to ethical guidelines to ensure the integrity and responsibility of our work. The ChineseMenuCSI dataset is created by scraping publicly available restaurant websites, ensuring that no private or sensitive information is collected. We obtain data in compliance with the terms of use of the websites and anonymise any identifying details of the restaurants. The human annotators involved in this study are fully informed about the nature of the research and provide their consent. We make the dataset available for research purposes under a license that respects the rights of the original content creators.

## References

Javier Franco Aixelá. 1996. Culture-specific items in translation. *Translation, power, subversion*, 8:52–78.

Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. Investigating cultural alignment of large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.

Román Álvarez and M Carmen Africa Vidal. 1996. *Translation, power, subversion*, volume 8. Multilingual Matters.

Kate Benedicta Amenador and Zhiwei Wang. 2022. The translation of culture-specific items (csis) in chinese-english food menu corpus: A study of strategies and factors. *SAGE Open*, 12(2):215824402210966.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Yong Cao, Yova Kementchedjhieva, Ruixiang Cui, Antonia Karamolegkou, Li Zhou, Megan Dare, Lucia Donatelli, and Daniel Hershcovich. 2024. Cultural Adaptation of Recipes. *Transactions of the Association for Computational Linguistics*, 12:80–99.

Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. Assessing cross-cultural alignment between ChatGPT and human societies: An empirical study. In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, pages 53–67, Dubrovnik, Croatia. Association for Computational Linguistics.

Alex J Chan, José Luis Redondo García, Fabrizio Silvestri, Colm O'Donnel, and Konstantina Palla. 2023. Harmonizing global voices: Culturally-aware models for enhanced content moderation. *arXiv preprint arXiv:2312.02401*.

Isabelle Chou, Victoria Lei, Defeng Li, and Yuanjian He. 2016. *Translational Ethics from a Cognitive Perspective: A Corpus-Assisted Study on Multiple English-Chinese Translations*, pages 159–173.

Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.

Georgiana Dinu, Prashant Mathur, Marcello Federico, and Yaser Al-Onaizan. 2019. Training neural machine translation to apply terminology constraints. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3063–3068, Florence, Italy. Association for Computational Linguistics.

Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.

Danilo Gallo, Jutta Willamowski, Yada Wisatekaew, Adrien Bruyat, and Antonietta Maria Grasso. 2021. Restaurant menu understanding: Illustrating the need for culturally augmented translation. In *22nd International Conference on Human-Computer Interaction with Mobile Devices and Services*, MobileHCI '20, New York, NY, USA. Association for Computing Machinery.

Federico Garcea, Margherita Martinelli, Maja Miličević Petrović, and Alberto Barrón-Cedeño. 2023. ! translate: When you cannot cook up a translation, explain. In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, pages 392–398.

HyoJung Han, Jordan Boyd-Graber, and Marine Carpuat. 2023. Bridging background knowledge gaps in translation with automatic explicitation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9718–9735, Singapore. Association for Computational Linguistics.

Daniel Hershcovich, Stella Frank, Heather Lent, Miryam de Lhoneux, Mostafa Abdou, Stephanie Brandl, Emanuele Bugliarello, Laura Cabello Piqueras, Ilias Chalkidis, Ruixiang Cui, Constanza Fierro, Katerina Margatina, Phillip Rust, and Anders Søgaard. 2022. Challenges and strategies in cross-cultural NLP. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6997–7013, Dublin, Ireland. Association for Computational Linguistics.

Junjie Hu, Hiroaki Hayashi, Kyunghyun Cho, and Graham Neubig. 2022. Deep: Denoising entity pretraining for neural machine translation. *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Jing Huang and Diyi Yang. 2023. Culturally aware natural language inference. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7591–7609, Singapore. Association for Computational Linguistics.

Vivek Iyer, Edoardo Barba, Alexandra Birch, Jeff Pan, and Roberto Navigli. 2023a. Code-switching with word senses for pretraining in neural machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12889–12901, Singapore. Association for Computational Linguistics.

Vivek Iyer, Pinzhen Chen, and Alexandra Birch. 2023b. Towards effective disambiguation for machine translation with large language models. In *Proceedings of the Eighth Conference on Machine Translation*, pages 482–495, Singapore. Association for Computational Linguistics.

Yova Kementchedjhieva, Di Lu, and Joel Tetreault. 2020. The ApposCorpus: a new multilingual, multi-domain dataset for factual appositive generation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1989–2003, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.

Fajri Koto, Haonan Li, Sara Shatnawi, Jad Doughman, Abdelrahman Sadallah, Aisha Alraeesi, Khalid Almubarak, Zaid Alyafeai, Neha Sengupta, Shady Shehata, Nizar Habash, Preslav Nakov, and Timothy Baldwin. 2024a. ArabicMMLU: Assessing massive multitask language understanding in Arabic. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 5622–5640, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Fajri Koto, Rahmad Mahendra, Nurul Aisyah, and Timothy Baldwin. 2024b. Indoculture: Exploring geographically-influenced cultural commonsense reasoning across eleven indonesian provinces. *arXiv preprint arXiv:2404.01854*.

Kai-Chee Lam, Man-Ling Ng, Lay-Hoon Ang, and Radina Mohamad Deli. 2018. Between concrete and abstract: The malaysian chinese way of naming dishes. *International Communication of Chinese Culture*, 5:247–259.

Cheng Li, Mengzhou Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. 2024a. Culturellm: Incorporating cultural differences into large language models. *arXiv preprint arXiv:2402.10946*.

Cheng Li, Damien Teney, Linyi Yang, Qingsong Wen, Xing Xie, and Jindong Wang. 2024b. Culturepark: Boosting cross-cultural understanding in large language models. *arXiv preprint arXiv:2405.15145*.

Daniel Liebling, Katherine Heller, Samantha Robertson, and Wesley Deng. 2022. Opportunities for human-centered evaluation of machine translation systems. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 229–240, Seattle, United States. Association for Computational Linguistics.

Chen Liu, Fajri Koto, Timothy Baldwin, and Iryna Gurevych. 2024. Are multilingual LLMs culturally-diverse reasoners? an investigation into multicultural proverbs and sayings. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2016–2039, Mexico City, Mexico. Association for Computational Linguistics.

Xiao Liu, Yansong Feng, Jizhi Tang, Chengang Hu, and Dongyan Zhao. 2022. Counterfactual recipe generation: Exploring compositional generalization in a realistic scenario. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7354–7370, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Reem I Masoud, Ziquan Liu, Martin Ferianc, Philip Treleaven, and Miguel Rodrigues. 2023. Cultural

alignment in large language models: An explanatory analysis based on hofstede's cultural dimensions. *arXiv preprint arXiv:2309.12342*.

Tarek Naous, Michael J Ryan, Alan Ritter, and Wei Xu. 2023. Having beer after prayer? measuring cultural bias in large language models. *arXiv preprint arXiv:2305.14456*.

Shravan Nayak, Kanishk Jain, Rabiul Awal, Siva Reddy, Sjoerd van Steenkiste, Lisa Anne Hendricks, Karolina Stańczak, and Aishwarya Agrawal. 2024. Benchmarking vision language models for cultural understanding. *arXiv preprint arXiv:2407.10920*.

P. Newmark. 1988. *A Textbook of Translation*. English language teaching. Prentice-Hall International.

Eugene Albert Nida and Charles R. Taber. 2003. *The Theory and Practice of Translation*. Leiden: Brill.

Valerie Pellatt and Eric Liu. 2010. *Thinking Chinese Translation : A Course in Translation Method: Chinese to English*. London: Routledge.

Denis Peskov, Viktor Hangya, Jordan Boyd-Graber, and Alexander Fraser. 2021. Adapting entities across languages and cultures. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3725–3750, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at trec-3. *Nist Special Publication Sp*, 109:109.

Pushpdeep Singh, Mayur Patidar, and Lovekesh Vig. 2024. Translating across cultures: Llms for intralingual cultural adaptation. *arXiv preprint arXiv:2406.14504*.

Melanie Soderstrom, Marisa Casillas, Elika Bergelson, Celia Rosemberg, Florencia Alam, Anne S Warlaumont, and John Bunce. 2021. Developing a cross-cultural annotation system and metacorpus for studying infants' real world language experience. *Collabra: Psychology*, 7(1):23445.

Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'14, page 3104–3112, Cambridge, MA, USA. MIT Press.

Y Tao, O Viberg, RS Baker, and RF Kizilcec. 2024. Cultural bias and cultural alignment of large language models. *arXiv preprint arXiv*, 2311.

Wenxuan Wang, Wenxiang Jiao, Jingyuan Huang, Ruyi Dai, Jen-tse Huang, Zhaopeng Tu, and Michael Lyu. 2024. Not all countries celebrate thanksgiving: On the cultural dominance in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6349–6384, Bangkok, Thailand. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2024. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.

Binwei Yao, Ming Jiang, Diyi Yang, and Junjie Hu. 2024. Benchmarking llm-based machine translation on cultural awareness. *Preprint*, arXiv:2305.14328.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2765–2781, Mexico City, Mexico. Association for Computational Linguistics.

# A Appendix

## A.1 Detailed Prompt Teamplates

---

**Prompt Strategy: Recipe + Equivalents**

**User:**

Similiar Recipe: [Recipe Instructions].

Based on the above recipe information, provide three translations for [Chinese dish name] based on the three translation strategies listed below and select the best one:

Cultural Equivalent: Substituting a source language term with a term from the target language that has similar cultural resonance and functionality.

Functional Equivalent: Rendering the source language's meaning, intent, and style into the target language in a culturally appropriate and understandable way. This strategy prioritizes the effect and function of the text in the target culture over a word-for-word translation, ensuring the translation fulfills the same purpose as the original.

Descriptive Equivalent: Providing an in-depth explanation of a term or concept that lacks a straightforward equivalent in the target language. The explanation could include details such as ingredients, culinary method, key characteristics, etc.

---

Figure 3: Recipe + Equivalents Detailed Prompt

---

**Prompt Strategy: Recipe + Neutralisation**

**User:**

Similiar Recipe: [Recipe Instructions].

Based on the above recipe information, provide a translation for [Chinese dish name] with the following translation strategy:

Menu Description Strategy: This strategy involves using culturally neutral language to describe or explain a cultural word, phrase, or rhetorical expression from the source text (ST). It answers the question, 'What is this?' and is similar to converting a metaphor to its literal meaning. The translations should include the key culinary method, ingredients, and characteristics.

---

Figure 4: Recipe + Neutralisation Detailed Prompt

## A.2 Menu Parser

We develop a heuristic parser to extract dish information from the bilingual menu images by detecting price tags and segmenting the raw text into aligned content. To achieve this, we utilise Google Cloud Vision OCR[6] to extract text and bounding boxes from the menu images. Price tags serve as unique indicators for each dish's content, as we observe that most menus included prices alongside their respective dishes. These price tags are identified using regular expressions, such as "dd.dd" or "£dd.dd".

Given that the position of price tags relative to dish names can vary across menus, we calculate

---

[6]Google Cloud Vision OCR: https://cloud.google.com/vision

---

alignment scores based on the cosine similarity and the gap distance between the potentially aligned Chinese and English text and select the alignment with the highest score from all possible combinations. Each entry undergoes manual review to ensure accuracy and errors are corrected before subsequent steps.

## A.3 Full Examples of Translation Prompts

| Strategy | Translation |
|---|---|
| Baseline | Not sticky in three ways |
| Reference | Sweet Egg Pudding |
| Recipe | Three Non-Stick Delicacy: Traditional Imperial Egg Yolk Treat |
| Recipe + EtT | Imperial Non-Stick Egg Delight |
| Equivalent | Sweet, Sticky and Chewy |
| Neutralisation | Non-Sticky Sweet and Savory Egg Custard |
| Recipe + Equivalent | Sweet Egg Pastry |
| Recipe + Neutralisation | Osmanthus Egg Custard |

Table 9: Comparison of GPT-4o Translations Across Different Strategies

1271

# Pitfalls and Outlooks in Using COMET

**Vilém Zouhar⋆ [1]   Pinzhen Chen⋆ [2]   Tsz Kin Lam[2]   Nikita Moghe[2]   Barry Haddow[2]**

[1]ETH Zurich      [2]University of Edinburgh

vzouhar@ethz.ch   {pinzhen.chen,tlam,nikita.moghe,bhaddow}@ed.ac.uk

## Abstract

The COMET metric has blazed a trail in the machine translation community given its strong correlation with human judgements of translation quality. Its success stems from being a pre-trained multilingual model finetuned for quality assessment. However, it being a neural metric also gives rise to a set of pitfalls that may not be widely known. We investigate these unexpected behaviours from three aspects: 1) technical: obsolete software versions and compute precision; 2) data: empty content, language mismatch, and translationese at test time as well as distribution and domain biases in training; 3) usage and reporting: multireference support and model referencing in the literature. All of these problems imply that COMET scores may be not incomparable between papers or technical setups and we put forward our perspective on fixing each issue. Furthermore, we release the sacreCOMET package that can generate a signature for the software and model configuration as well as an appropriate citation. The goal of this work is to help the community make more sound use of the COMET metric.

## 1   Introduction

Automated metrics provide a cheap and scalable way of evaluating and benchmarking NLP models. In machine translation (MT), the evaluation protocol has moved from string matching metrics (BLEU, TER, chrF, inter alia; Papineni et al., 2002; Snover et al., 2006; Popović, 2015) to trained neural metrics (Shimanaka et al., 2018; Takahashi et al., 2020; Rei et al., 2020a; Sellam et al., 2020) with COMET being widely adopted. The trained metrics have been shown to correlate much better with human judgement (Freitag et al., 2021, 2022b, 2023), making them more reliable in estimating translation quality and ranking translation systems.

⋆Equal contributions.
[0]Code: github.com/PinzhenChen/sacreCOMET

Nonetheless, the solution to translation evaluation is yet to be perfected. One problem is the haphazard use of the metric. Previously, Post (2018) showed that different usages and implementations of BLEU, e.g. tokenization and smoothing, lead to inconsistencies in scores. We suspect that the use of COMET might be sensitive to misconfigurations too, resulting in unexpected behaviours. Furthermore, trained MT metrics are optimized on a limited amount of data (usually valid machine translations), leading to overfitting and reduced robustness against corner cases. Contributions of this work are listed as follows:

- we reveal nine problems spanning technical issues, data biases, and model reporting;
- we show that inconsistent use of COMET leads to non-comparable scores across papers or setups;
- we release the sacreCOMET package for better reporting and reproducibility;
- we provide directions for future work on building learned metrics.

## 2   Background and Setup

**Metric background.**   Publicly available human judgements of translation quality come from shared task annotation campaigns, where translations are evaluated with some annotation protocol. From 2017, in WMT, the protocol was a variant of direct assessment (DA; Graham et al., 2013) which has annotators providing a number from 0 (lowest) to 100 (highest) as the segment quality. This has been subsequently replaced by MQM and ESA protocols (Lommel et al., 2014; Kocmi et al., 2024d), though DA remains the most abundant data source for neural metric training.

Automated metrics aim to yield scores that correlate with human judgements of translations. Most metric scores are computed at the segment level and then aggregated at the system level to e.g. obtain system comparison. The evaluation of metrics is done with respect to the human judgements.

**COMET models.** Metrics such as chrF or BLEU are heuristic algorithms that match n-grams between the translation and the reference to compute a score. In contrast, COMET is a machine learning model fine-tuned from a pre-trained multilingual language model, e.g. XLM-R (Conneau et al., 2020), with an additional regression head. A reference-based COMET model learns to regress from a tuple of [source, hypothesis, reference] to the human judgement score (from previous evaluation campaigns) at the segment level. The quality estimation (reference-free) version of COMET is prepared by omitting the reference from the input.

Most issues in this work are demonstrated using two COMET checkpoints unless noted otherwise: the reference-based $COMET_{22}^{DA}$ and the reference-free $COMET_{22}^{kiwiDA}$ (Rei et al., 2022a). Both metrics output to a normalized range between 0 and 1. The COMET framework unbabel-comet is of version 2.2.2 except when we test how different software versions affect COMET scores.

**Data setup.** We base our experiments on the general domain translation and metrics shared tasks of WMT from 2023 (Kocmi et al., 2023; Freitag et al., 2023). The translation directions in the paper are centred around En↔De and En↔Zh, though we occasionally include other translation directions for demonstrative purposes.

Whenever possible, we divide all scores, including DA and model outputs, such that their output is between 0 and 1.

## 3 Problems

In this section, we identify and test nine possible pitfalls or curious behaviours with COMET which are not all well-studied. In three groups, these are:

- **Technicality**: obsolete Python and COMET software versions as well as compute precisions could lead to inaccurate score computation.
- **Training and test data**: COMET as a neural metric, might be derailed by empty hypotheses, language mismatch, and translationese at test time. It may also follow the training data biases.
- **Tool usage and score interpretation**: COMET has no defined way of equipping multiple references when available which leaves room for research. From a bibliometric perspective, we reveal that some literature omits a clear reference to the checkpoint version or citation.

In addition, we discuss some final issues that need more attention from the community.

| Python | 3.7.16 | 3.8.11 | 3.12.4 |
|---|---|---|---|
| unbabel-comet | 1.1.2 | 2.2.2 | 2.2.2 |
| En→De | 0.796 | 0.837 | 0.837 |
| En→Zh | 0.911 | 0.862 | 0.862 |
| De→En | 0.851 | 0.855 | 0.855 |
| Zh→En | 0.795 | 0.803 | 0.803 |

Table 1: $COMET_{22}^{DA}$ scores for WMT 23 Online-A under different package versions.

### 3.1 Software versions [technical]

The official installation of the COMET package requires Python 3.8 or above.[1] We demonstrate that neglecting this would lead to unexpected scores because the same COMET checkpoint can produce vastly different scores with previous COMET framework versions that are no longer supported.

Under several Python versions, executing the following code leads to different COMET package (unbabel-comet) versions being installed. Running the framework for translation evaluation will subsequently result in false conclusions as shown in Table 1's evaluation on WMT23 tests. The direct cause is that Python 3.7, which has been discontinued, only supports unbabel-comet versions up to 1.1.2. Nonetheless, we caution that the underlying factor is the version of unbabel-comet rather than Python.

```
$ pip install pip --upgrade
$ pip install unbabel-comet --upgrade

# will install
# unbabel-comet==1.1.2 under Python 3.7.16
# unbabel-comet==2.2.2 under Python 3.8.11
# unbabel-comet==2.2.2 under Python 3.12.4
```

**Recommendation.** Updating both Python and unbabel-comet to their latest versions is helpful and reporting the toolkit version aids reproducibility.

### 3.2 Numerical precision [technical]

Model quantization represents a model using lower-numerical precision data types so that the model consumes less memory and model passes can be computed faster. Such improvement in inference is directly beneficial to deployment efficiency; it is also useful in other complex procedures involving COMET scoring, such as data filtering, re-ranking, and Minimum Bayes Risk (MBR) decoding (Kumar and Byrne, 2004).

Despite the aforementioned advantages, model quantization is not a feature supported by the cur-

---

[1] github.com/Unbabel/COMET 332dfb0 as of Aug 2024.

| | Precision | COMET$_{22}^{DA}$ | MAE | $\tau_c$ | Acc | Time (s) |
|---|---|---|---|---|---|---|
| **En→De** | GPU FP32 | 0.822 | 10.4 | 0.274 | 0.885 | 113 |
| | GPU FP16 | 0.822 | 10.4 | 0.274 | 0.885 | 55 |
| | FP32 | 0.822 | 10.4 | 0.274 | 0.885 | 2262 |
| | CPU FP16 | 0.822 | 10.4 | 0.274 | 0.885 | 2403 |
| | QINT8 | 0.852 | 10.8 | 0.109 | 0.385 | 1856 |
| **De→En** | GPU FP32 | 0.841 | 9.98 | 0.296 | 0.901 | 87 |
| | GPU FP16 | 0.841 | 9.98 | 0.296 | 0.901 | 48 |
| | FP32 | 0.841 | 9.98 | 0.296 | 0.901 | 1674 |
| | CPU FP16 | 0.841 | 9.99 | 0.296 | 0.901 | 1758 |
| | QINT8 | 0.860 | 10.7 | 0.164 | 0.516 | 1249 |
| **En→Zh** | GPU FP32 | 0.842 | 11.7 | 0.290 | 0.933 | 111 |
| | GPU FP16 | 0.842 | 11.7 | 0.290 | 0.933 | 102 |
| | FP32 | 0.842 | 11.7 | 0.290 | 0.933 | 1751 |
| | CPU FP16 | 0.842 | 11.7 | 0.290 | 0.933 | 1710 |
| | QINT8 | 0.881 | 13.2 | 0.031 | 0.608 | 1258 |
| **Zh→En** | GPU FP32 | 0.799 | 9.95 | 0.153 | 0.717 | 113 |
| | GPU FP16 | 0.799 | 9.95 | 0.153 | 0.725 | 86 |
| | FP32 | 0.799 | 9.95 | 0.153 | 0.717 | 1936 |
| | CPU FP16 | 0.799 | 9.95 | 0.153 | 0.717 | 1995 |
| | QINT8 | 0.872 | 10.5 | 0.081 | 0.475 | 1351 |

Table 2: System ranking with quantization on GPU and CPU. COMET$_{22}^{DA}$ is the absolute model score; MAE, $\tau_c$, and Acc are mean average error, correlation, and accuracy with respect to human judgements; Time refers to computation time in seconds.

rent COMET framework except in a concurrent work (Gowda et al., 2024). We make minimal modifications to the software and investigate the effect of numerical precision on COMET scores on both CPU (FP32, FP16, and QINT8) and GPU (FP32 and FP16). When using FP16, we first load the model weights to FP32, followed by `.half()` call. This is because loading the weights directly in FP16 still incorrectly results in FP32 precision. For CPU inference with dynamic QINT8, we apply the quantization module `torch.ao.quantization` from PyTorch.

We use AMD Ryzen 9 5900X with NVIDIA GeForce RTX 3090 for GPU inference and a batch size of 8 in all settings (in practice a quantized model makes room for a larger batch size). Table 2 summarises the effect of numerical precision. In addition to reporting COMET scores, we also report (1) inference time in seconds (sec) as an efficiency measure; and (2) segment-level mean absolute error (MAE), segment-level Kendall's tau-c ($\tau_c$), and system-level pairwise accuracy (Acc). Everything is compared to the human DA scores either on segment- or system-level. Technically, Kendall's $\tau_c$ calculates rank correlation on an ordinal scale with adjustments for ties and pairwise accuracy computes the proportion of system pairs that have the same ordering by a metric as by humans.

Our results show that there is no meaningful difference between FP32 and FP16 in both CPU and GPU devices up to 3 significant figures. On GPU, FP16 is about 30% faster in time, but unsurprisingly it does not provide any speed-up on CPU. Interestingly, on the CPU, dynamic QINT8 gives systematically higher COMET scores and shorter running times than FP32 and FP16. However, the much lower $\tau_c$ and pairwise accuracy indicate the lack of reliability at this precision. In addition to precision, we explored the effect of batch size and the choice of GPU or CPU during inference with results listed Appendix A. Whilst there are some fluctuations, they are mostly negligible. However, lower precision allows for higher batch size which usually directly corresponds to speed-up.

**Recommendation.** If GPU is available, it is feasible to run COMET with FP16 with a larger batch size for much faster inference without any quality loss. Otherwise, FP32 should be used.

### 3.3 Empty hypothesis [data]

An empty translation (a string of length 0) gets penalized heavily by string-based metrics because an empty string has zero surface overlap with the reference. However, neural metrics provide no such guarantee. We show that COMET assigns a positive instance-level score even if the hypothesis is an empty string as corroborated by Lo et al. (2023). Problematically, this score can even occasionally be higher than that of a genuine system translation.

In Table 3, we list COMET scores for system Online-A's hypotheses at WMT23 and a file full of empty lines. Furthermore, we compare them with completely incorrect translations to explicate the score magnitude in two ways:

- *Random hypothesis:* we shuffle WMT22's reference files at the sentence level in the respective translation directions. This provides us with high-quality human-written sentences. We sub-sample or over-sample if the number of lines in WMT22 is larger or smaller than the WMT23 size.
- *Random hypothesis (shuffled words):* we further shuffle the words at each line in the sentence-shuffled files, generating nonsensical sentences.

Sentence-shuffled hypotheses can be seen as fluent but extremely inadequate sentences whereas word-shuffled sentences are neither fluent nor adequate.

We observe that sentence-shuffled hypotheses attain comparable scores to empty ones, but word-shuffled hypotheses have the lowest scores across

| Hypothesis (↓) | COMET$_{22}^{DA}$ | | | | COMET$_{22}^{kiwiDA}$ | | | |
|---|---|---|---|---|---|---|---|---|
| | En→De | En→Zh | De→En | Zh→En | En→De | En→Zh | De→En | Zh→En |
| Real system (Online-A) | 0.837 | 0.862 | 0.855 | 0.803 | 0.800 | 0.791 | 0.794 | 0.787 |
| Empty hypothesis | 0.335 | 0.392 | 0.353 | 0.374 | 0.315 | 0.319 | 0.537 | 0.467 |
| Random hypothesis (fluent) | 0.373 | 0.434 | 0.334 | 0.350 | 0.333 | 0.341 | 0.447 | 0.391 |
| Random hypothesis (shuffled words) | 0.244 | 0.419 | 0.264 | 0.347 | 0.232 | 0.325 | 0.307 | 0.385 |

Table 3: Absolute average COMET scores for WMT23 Online-A, empty hypotheses, and random sentences. Random sentences are either fluent but irrelevant or perturbed with words shuffled and thus non-fluent.

the majority of the translation directions. Empty and shuffled hypotheses, despite having much lower COMET than the valid translations, would not be assigned zero scores by COMET$_{22}^{DA}$ or COMET$_{22}^{kiwiDA}$, showing that COMET is more lenient than string overlap-based metrics in penalizing such irregularities.

We count the number of empty lines that score better than a translation from Online-A in Table 4. We observe roughly 0.25% such cases for most translation directions except for De→En's COMET$_{22}^{kiwiDA}$ score situating at 1.45%. Further, in Figure 1 we plot the distributions of COMET scores for Zh→En's empty and genuine translations with other translation directions in Appendix B. Noticeable overlaps are observed for COMET$_{22}^{kiwiDA}$ when translating into English.

Adhering to the DA protocol guidelines, this is not the proper behaviour because an empty hypothesis should receive a score of 0, to match *no meaning preserved*. This however is unsurprising with COMET, which has likely not seen empty hypotheses during training that would have received a score of 0 from a human annotator. Finally, even by relaxing the 0-score expectation, the metric should still assign the same score to all empty hypotheses regardless of the source. Since the distributions of empty hypotheses are nowhere close to a single vertical bar in Figure 1, it exposes the issue that segment-level COMET scores oddly hinge on the source sentence, as noted by Sun et al. (2020).

**Recommendation.** Force empty hypotheses to have 0 scores before aggregating. Also, a string-based metric like BLEU or chrF should be used to catch similarly malformed hypotheses.

### 3.4 Hypothesis language mismatch [data]

String overlap-based metrics can also score a hypothesis in a language different from the reference almost zero, especially with script mismatch. However, even for the reference-based COMET, there is no explicit way to enforce the intended target



Figure 1: Distribution of instance-level scores for empty and baseline translations (x-axis: score; y-axis: count). See other translation directions in Appendix B.

| | translation < empty | |
|---|---|---|
| | COMET$_{22}^{DA}$ | COMET$_{22}^{kiwiDA}$ |
| En→De | 0 / 558 | 2 / 558 |
| En→Ru | 1 / 2075 | 0 / 2075 |
| En→Uk | 2 / 2075 | 1 / 2075 |
| En→Zh | 6 / 2075 | 2 / 2075 |
| De→En | 1 / 550 | 8 / 550 |
| Ru→En | 5 / 1724 | 6 / 1724 |
| Uk→En | 1 / 1827 | 5 / 1827 |
| Zh→En | 5 / 1977 | 1 / 1977 |

Table 4: Proportion of WMT23 Online-A's translations that are worse than an empty line for the same source, displayed as "empty/total".

language. This poses an increasingly pronounced problem, especially for multilingual translation models as well as the recent large language models, in which the generated language cannot be as easily controlled (Zhang et al., 2023).

We conduct experiments to understand if translation outputs in an incorrect language impact the score, and whether different mismatching languages can lead to distinct patterns. We use the translation directions En→Ru, En→Uk, and En→Zh in WMT23 which share the same English source input. Having Online-A's output in all three directions, we substitute hypotheses in a particular translation direction with those from another direction. A similar hypothesis was presented by Amrhein et al. (2022) which suggested that COMET metrics are not robust to hypothesis language mis-

|  | En→Ru | En→Uk | En→Zh |
|---|---|---|---|
| Correct lang. | 0.853 | 0.832 | 0.862 |
| Incorrect target lang. | Uk: 0.797 Zh: 0.536 | Ru: 0.807 Zh: 0.540 | Ru: 0.655 Uk: 0.644 |
| Empty hyp. | 0.316 | 0.329 | 0.391 |
| Random hyp. | 0.463 | 0.472 | 0.435 |

Table 5: COMET$_{22}^{DA}$ scores for WMT23 Online-A's output in (1) correct, (2, 3) incorrect language, (4) empty outputs, (5) random, but fluent, output.

match. Our experiment setup offers a more detailed evaluation setup than their contrastive setup.

Table 5 presents COMET$_{22}^{DA}$ scores for translations in correct and incorrect languages, as well as empty lines and random sentences in the correct language as "baselines", deemed as completely wrong translations. The pattern shows that when the hypothesis is in a language distant from the reference, the COMET score declines much more than when the hypothesis is in a similar language. A more concerning issue is that even when the hypotheses' language is completely wrong, the resulting COMET score can still be vastly higher than empty hypotheses or random sentences in the correct language. We omit COMET$_{22}^{kiwiDA}$ because it does not have a mechanism to read a reference (language) making it inherently incapable of distinguishing output languages.

**Recommendation.** Run language identification and set hypotheses in an unexpected language to have a 0 COMET score before aggregating them system-level. Also, check with a string overlap-based metric like BLEU or chrF.

| Lang | Score | Lang | Score | Lang | Score |
|---|---|---|---|---|---|
| De→En | 0.754 | En→Is | 0.666 | Hi→Bn | 0.910 |
| Ps→En | 0.670 | En→Lt | 0.600 | De→Fr | 0.792 |
| Is→En | 0.724 | En→Ru | 0.765 | Fr→De | 0.834 |
| Pl→En | 0.761 | En→Iu | 0.720 | Zu→Xh | 0.639 |
| Ru→En | 0.771 | En→Ha | 0.768 | De→Cs | 0.510 |
| Ja→En | 0.663 | En→Ja | 0.745 | Xh→Zu | 0.574 |
| Ta→En | 0.655 | En→Pl | 0.706 | Bn→Hi | 0.770 |
| Zh→En | 0.743 | En→Gu | 0.514 | | |
| Ha→En | 0.641 | En→Cs | 0.767 | | |
| Km→En | 0.659 | En→Zh | 0.775 | | |
| Lt→En | 0.726 | En→Fi | 0.616 | | |
| Cs→En | 0.740 | En→Ta | 0.709 | | |
| Gu→En | 0.575 | En→De | 0.841 | | |
| Kk→En | 0.649 | En→Kk | 0.574 | | |
| Iu→En | 0.724 | | | | |
| Fi→En | 0.719 | | | | |

Table 6: Average human DA score for each translation direction in WMT data up to 2023 (inclusive).



Figure 2: Setup of an experiment with bottom 75% of En→Zh scores which creates a bias in COMET$_{22}^{DA}$. In the new data for En→Zh (bottom right) there are no translations with perfect scores. En→De data are unaffected.

### 3.5 Score distribution bias [data]

As the COMET metric is a machine learning model trained on human ratings of existing machine translations, it inherits many properties of statistical learning such as data (output score) distribution bias. The yearly WMT shared task receives submissions with varying quality and potentially varying quality ranges for different translation directions per year (from Koehn and Monz, 2006 to Freitag et al., 2023). This is attributed to diverse factors: the availability of data, source-target language similarity, the level of interest in languages, etc. The gap in translation quality will then propagate into skewed human judgement scores across translation directions. When a single COMET model learns to score all translation directions, it can overfit the score distribution w.r.t. a translation direction in addition to the quality of a translation hypothesis.

We first verify this in Table 6 which shows that WMT translation directions are associated with vastly different human DA scores (from 0.51 to 0.91). Empirically, we illustrate this issue using two high-resource directions En→De and En→Zh. As shown in Figure 2, we keep either the top- or bottom-75% of all scored translations to alter the score distribution for each direction, simulating the scenario where low- and high-performing system submissions are received for different directions. We then train different COMET models on the human train data before and after alteration as per Figure 2. Finally, we evaluate those checkpoints on the same test set and report results in Table 7.

As expected, for both En→De and En→Zh, training on the top or bottom-scoring data leads to increased or decreased COMET scores on the same set of hypotheses. Besides, we observe that altering scores in a particular translation direction incon-

| Training data | En→De | En→Zh |
|---|---|---|
| All | 0.770 | 0.770 |
| Top-75% of En→De | 0.790 | 0.770 |
| Bot-75% of En→De | 0.765 | 0.764 |
| Top-75% of En→Zh | 0.783 | 0.789 |
| Bot-75% of En→Zh | 0.772 | 0.751 |

Table 7: Average scores from $COMET_{22}^{DA}$ trained on data with top- or bottom-75% of scores kept in a particular direction.

sistently affects scores in another direction systematically. For example, removing the bottom 25% En→De scores "improves" the test-time score from 0.770 to 0.790 whilst En→Zh remains unaffected.

Finally, an empty translation should have the same score irrespective of the original source. Due to the statistical learning nature of the metric, this is not the case, as found in Section 3.3. Sun et al. (2020) and Zouhar et al. (2023) show that some meaningful correlation with human scores can be attained with just the source as the input. This shows, that there is a learned bias based on the prior difficulty of the source segment, which is undesirable for an objective evaluation metric.

**Implication.** A trivial conclusion is that COMET scores for different translation directions are not comparable. Nevertheless, we caution that the same phenomenon could happen for other features such as the domain, output style, etc. Although z-score re-scaling *could* mitigate this problem, it has not been a common practice since WMT22 and it would further contribute to non-objective scores (Knowles, 2021). Moreover, while z-scoring is straightforward for the translation direction, other latent, language-agnostic biases still exist.

### 3.6 Domain bias [data]

Neural metrics like COMET are biased towards particular domains, manifested by worse test performance on unseen domains (Zouhar et al., 2024a). Taking inspiration from previous work and our discussions on "latent biases", we now raise a question—can we create adversarial hypotheses at test time to exploit the domain bias in training time? Specifically, different domains in the training data are associated with different score ranges. By pretending that a translation is in a particular domain, it might manipulate its COMET score.

To make it explicit to COMET during training, we prepend the target translation with a tag of its domain—in our case, the year the WMT data originated. Note that in each iteration of WMT, systems

| | |
|---|---|
| **Train** | 2020 Fire prevented from spreading |
| **Test** | 2019 Now I have to tell you a nice story. |
| | 2020 Now I have to tell you a nice story. |
| | 2021 Now I have to tell you a nice story. |
| | 2022 Now I have to tell you a nice story. |
| | 2023 Now I have to tell you a nice story. |

Table 8: An illustration of year-as-a-domain tagging during training and testing.

| Tag | Train | Test |
|---|---|---|
| 2018 | unseen | 0.736 |
| 2019 | 0.721 | 0.737 |
| 2020 | 0.735 | 0.744 |
| 2021 | 0.749 | 0.749 |
| 2022 | unseen | 0.747 |
| 2023 | unseen | 0.747 |
| 2024 | unseen | 0.739 |
| 2025 | unseen | 0.747 |

Table 9: Average $COMET_{22}^{DA}$ scores for subsets in training and predictions on test data. During testing, the whole test set had a single tag, e.g. 2024, irrespective of the data origin.

get higher overall DA scores (e.g. 0.721 for 2019 and 0.749 for 2021). Table 8 illustrates our setup: during training, we tag the scored translation data with its year; during testing, we trial various year prefixes to understand the effect.

One would expect the metric to produce the same score based solely on the translation quality. However, as shown in Table 9, by merely changing the year tag, we can influence the average score of the test set. During training, the model would be able to observe that 2019 is associated with the worst score and 2021 the best. During test time, the model follows this bias and also extrapolates it to upcoming years where it predicts an improvement in the average DA scores. While the differences appear small, they are on the same scale as the differences between years in the training data.

**Implication.** Our year-as-a-domain setting might be overly simple, but the vulnerability of COMET to latent biases cannot be neglected. Although Amrhein and Sennrich (2022) has shown that COMET is not sensitive to numbers, this work reveals that it can be systematically exploited in an artificial setting. We offer a more practical (adversarial) example that one may disguise biomedical domain translations as news translations to game COMET.

### 3.7 Lack of multi-reference support [usage]

In machine translation, there usually exist many valid translations for the same input. An effective metric should incorporate multiple ground truths

1277

| WMT23 | | COMET$_{22}^{DA}$ | | | | |
|---|---|---|---|---|---|---|
| | | ref | ref (alt) | avg | max | agg |
| MQM | He→En | 0.885 | 0.897 | 0.910 | 0.910 | **0.949** |
| | En→De | 0.974 | 0.936 | 0.974 | **0.987** | 0.974 |
| | Zh→En | 0.783 | 0.908 | 0.850 | 0.858 | **0.950** |
| DA | En→De | 0.885 | **0.949** | 0.910 | 0.897 | 0.910 |
| | Zh→En | 0.717 | 0.875 | 0.783 | 0.775 | **0.883** |
| | De→En | 0.901 | 0.912 | **0.923** | **0.923** | 0.912 |
| | En→Zh | **0.933** | 0.817 | 0.900 | 0.867 | **0.933** |
| | Cs→Uk | 0.846 | 0.802 | **0.901** | 0.890 | 0.868 |
| | En→Cs | 0.858 | **0.875** | 0.858 | 0.858 | 0.867 |
| | En→Ja | **0.941** | 0.824 | 0.934 | 0.934 | 0.926 |
| | Ja→En | 0.922 | 0.915 | **0.928** | **0.928** | 0.922 |

Table 10: Pairwise system-level accuracy for different strategies incorporating multiple references into COMET$_{22}^{DA}$. Evaluation is carried out on WMT23 with DA or MQM scores as human ground truths.

if available, thereby enhancing the accuracy and robustness of its evaluation. Existing metrics like BLEU or chrF rely on surface-level overlap to capture the ground truth space from multiple references while metrics like ParBLEU (Bawden et al., 2020) can automatically generate paraphrases of a given reference to be included during evaluation.

By design, only one reference can be used in COMET. Whilst one may argue that representing a text reference in the neural space can ease the restriction on word choices, it might still be beneficial to use multiple references to overcome defects in the base embedding model. Therefore, we test whether COMET can explicitly and reliably leverage multiple references. We identify three distinct ways in which multiple references have been incorporated in COMET in previous literature (Rei et al., 2020b; Zouhar and Bojar, 2024):

- **max**: Taking the maximum over the scores from multiple passes with different references.
- **avg**: Averaging the scores from multiple passes with different references.
- **agg**: Obtaining an aggregate score per example as follows. A quadruplet of source $\mathbf{s}$, hypothesis $\mathbf{h}$, reference $\mathbf{r}$, and alternative reference $\hat{\mathbf{r}}$ is fed to COMET six times in different [src, hyp, ref] arrangements: $[\mathbf{s},\mathbf{h},\mathbf{r}], [\mathbf{r},\mathbf{h},\mathbf{s}], [\mathbf{s},\mathbf{h},\hat{\mathbf{r}}], [\hat{\mathbf{r}},\mathbf{h},\mathbf{s}], [\mathbf{r},\mathbf{h},\hat{\mathbf{r}}]$, as well as $[\hat{\mathbf{r}},\mathbf{h},\mathbf{r}]$. Then, the average score from these six passes is multiplied by $(1-\sigma)$ where $\sigma$ denotes the standard deviation.

We use additional references from the WMT23 test set if available (He→En) or the outputs from the best-scoring system in each direction in the metrics shared task (Freitag et al., 2023) as an alter-

native reference. We report pairwise system-level accuracy (Kocmi et al., 2021) for various translation directions in Table 10. Our results suggest that there is no single method that can consistently take advantage of the inclusion of multiple references with the existing COMET implementation. At a higher inference cost, the six-pass aggregation with COMET might have a tiny edge over other methods when MQM is treated as human ground truths, but it is also outperformed under DA by single-reference or other multi-reference methods.

As translation systems have greatly improved lately, the above pattern might be explained by Freitag et al. (2020)'s finding that high-quality translation outputs do not benefit from multi-reference evaluation. We also caution that these observations are highly dependent on the quality of the underlying references. As studied previously, obtaining high-quality references is not trivial (Freitag et al., 2020, 2023; Zouhar and Bojar, 2024). Our use of the top-performing system outputs as alternate references is fit for the purpose but not optimal.

**Recommendation.** Our recommendations for the inclusion of multiple references into COMET or even other neural metrics are aspirational as this topic warrants further investigation. Extending unified pre-training (Wan et al., 2022) with multiple references in the architecture as well as using training objectives more suitable for handling more than one references (Zheng et al., 2018; Fomicheva et al., 2020a) can be helpful.

### 3.8 Translationese [data]

COMET has been trained with human translations as references and machine translations as hypotheses, where both could be deemed "translationese" to a certain extent (Gellerstam, 1986).

**Translationese in references.** We first conduct an experiment to see if the translationese present in the reference would undermine system evaluation with COMET. We consider WMT's official reference as a standard version and Freitag et al. (2020)'s paraphrased reference as a less translationese reference (we use their "paraphrased as-much-as-possible" version). Experiments are carried out under two settings: (1) WMT19 En→De submissions scored by COMET$_{22}^{DA}$, and (2) WMT20 En→De submissions scored by COMET$_{20}^{DA}$ (Rei et al., 2020b). These two settings cover two scenarios—whether the test suite has been used in training the COMET model, or not. A breakdown

```
Please paraphrase the following text as
much as possible. Provide the paraphrase
without any explanation:

$HYPOTHESIS
```

Figure 3: Prompt template used to request a paraphrase from GPT-4o, where $HYPOTHESIS is replaced by individual hypotheses.

of COMET scores and rankings for individual systems are listed in Appendix C Tables 12 and 13.

The COMET scores decline dramatically when we switch the reference from the original one to the paraphrased one—aiming to reduce translationese. It means that COMET is indeed sensitive to such changes in the reference. Yet interestingly, the overall system ranking in either setting remains rather stable. We find a very high Kendall's $\tau_a$ of 0.9827 and 0.9833 on the system rankings in the two settings; pairwise accuracy computed against human judgements also maintained at 0.924. We conclude that translationese in the reference impacts absolute COMET scores but not system ranking. These patterns are consistent regardless of whether the model has been exposed to the test set.

**Translationese in hypotheses.** We then attempt to understand if a varying degree of translationese in the system outputs will influence system ranking by COMET. We run $\text{COMET}_{22}^{\text{DA}}$ and $\text{COMET}_{22}^{\text{kiwiDA}}$ on WMT19 En→De system outputs as well as their corresponding rephrased outputs against the same source and reference. To acquire paraphrases affordably, we shortlist the top-10 systems' translations in the previous experiment and we prompt GPT-4o using the prompt outlined in Figure 3.[2] We do not feed the source sentence to prevent the model from revising the quality.

Appendix Table 14 shows that both models yield the same system ranking when the original hypotheses are scored. After substituting the hypotheses with their paraphrases, the ranking has changed more under $\text{COMET}_{22}^{\text{kiwiDA}}$ which witnesses much lower pairwise accuracy and Kendall's $\tau_c$ compared to $\text{COMET}_{22}^{\text{DA}}$. This suggests that $\text{COMET}_{22}^{\text{kiwiDA}}$ is more sensitive to potential changes in the degree of translationese than $\text{COMET}_{22}^{\text{DA}}$.

**Considerations.** We note the limitations of our experiments. First, we assume that the paraphrased references are as good as the original ones and less translationese, but we did not verify this when para-

phrasing the hypotheses. If the hypothesis quality has been affected, we also assume that the LLM paraphrasing process affects all system outputs in an equal magnitude. Second, the evaluations that anchor to human judgements assume that human evaluators provide assessment solely on the quality and do not overly insist on adequacy/translationese. Third, our comparison between $\text{COMET}_{22}^{\text{DA}}$ and $\text{COMET}_{22}^{\text{kiwiDA}}$ only shows that they do not behave the same in dealing with change in the degree of translationese in hypotheses.

### 3.9 Model reporting [usage]

Different COMET models can yield distinct results. Therefore it is important to always specify the specific model for sensible score interpretation and comparison. In this section, we examine to what extent this holds up in scientific literature.

We automate this bibliometric task with SemanticScholar API (Kinney et al., 2023). Starting with 1100 papers that cite one of the COMET papers, we obtain 417 papers from 2021 to 2024 that have an easily accessible PDF version.[3] We check if any of the *tables* contains the string comet. Within those papers, we check whether the COMET model information is contained in the PDF using a regular expression.[4] After further manual validation, we found that 50 of the examined papers do not report a specific COMET version. This establishes that *at least 12% of papers report COMET scores without specific model information*.

In addition, out of the almost 1000 papers running COMET in their evaluation, most only cite the first COMET paper (Rei et al., 2020a) instead of the paper that describes the specific models that are being used (Rei et al., 2020b, 2022b,a,c, 2023a,b; Glushkova et al., 2021; Wan et al., 2022; Alves et al., 2024; Guerreiro et al., 2023).

**Recommendation.** Always report the COMET version, ideally with a link. Also, cite the affiliated COMET paper as opposed to the first paper (Rei et al., 2020a), because different checkpoints have variations in training regimes that might be crucial in analysing the evaluation outcome.

---

[2]We accessed gpt-4o-2024-08-06 via API in Aug 2024.

[3]Papers in 2020 did not have to report the specific model as there was only one available at the time. Further, we acknowledge potential bias to only papers with available PDFs.

[4]Case-insensitive: "comet[ \-](da|20|21|22|23)|wmt(20|21|22|23)\-comet|xcomet\-|wmt\-da\-estimator"

### 3.10 Discussions on other issues

**Interpretation of significance.** Statistical hypothesis test merely shows how likely the difference between the average of two model's scores on the same test set is caused by random fluctuation. Kocmi et al. (2024c) shows the significance of a difference between two metric scores can be made arbitrarily high and one can force $p \to 0$ by using a sizeable test set. This tells little about whether this difference is meaningful to a human reader. For this reason, we stress the use of `mt-thresholds` that converts differences in metric scores to how perceivable they are by human annotators.[5]

**Averaging and subtracting COMET scores.** Research nowadays favours experiments on multiple translation directions, as multilingual translation models and large language models become trendy. Recent papers are more often seen to report a (macro-)average COMET score as an aggregate measure across many directions, usually in Xx→En, En→Xx, and All→All. Whilst indicative, this is not entirely scientific because (1) the score range is inherently distinct for each translation direction and (2) there is no assurance that the scores are on a linear scale. Consequently, an outlying score in a single direction can distort the average, leading to a false claim. Likewise, absolute COMET score differences are not comparable if from different base numbers or directions. We suggest that, when multiple translation directions are of interest, in addition to averaged scores, practitioners can report the number of wins (against another system) as another aggregation of individual scores.

**Optimizing to COMET.** Owing to COMET's strong correlation with human judgement, recent works investigate the feasibility of using it in translation modelling directly. These strategies either include COMET in a distillation workflow (Finkelstein and Freitag, 2024; Guttmann et al., 2024), as a data filtering method (Peter et al., 2023), as a decoding method (Freitag et al., 2022a; Fernandes et al., 2022; Vernikos and Popescu-Belis, 2024) or as a training objective (Yan et al., 2023).

Nevertheless, COMET may cease to be a good measure if practitioners over-optimize a system towards it. As Yan et al. (2023) demonstrated, a model trained towards COMET can generate "universal translations" (hallucinations) preferred by COMET regardless of the source sentence. Re-

---

cently, using COMET-based MBR decoding has become prevalent in shared tasks. For example, Unbabel-Tower70B at WMT24 also used MBR and dominated all automatic metrics but not so much under human evaluation (Kocmi et al., 2024a,b). MBR decoding could be seen as an automatic way to exploit bias in the scoring method, (currently COMET in most cases), so practitioners need to be aware of its shortcomings, disclose the use of such, and base system building on multiple (less correlated) metrics (Jon et al., 2023).

A novel issue is using automated metrics in human evaluation (Zouhar et al., 2024b) that collects data for metric training. This might create a similar effect as translationese in machine translation. The data could be biased by the particular quality estimator that is assisting annotators in the data collection process.

**Sensitivity to sentence segmentation.** Like most MT metrics, COMET works at the sentence level, but sometimes sentence-segmented input is not available. This is often the case in speech translation (ST) where sentence segmentation is treated as part of the task (Ahmad et al., 2024). To address the problem of mismatching segmentation between the system output and the reference, a common solution in ST is to re-segment the output using a minimum error rate method (Matusov et al., 2005) in order to force-align it with the reference. Forced alignment can introduce segmentation errors resulting in truncated (thus grammatically incorrect) sentences. There is evidence that COMET, as a metric reliant on sentence embeddings, is more sensitive to segmentation errors than string-based metrics, like BLEU, which rely purely on $n$-gram overlaps with no linguistic notion of a sentence (Amrhein and Haddow, 2022). In a recent comparison of COMET with human ranking, Sperber et al. (2024) suggested that COMET-based ranking is robust to segmentation errors but that a "more thorough study of this issue is needed".

**Other metrics.** Our work focused on COMET, the current most popular family of MT metrics. Nonetheless, our recommendations could apply to other neural metrics, like MetricX-23 (Juraska et al., 2023), because many issues we outlined are due to their statistical learning nature. Even beyond this, metric reporting and score interpretation in practice, e.g. software usage or averaging across multiple directions, can be problematic for string-matching metrics like BLEU or chrF too.

## 4 The `SacreCOMET` Package

To help alleviate problems in Sections 3.1 (software version), 3.2 (compute precision), and 3.9 (model reporting), we release a simple package sacreCOMET with two functionalities. Given a model name, the first functionality attempts to find the appropriate citation including a link to the paper and a BibTeX:

```
$ pip install sacrecomet
$ sacrecomet cite Unbabel/xcomet-xl

https://arxiv.org/abs/2310.10482
@misc{guerreiro2023xcomet,
 title={xCOMET: Transparent Machine
     Translation Evaluation through Fine-
     grained Error Detection},
 ...
```

The second functionality semi-automatically detects the local software versions to generate a signature for better reproducibility. Both functionalities can also be run in interactive mode.

```
$ sacrecomet --model unite-mup --prec fp32

Python3.11.8|Comet2.2.2|fp32|unite-mup
```

## 5 Future Work on Learned Metrics

- **Fixing data bias**: Learned metrics are sensitive to the training data distribution. Future metrics should aim to reduce the bias caused by the data selection process such that they are applicable to a range of MT systems.

- **Interpretability across languages**: Currently, practitioners cannot compare, or pedantically, aggregate scores in different translation directions. It would be useful to unify the scores to a single scale that can be interpreted independent of the language (similar to Kocmi et al., 2024c), e.g. to indicate X% of segments are production-ready.

- **Confidence-aware metrics:** As seen in works of Glushkova et al. (2021); Fomicheva et al. (2020b), it is possible to build metrics that output a confidence interval, though its usage in evaluation and modelling remains scarce.

- **Inference speed:** Learned metrics are getting better but at the cost of bulky models and increased inference time. These overheads should be taken into account when developing new models, such as the work of Rei et al. (2022b).

- **Representations:** Current COMET models are built upon off-the-shelf multilingual encoder models which are likely trained on human-written texts. However, this could bring in a domain mismatch—when translation hypotheses act as the input to metric models, they are not human-written but machine-translated.

- **Robustness**: Metrics should have the correct behaviour even in corner cases, be it empty output or incorrect language. Mapping all inputs (Amrhein et al., 2022, *inter alia)*, including partial or adversarial ones, evaluating the metrics, and coming up with methods to make them more robust would increase the metrics' adoption and trust.

- **Built-in QE:** In production, machine translation and quality estimation are commonly two different processes. In many applications, however, a single QE model is used for a single MT model. Quantifying how much is QE adaptation to a particular MT model useful is beneficial for a holistic understanding of QE metrics. Further, proposing methods for supervised quality estimation built into the MT could ease industry adoption. Beyond the work of Tomani et al. (2024), this remains largely unexplored.

- **Noise-aware training:** Human annotations are notoriously noisy. At the scale of WMT data, poor-quality annotations are unavoidable. The inter-annotator agreement for even robust annotations, such as ESA, remains low at $\tau_c \approx 0.3$. The effect of data quality on learned metrics is so far unknown and methods for noise/uncertainty-aware training are under-studied.

## 6 Conclusion

COMET is currently one of the most powerful automatic metrics/quality estimators for machine translation, consistently achieving the top correlation with human judgement. In comparison to previous metrics, it is a statistical learning model and thus inherits all the related problems in addition to possible technical misconfigurations. We urge practitioners to consider more deeply the use of COMET in non-standard scenarios especially where such training bias might come into play. Beyond these issues, there has been confusion in the literature in reporting the correct COMET model and its correct setting. For improved consistency, we release an easy-to-use tool to assist practitioners.

## Acknowledgments

## References

Ibrahim Said Ahmad, Antonios Anastasopoulos, Ondřej Bojar, Claudia Borg, Marine Carpuat, Roldano Cattoni, Mauro Cettolo, William Chen, Qianqian Dong, Marcello Federico, Barry Haddow, Dávid Javorský, Mateusz Krubiński, Tsz Kin Lam, Xutai Ma, Prashant Mathur, Evgeny Matusov, Chandresh Maurya, John McCrae, Kenton Murray, Satoshi Nakamura, Matteo Negri, Jan Niehues, Xing Niu, Atul Kr. Ojha, John Ortega, Sara Papi, Peter Polák, Adam Pospíšil, Pavel Pecina, Elizabeth Salesky, Nivedita Sethiya, Balaram Sarkar, Jiatong Shi, Claytone Sikasote, Matthias Sperber, Sebastian Stüker, Katsuhito Sudoh, Brian Thompson, Alex Waibel, Shinji Watanabe, Patrick Wilken, Petr Zemánek, and Rodolfo Zevallos. 2024. Findings of the IWSLT 2024 evaluation campaign. In *Proceedings of the 21st International Conference on Spoken Language Translation (IWSLT 2024)*.

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *Preprint*, arXiv:2402.17733.

Chantal Amrhein and Barry Haddow. 2022. Don't discard fixed-window audio segmentation in speech-to-text translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*.

Chantal Amrhein, Nikita Moghe, and Liane Guillou. 2022. ACES: Translation accuracy challenge sets for evaluating machine translation metrics. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*.

Chantal Amrhein and Rico Sennrich. 2022. Identifying weaknesses in machine translation metrics through minimum Bayes risk decoding: A case study for COMET. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*.

Rachel Bawden, Biao Zhang, Andre Tättar, and Matt Post. 2020. ParBLEU: Augmenting metrics with automatic paraphrases for the WMT'20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

Mara Finkelstein and Markus Freitag. 2024. MBR and QE finetuning: Training-time distillation of the best and most expensive decoding methods. In *The Twelfth International Conference on Learning Representations*.

Marina Fomicheva, Lucia Specia, and Francisco Guzmán. 2020a. Multi-hypothesis machine translation evaluation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020b. Unsupervised Quality Estimation for Neural Machine Translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.

Markus Freitag, David Grangier, and Isaac Caswell. 2020. BLEU might be guilty but references are not innocent. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022a. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022b. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej

Bojar. 2021. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*.

Martin Gellerstam. 1986. Translationese in swedish novels translated from english. In L. Wollin and H. Lindquist, editors, *Translation studies in Scandinavia: Poceedings from the Scandinavian Symposium on Translation Theory (SSOTT) II*, number 75 in Lund Studies in English, page 88–95. CWK Gleerup, Lund.

Taisiya Glushkova, Chrysoula Zerva, Ricardo Rei, and André F. T. Martins. 2021. Uncertainty-aware machine translation evaluation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*.

Thamme Gowda, Roman Grundkiewicz, Elijah Rippeth, Matt Post, and Marcin Junczys-Dowmunt. 2024. Pymarian: Fast neural machine translation and evaluation in python. *Preprint*, arXiv:2408.11853.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *Preprint*, arXiv:2310.10482.

Kamil Guttmann, Mikołaj Pokrywka, Adrian Charkiewicz, and Artur Nowakowski. 2024. Chasing COMET: Leveraging minimum bayes risk decoding for self-improving machine translation. *arXiv preprint arXiv:2405.11937*.

Josef Jon, Martin Popel, and Ondřej Bojar. 2023. CUNI at WMT23 general translation task: MT and a genetic algorithm. In *Proceedings of the Eighth Conference on Machine Translation*.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*.

Rodney Michael Kinney, Chloe Anastasiades, Russell Authur, Iz Beltagy, Jonathan Bragg, Alexandra Buraczynski, Isabel Cachola, Stefan Candra, Yoganand Chandrasekhar, Arman Cohan, Miles Crawford, Doug Downey, Jason Dunkelberger, Oren Etzioni, Rob Evans, Sergey Feldman, Joseph Gorney, David W. Graham, F.Q. Hu, Regan Huff, Daniel King, Sebastian Kohlmeier, Bailey Kuehl, Michael Langan, Daniel Lin, Haokun Liu, Kyle Lo, Jaron Lochner, Kelsey MacMillan, Tyler C. Murray, Christopher Newell, Smita R Rao, Shaurya Rohatgi, Paul Sayre, Zejiang Shen, Amanpreet Singh, Luca Soldaini, Shivashankar Subramanian, A. Tanaka, Alex D Wade, Linda M. Wagner, Lucy Lu Wang, Christopher Wilhelm, Caroline Wu, Jiangjiang Yang, Angele Zamarron, Madeleine van Zuylen, and Daniel S. Weld. 2023. The semantic scholar open data platform. *ArXiv*, abs/2301.10140.

Rebecca Knowles. 2021. On the stability of system rankings at WMT. In *Proceedings of the Sixth Conference on Machine Translation*.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024a. Preliminary wmt24 ranking of general mt systems and llms. *Preprint*, arXiv:2407.19884.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024b. Findings of the WMT24 general machine translation shared task: The llm era is here but MT is not solved yet. To be published at WMT 2024.

Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In *Proceedings of the Sixth Conference on Machine Translation*.

Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024c. Navigating the metrics maze: Reconciling score magnitudes and accuracies. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Tom Kocmi, Vilém Zouhar, Eleftherios Avramidis, Roman Grundkiewicz, Marzena Karpinska, Maja Popović, Mrinmaya Sachan, and Mariya Shmatova. 2024d. Error span annotation: A balanced approach for human evaluation of machine translation. *Preprint*, arXiv:2406.11580.

Philipp Koehn and Christof Monz. 2006. Manual and automatic evaluation of machine translation between European languages. In *Proceedings on the Workshop on Statistical Machine Translation*.

Shankar Kumar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*.

Chi-kiu Lo, Samuel Larkin, and Rebecca Knowles. 2023. Metric score landscape challenge (MSLC23): Understanding metrics' performance on a wider landscape of translation quality. In *Proceedings of the Eighth Conference on Machine Translation*.

Arle Lommel, Aljoscha Burchardt, Maja Popović, Kim Harris, Eleftherios Avramidis, and Hans Uszkoreit. 2014. Using a new analytic measure for the annotation and analysis of MT errors on real data. In *Proceedings of the 17th Annual Conference of the European Association for Machine Translation*.

E. Matusov, G. Leusch, O. Bender, and H. Ney. 2005. Evaluating Machine Translation Output with Automatic Sentence Segmentation. In *Proceedings of the 2nd International Workshop on Spoken Language Translation (IWSLT)*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*.

Jan-Thorsten Peter, David Vilar, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, and Markus Freitag. 2023. There's no data like better data: Using QE metrics for MT data filtering. In *Proceedings of the Eighth Conference on Machine Translation*.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*.

Ricardo Rei, Ana C Farinha, José G.C. de Souza, Pedro G. Ramos, André F.T. Martins, Luisa Coheur, and Alon Lavie. 2022b. Searching for COMETINHO: The little metric that could. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André F. T. Martins. 2023a. Scaling up cometkiwi: Unbabel-ist 2023 submission for the quality estimation shared task. *Preprint*, arXiv:2309.11925.

Ricardo Rei, Nuno M. Guerreiro, Marcos Treviso, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2023b. The inside story: Towards better understanding of machine translation neural evaluation metrics. *Preprint*, arXiv:2305.11806.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020a. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020b. Unbabel's participation in the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022c. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Hiroki Shimanaka, Tomoyuki Kajiwara, and Mamoru Komachi. 2018. RUSE: Regressor using sentence embeddings for automatic machine translation evaluation. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*.

Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*.

Matthias Sperber, Ondřej Bojar, Barry Haddow, Dávid Javorský, Xutai Ma, Matteo Negri, Jan Niehues, Peter Polák, Elizabeth Salesky, Katsuhito Sudoh, and Marco Turchi. 2024. Evaluating the IWSLT2023 speech translation tasks: Human annotations, automatic metrics, and segmentation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*.

Shuo Sun, Francisco Guzmán, and Lucia Specia. 2020. Are we estimating or guesstimating translation quality? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Kosuke Takahashi, Katsuhito Sudoh, and Satoshi Nakamura. 2020. Automatic machine translation evaluation using source language inputs and cross-lingual language model. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Christian Tomani, David Vilar, Markus Freitag, Colin Cherry, Subhajit Naskar, Mara Finkelstein, Xavier Garcia, and Daniel Cremers. 2024. Quality-aware translation models: Efficient generation and quality estimation in a single model. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Giorgos Vernikos and Andrei Popescu-Belis. 2024. Don't rank, combine! combining machine translation hypotheses using quality estimation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Yu Wan, Dayiheng Liu, Baosong Yang, Haibo Zhang, Boxing Chen, Derek Wong, and Lidia Chao. 2022. UniTE: Unified translation evaluation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Yiming Yan, Tao Wang, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Mingxuan Wang. 2023. BLEURT has universal translations: An analysis of automatic metrics by minimum risk training. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: a case study. In *Proceedings of the 40th International Conference on Machine Learning*.

Renjie Zheng, Mingbo Ma, and Liang Huang. 2018. Multi-reference training with pseudo-references for neural translation and text generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.

Vilém Zouhar and Ondřej Bojar. 2024. Quality and quantity of machine translation references for automatic metrics. In *Proceedings of the Fourth Workshop on Human Evaluation of NLP Systems (HumEval) @ LREC-COLING 2024*.

Vilém Zouhar, Shehzaad Dhuliawala, Wangchunshu Zhou, Nico Daheim, Tom Kocmi, Yuchen Eleanor Jiang, and Mrinmaya Sachan. 2023. Poor man's quality estimation: Predicting reference-based MT metrics without the reference. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*.

Vilém Zouhar, Shuoyang Ding, Anna Currey, Tatyana Badeka, Jenyuan Wang, and Brian Thompson. 2024a. Fine-tuned machine translation metrics struggle in unseen domains. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*.

Vilém Zouhar, Tom Kocmi, and Mrinmaya Sachan. 2024b. AI-assisted human evaluation of machine translation. *Preprint*, arXiv:2406.12419.

# A Batch size and GPU/CPU

We run the test inference on a combination of GPU or CPU with varying batch sizes (BS, 1 or 100). Results in Table 11 demonstrate that the tiny effects of these choices are negligible for COMET reporting.

| Difference | MAE |
|---|---|
| BS=1, GPU  – BS=1, GPU | 0 |
| BS=1, GPU  – BS=64 GPU | $2 \times 10^{-7}$ |
| BS=1, GPU  – BS=1, CPU | $4 \times 10^{-7}$ |
| BS=64, GPU – BS=64, CPU | $4 \times 10^{-7}$ |

Table 11: MAE between segment-level $COMET_{22}^{DA}$ scores under various inference settings. The "BS=1, GPU" setting in the first row was run twice.

# B Distribution of COMET scores for empty and valid hypothesis



Figure 4: Distribution of instance-level scores for empty and baseline translations (x-axis: score; y-axis: count).

## C System rankings before and after paraphrasing references or hypotheses

| WMT20 En→De | COMET$_{20}^{DA}$ | | Ranking | |
|---|---|---|---|---|
| | original ref | paraphrased ref | original ref | paraphrased ref |
| Sys-1069 | 0.508 | 0.305 | 9 | 9 |
| Sys-832 | 0.540 | 0.333 | 8 | 8 |
| Sys-1535 | 0.560 | 0.356 | 5 | 5 |
| Online-A | 0.499 | 0.288 | 10 | 10 |
| Online-B | 0.554 | 0.351 | **7** | **6** |
| Online-G | 0.268 | 0.052 | 14 | 14 |
| Online-Z | 0.329 | 0.121 | 13 | 13 |
| Sys-73 | 0.405 | 0.192 | 12 | 12 |
| Sys-1520 | 0.563 | 0.360 | 4 | 4 |
| Sys-890 | 0.578 | 0.371 | 3 | 3 |
| Sys-1136 | 0.472 | 0.264 | 11 | 11 |
| Sys-388 | 0.102 | -0.074 | 16 | 16 |
| Sys-737 | 0.555 | 0.351 | **6** | **7** |
| Ref-A | 0.878 | 0.526 | 1 | 1 |
| Ref-B | 0.591 | 0.446 | 2 | 2 |
| Sys-179 | 0.189 | 0.000 | 15 | 15 |
| | | | $\tau_c = 0.9833$ | |
| | | | Acc = 0.924 | Acc = 0.924 |

Table 12: Results for WMT20 En→De submissions evaluated against the original or human-paraphrased reference. Kendall's $\tau_c$ is measured between two evaluations based on original and paraphrased references; pairwise system-level accuracy (Acc) is measured against human DA scores.

| WMT19 En→De | COMET$_{22}^{DA}$ | | Ranking | |
|---|---|---|---|---|
| | original ref | paraphrased ref | original ref | paraphrased ref |
| Sys-6862 | 0.867 | 0.817 | 2 | 2 |
| Sys-6820 | 0.834 | 0.780 | 12 | 12 |
| Sys-6819 | 0.843 | 0.789 | 9 | 9 |
| Sys-6651 | 0.847 | 0.794 | 7 | 7 |
| Sys-6926 | 0.852 | 0.797 | **4** | **5** |
| Sys-6808 | 0.869 | 0.818 | 1 | 1 |
| Sys-6785 | 0.837 | 0.785 | 11 | 11 |
| Sys-6974 | 0.866 | 0.814 | 3 | 3 |
| Sys-6763 | 0.851 | 0.797 | **5** | **6** |
| Sys-6674 | 0.811 | 0.756 | 17 | 17 |
| Sys-6508 | 0.804 | 0.752 | 19 | 19 |
| Sys-6731 | 0.850 | 0.797 | **6** | **4** |
| Sys-6871 | 0.809 | 0.756 | 18 | 18 |
| Sys-6479 | 0.833 | 0.777 | 13 | 13 |
| Sys-6823 | 0.845 | 0.792 | 8 | 8 |
| Sys-6790 | 0.386 | 0.364 | 22 | 22 |
| Sys-6981 | 0.826 | 0.774 | 14 | 14 |
| Online-A | 0.815 | 0.761 | 16 | 16 |
| Online-B | 0.838 | 0.784 | 10 | 10 |
| Online-G | 0.795 | 0.738 | 20 | 20 |
| Online-X | 0.728 | 0.673 | 21 | 21 |
| Online-Y | 0.821 | 0.762 | 15 | 15 |
| | | | $\tau_c = 0.9827$ | |
| | | | Acc = 0.875 | Acc = 0.845 |

Table 13: Results for WMT19 En→De submissions evaluated against the original or human-paraphrased reference. Kendall's $\tau_c$ is measured between two evaluations based on original and paraphrased references; pairwise system-level accuracy (Acc) is measured against human DA scores.

| WMT19 En→De | COMET$_{22}^{DA}$ | | Ranking | | COMET$_{22}^{kiwiDA}$ | | Ranking | |
|---|---|---|---|---|---|---|---|---|
| | orig. hyp | para. hyp | orig. hyp | para. hyp | orig. hyp | para. hyp | orig. hyp | para. hyp |
| Sys-6823 | 0.845 | 0.840 | **8** | **5** | 0.821 | 0.824 | **8** | **4** |
| Sys-6862 | 0.867 | 0.842 | 2 | 2 | 0.840 | 0.828 | **2** | **3** |
| Sys-6819 | 0.843 | 0.835 | 9 | 9 | 0.815 | 0.817 | **9** | **8** |
| Sys-6808 | 0.869 | 0.843 | 1 | 1 | 0.840 | 0.828 | **1** | **2** |
| Sys-6974 | 0.866 | 0.842 | 3 | 3 | 0.838 | 0.829 | **3** | **1** |
| Sys-6651 | 0.847 | 0.835 | **7** | **8** | 0.821 | 0.817 | **7** | **9** |
| Sys-6926 | 0.852 | 0.839 | **4** | **6** | 0.824 | 0.822 | **4** | **6** |
| Sys-6763 | 0.851 | 0.841 | **5** | **4** | 0.823 | 0.824 | 5 | 5 |
| Online-B | 0.838 | 0.830 | 10 | 10 | 0.805 | 0.810 | 10 | 10 |
| Sys-6731 | 0.850 | 0.837 | **6** | **7** | 0.822 | 0.820 | **6** | **7** |
| | | | $\tau_c = 0.822$ | | | | $\tau_c = 0.644$ | |
| | | | Acc = 0.911 | | | | Acc = 0.822 | |

Table 14: Results for WMT19 En→De system outputs and LLM-paraphrased outputs evaluated against the original reference. Both Kendall's $\tau_c$ and pairwise system-level accuracy (Acc) are measured between two evaluations based on original and paraphrased references.

# Post-edits Are Preferences Too

**Nathaniel Berger**[a] **and Stefan Riezler**[ab]
Computational Linguistics[a] & IWR[b]
Heidelberg University
69120 Heidelberg, Germany
berger@cl.uni-heidelberg.de
riezler@cl.uni-heidelberg.de

**Miriam Exel**[c] **and Matthias Huck**[c]
SAP SE[c]
Dietmar-Hopp-Allee 16
69190 Walldorf, Germany
miriam.exel@sap.com
matthias.huck@sap.com

## Abstract

Preference Optimization (PO) techniques are currently one of the state of the art techniques for fine-tuning large language models (LLMs) on pairwise preference feedback from human annotators. However, in machine translation, this sort of feedback can be difficult to solicit. Additionally, Kreutzer et al. (2018) have shown that, for machine translation, pairwise preferences are less reliable than other forms of human feedback, such as 5-point ratings.

We examine post-edits to see if they can be a source of reliable human preferences by construction. In PO, a human annotator is shown sequences $s_1$ and $s_2$ and asked for a preference judgment, while for post-editing, editors *create* $s_1$ and know that it should be better than $s_2$. We attempt to use these implicit preferences for PO and show that it helps the model move towards post-edit-like hypotheses and away from machine translation-like hypotheses. Furthermore, we show that best results are obtained by pre-training the model with supervised fine-tuning (SFT) on post-edits in order to promote post-edit-like hypotheses to the top output ranks.

## 1 Introduction

The current state of the art methods for training large language models offline on human preference data are Direct Preference Optimization (DPO) (Rafailov et al., 2023) or Identity Preference Optimization (IPO) (Gheshlaghi Azar et al., 2024). Instead of training a separate reward model and then performing reinforcement learning, these methods train directly on the collected preference data by deriving a directly optimizable loss function from the preference model.

However, in some domains, the pairwise preference annotations required for using these methods have been found to be less reliable than other annotation schemes. Kreutzer et al. (2018) find that inter-rater reliability for pairwise ranking of



Figure 1: The generative process for preference optimization is that two sequences $s_1$ and $s_2$ are given, and a preference judgment $s_1 > s_2$ is generated (upper graph). The data generating process of post-editing yields reliable preferences by construction: Given $s_2$ and the implicit preference that $s_1 > s_2$, create $s_1$ (lower graph). We propose using the implicit preferences from post-editing for preference optimization.

machine translation outputs to be less than that of 5-point rating. In the field of translation, there are many different dimensions on which one translation may be better than another, e.g. fluency, faithfulness, formality, terminology, etc. (Lommel et al., 2013). This poses a problem for human annotators when they are presented with two plausible translations.

We propose using the data generated by post-editing to yield reliable preferences by construction. The current generative process for preference data is that two sequences $s_1$ and $s_2$ are given, and a preference judgment $s_1 > s_2$ is sought, yielding the generative process $s_1 \rightarrow s_1 > s_2 \leftarrow s_2$. We propose using data generated by the following process: Given $s_2$ and the implicit preference that $s_1 > s_2$, create $s_1$, yielding the generative process $s_1 > s_2 \rightarrow s_1 \leftarrow s_2$ (see Fig. 1).

Post-editing is already a common practice in the translation community to clean up raw-MT outputs before publishing. Post-editors create new sequences that they prefer with regards to the qual-

ity expected in their domain. Typically, the original raw-MT output is discarded and the post-edit is published. If this data is used for training, the post-edit is treated as a new reference for supervised fine-tuning (SFT). This ignores, however, the fact that the post-edit is not just a new reference translation but also a quality judgment of what in the raw-MT was erroneous. Using PO objectives allows us to fine-tune an LLM to translate in a way that is more in line with the post-editors' implicit preferences. However, PO does not necessarily promote the preferred sequence to become the argmax output of the model, but rather re-ranks sequences within the model's probability space. If the two sequences are both unlikely under the model's output distribution, they will remain unlikely but their relative probability will respect the preferences. We show that best results are obtained by pre-training the model on post-edits with SFT, promoting post-edits to the top ranks, followed by fine-tuning with a PO loss. This combined training teaches the model to prefer and promote post-edits such that reference-like translations are produced but also dispreferred machine-translation-like hypotheses are avoided.

## 2 Related Work

Kreutzer et al. (2018) gather human feedback on machine translation outputs in the form of 5-point ratings and as pairwise preferences. They then use this feedback to train two reward models, one that is trained on the 5-point ratings and is trained with a regression loss to directly predict a reward value and one that is trained on pairwise preferences by fitting a Bradley-Terry model (Bradley and Terry, 1952) to the preferences as had been done by Christiano et al. (2017). These reward models are then used to train machine translation models. Kreutzer et al. (2018) find that ratings are more reliable than rankings and that reinforcement learning with a ratings-trained reward estimator yields better results than using rankings-trained reward estimates.

Berger et al. (2023) fine-tune a pre-trained NMT model on post-editing data by presenting the model with both the post-edit and the current MT hypothesis. At each epoch, the NMT model being trained generates translations for all training data. These generated outputs are then compared to the original post-edits with a token-level diff. Both sequences are then used as training exam-

ples for the NMT system. However, tokens that appear in the hypothesis but not the post-edit are given a negative weight in the loss function. On examples where the two sequences differ, the model gets both negative feedback, where the probability of that token is to be decreased, and positive feedback, where the probability should be increased.

Xu et al. (2024b) similarly present the model with a positive and negative example of machine translation outputs during training but use a modified version of the DPO (Rafailov et al., 2023) loss to optimize it. Their change to DPO adds an SFT term. The SFT term promotes the preferred sequence to be the argmax output of the model while the DPO part of the loss establishes the distance between the two sequences in log-probability space. The MT hypotheses that they generate come from two different LLMs; ALMA-13B-LoRA (Xu et al., 2024a) and GPT-4 (OpenAI et al., 2024). Additionally, they use reference translations from the original dataset. The preferences that they use are predicted by open-source quality estimation models KIWI-XXL (Rei et al., 2023) and XCOMET (Guerreiro et al., 2023).

## 3 Preference Optimization

### 3.1 Background

Using reinforcement learning with human feedback (RLHF) has recently re-emerged as a method for training LLMs to generate outputs that are preferred by human annotators (Ziegler et al. (2019), Ouyang et al. (2022), Bai et al. (2022), inter alia) without requiring handwritten demonstrations of preferred behavior which would be required for supervised fine-tuning (SFT). The general recipe is as follows: pre-train an LLM on in-domain data; generate multiple completions $y$ for a single input $x$ (or prompt); have human annotators rank or rate the completions; train a reward model to predict rankings or ratings given inputs and completions; use the trained reward model to predict rewards for reinforcement learning, frequently with proximal policy optimization (Schulman et al., 2017). Training a separate model to predict rewards for reinforcement learning is known as an actor-critic method.

The reward model in the previous works is structured as a Bradley-Terry model (Bradley and Terry, 1952), where the probability of preferring $y_1$ over $y_2$ is given by

$$p(y_1 \succ y_2|x) = \sigma(r_\theta(x, y_1) - r_\theta(x, y_2))$$

where $\sigma$ is the logistic function and $r_\theta$ is the reward model that is trained on the pairwise rankings to give the preferred sequence a higher value. The reward model can then be used to estimate rewards for outputs sampled during online training.

This process requires training an additional model and hiring human annotators to perform ranking. DPO (Rafailov et al., 2023) is a technique that obviates the need for a secondary model by instead giving the model both the preferred and dispreferred sequences and optimizing a distance between the two sequence in log-probability space.

If handwritten demonstrations of preferred sequences are available, then SFT would typically be performed. The goal of SFT is to maximize the probability of the demonstrations under the model. For text generation, this is done by minimizing the negative log-probability of each token given all previous tokens in the sequence.

$$\mathcal{L}_{SFT}(y) = -\sum_{i=0}^{|y|} \log(\pi(y_i|y_{0:i-1}))$$

Minimizing this loss promotes the sequence $y$ to be the argmax output of the model, while reinforcement learning increases or decreases the probability of a sequence with regard to the magnitude of its reward.

### 3.2 PO Objectives

The DPO loss (Rafailov et al., 2023) is based on the Bradley-Terry model of human preferences but, unlike actor-critic reinforcement learning techniques, it does not train a separate reward model. Instead they rewrite the reward function $r$ in terms of the optimal policy and the baseline model. They notice that the theoretically optimal policy $\pi_r$, with a KL-divergence constraint, is equal to the baseline model with its output distribution re-weighted according to the reward function

$$\pi_r(y|x) = \frac{1}{Z(x)} \pi_{ref}(y|x) \exp\left(\frac{1}{\beta} r(x,y)\right)$$

where $Z$ is the partition function, which normalizes the function to be a proper probability distribution. This formula can also be solved for the reward function $r$, such that rewards are expressed as the difference between two models' probability ratios. If this $r$ is then inserted back into the

Bradley-Terry model, it becomes

$$p(y_w \succ y_l|x) =$$
$$\sigma\left(\beta \log\left(\frac{\pi^*(y_w|x)}{\pi_{ref}(y_w|x)}\right) - \beta \log\left(\frac{\pi^*(y_l|x)}{\pi_{ref}(y_l|x)}\right)\right)$$

where $y_w$ and $y_l$ denote the preferred and dispreferred completion, respectively, and $\pi^*$ is now the model we are training to be optimal under the reward function. This probability, $p(y_w \succ y_l|x)$ can be optimized by minimizing the negative log-probability. With regards to output probability, the loss monotonically decreases as $y_w$ becomes more probable than $y_l$. Increasing the difference between the two always decreases the loss.

Gheshlaghi Azar et al. (2024) re-derive a similar loss with some theoretical advantages. Instead of training the model to be optimal under the Bradley-Terry derived reward function, they train the model to separate the two outputs by a fixed difference in log-probability space.

$$\mathcal{L}_{IPO}(y_w, y_l, x) =$$
$$+\left(\left(\log\left(\frac{\pi^*(y_w|x)}{\pi_{ref}(y_w|x)}\right) - \log\left(\frac{\pi^*(y_l|x)}{\pi_{ref}(y_l|x)}\right)\right) - \frac{1}{2\beta}\right)^2$$

Because this loss function is minimized when the log-probability ratio difference is exactly $(2\beta)^{-1}$, and will increase when the outputs move further apart in log-probability space, the authors claim an advantage for deterministic preferences, where the same preferences are seen multiple times during training. Because the preferences that we use are deterministic, we opt for the IPO paradigm of PO.

A follow-up work to DPO, focused specifically on machine translation, additively combines the DPO loss and the SFT loss (Xu et al., 2024b), which the authors call Contrastive Preference Optimization (CPO). Additionally, they perform an ad-hoc modification of the DPO loss by dropping the normalizer $\pi_{ref}$ in the DPO loss so as to not perform a second forward pass on the reference model.

$$\mathcal{L}_{CPO}(y_w, y_l, x) = -\log(\pi^*(y_w|x))$$
$$-\log(\sigma(\beta \log(\pi^*(y_w|x)) - \beta \log(\pi^*(y_l|x))))$$

We use a reformulation of the CPO loss with the IPO training objective for our experiments because our preferences are deterministic. Additionally, we keep the normalizers in the IPO loss because these can be pre-computed in advance, instead of in a second forward pass, and incur only a negligible memory and speed penalty.

Our modified variant of the CPO loss, which we call dDPO for the deterministic preferences involved in post-editing, is

$$
\mathcal{L}_{dCPO}(y_w, y_l, x) = -\log(\pi^*(y_w|x))
$$
$$
+\left(\left(\log\left(\frac{\pi^*(y_w|x)}{\pi_{ref}(y_w|x)}\right) - \log\left(\frac{\pi^*(y_l|x)}{\pi_{ref}(y_l|x)}\right)\right) - \frac{1}{2\beta}\right)^2
$$

which is the SFT objective added to the IPO objective. When we refer to dCPO later in this paper, we are referring to this modified version of the CPO objective.

## 4 Experiments

| Data | Split | BLEU | TER | CHRF |
|------|-------|------|-----|------|
| En→DE | Train | 49.4 | 37.6 | 71.6 |
| | Dev | 50.9 | 36.5 | 72.5 |
| | Test | 50.8 | 36.4 | 72.8 |
| En→Ru | Train | 80.9 | 13.6 | 89.9 |
| | Dev | 80.2 | 14.9 | 89.0 |
| | Test | 76.3 | 17.4 | 87.2 |

Table 1: Token level metrics comparing the WMT APE datasets' machine translations to the post-edits.

We fine-tune an LLM for the task of machine translation under five different conditions: SFT, IPO, dCPO, and pre-training with SFT followed by either IPO or dCPO, denoted as SFT→IPO and SFT→dCPO, respectively, and evaluate with the neural metrics XCOMET (Guerreiro et al., 2023) and MetricX (Juraska et al., 2023). Rafailov et al. (2023) pre-train their large language models on in-domain data such that they are already able to perform the requested task to begin with. For the task single-turn dialogue, they use the Anthropic Helpful and Harmless dialogue dataset but because no pre-training data is available, they perform SFT on the helpful answers as a pre-training step. This is similar to our conditions SFT→IPO and SFT→dCPO.

The LLM that we choose to fine-tune is Tower-Base by Alves et al. (2024). We make this choice because it is a multi-lingual LLM pre-trained on all languages we intend to work with and because the size of the model is still small enough to perform a full fine-tune with our resources[1]. We opt for Tower-Base instead of Tower-Instruct

---

[1]We train on a server with 4x Nvidia A40 GPUs with 48 GB of memory each. The system contains 256GB of RAM and 64 CPU cores.

because Tower-Instruct has been instruction fine-tuned for various down-stream tasks and not just for machine translation. Using Tower-Base instead allows us to perform a SFT step on our own. We fine-tune in all scenarios with a minimal prompt "Translate English to German.\nEnglish: {Source}\nGerman:" for our German examples. Our Russian examples use a prompt with the language name changed. This prompt is used for both SFT and PO training objectives.

Our post-edits come from WMT Automatic Post-Editing (APE) shared tasks of previous years. These datasets contain triples of source, machine-translation (MT), and post-edit (PE). We focus on the language pairs En→De from 2020 and En→Ru from 2019. The En→De source data comes from Wikipedia and is translated by a black-box NMT system (Chatterjee et al., 2020). The En→Ru data comes from the information technology domain from Microsoft Office localization work and was translated by Microsoft's production NMT system (Chatterjee et al., 2019). The En→Ru data contains base64 encoded data and sequences long enough to cause out of memory errors. We therefore filter out sequences with fewer than 4 tokens, more than 128 tokens, or more than 500 characters from the En→Ru training data, leaving 9290 (source, mt, pe) triples for training. The En→De training data was already clean and all 7000 (source, mt, pe) triples were kept for training.

Table 1 shows the performance of the datasets' machine translations when compared to their post-edits in terms of token based metrics, BLEU, TER, and CHRF (Papineni et al. (2002), Snover et al. (2006), and Popović (2015)). We see that more edits were made to the German machine translations compared to the Russian machine translations. The Russian data has far more unedited sequences—of the 9290 examples we have in our Russian training data after filtering, 5263 are unedited or 56.7%. To compare, of the 7000 German training examples, 448 are unedited or just 6.4%. We keep the unedited data for training as the SFT and dCPO objectives can still take advantage of unedited data, but filter it out for our analysis later as it is impossible for a model to prefer an unedited "post-edit" over the machine translation.

We train with fully-sharded data parallelism (FSDP) in PyTorch using Accelerate (Gugger et al., 2022) across four GPUs with an effective batch size of 256 sequences. When possi-

**Without References**

| | Model | En→DE | | En→Ru | |
|---|---|---|---|---|---|
| | | XCOMET-XL | XCOMET-XXL | XCOMET-XL | XCOMET-XXL |
| a | APE Data MT | 92.78 | 94.47 | $93.07^{ce}$ | $91.35^{ce}$ |
| b | APE Data PE | $95.55^{ac}$ | $97.01^{ac}$ | $95.29^{ace}$ | $93.78^{ace}$ |
| c | Tower Base | $94.33^{a}$ | 94.75 | 85.50 | 65.07 |
| d | SFT | $95.63^{ac}$ | $97.01^{ac}$ | $95.29^{ace}$ | $93.55^{ace}$ |
| e | IPO | $95.87^{ac}$ | $97.18^{ac}$ | $89.65^{c}$ | $72.90^{c}$ |
| f | dCPO | $95.67^{ac}$ | $97.51^{abcde}$ | $95.55^{ace}$ | $93.73^{ace}$ |
| g | SFT→IPO | $95.87^{abcd}$ | $97.48^{abcde}$ | $95.62^{acde}$ | $94.40^{abcdef}$ |
| h | SFT→dCPO | $95.91^{abcdef}$ | $97.57^{abcde}$ | $95.85^{abcdefg}$ | $94.76^{abcdefg}$ |

**With References**

| | Model | En→DE | | En→Ru | |
|---|---|---|---|---|---|
| | | XCOMET-XL | XCOMET-XXL | XCOMET-XL | XCOMET-XXL |
| a | APE Data MT | 92.80 | 94.20 | $94.99^{bd}$ | $92.68^{bd}$ |
| b | Tower Base | $93.90^{a}$ | 94.44 | 83.65 | 65.48 |
| c | SFT | $95.57^{ab}$ | $96.77^{ab}$ | $95.36^{bd}$ | $93.30^{abd}$ |
| d | IPO | $95.56^{ab}$ | $96.85^{ab}$ | $88.15^{b}$ | $72.64^{b}$ |
| e | dCPO | $95.67^{ab}$ | $97.20^{abcd}$ | $95.36^{bd}$ | $93.06^{bd}$ |
| f | SFT→IPO | $95.94^{abcde}$ | $97.31^{abcd}$ | $95.77^{abcde}$ | $93.91^{abcde}$ |
| g | SFT→dCPO | $96.00^{abcde}$ | $97.36^{abcd}$ | $96.11^{abcdef}$ | $94.14^{abcde}$ |

Table 2: XCOMET-XL and -XXL on the WMT 2020 En->DE and 2019 En->Ru test sets. Higher values are better. Superscripts indicate which system the given line is significantly better than with $\alpha < 0.05$ according to pair-wise bootstrap resampling. We see that initializing with an SFT model and then performing PO yields the best results.

| Model | En→De | | En→Ru | |
|---|---|---|---|---|
| | w Ref | w/o Ref | w Ref | w/o Ref |
| Tower Base | 1.1396 | 1.4224 | 4.1858 | 8.1576 |
| SFT | 0.9174 | 1.1757 | 1.2706 | 1.4246 |
| IPO | 0.8240 | 0.9484 | 3.0420 | 5.3263 |
| dCPO | 0.8286 | 1.0092 | 1.2873 | 1.4575 |
| SFT→IPO | 0.7985 | 0.9476 | 1.1554 | 1.3335 |
| SFT→dCPO | 0.7978 | 0.9558 | 1.1110 | 1.2607 |

Table 3: MetricX 23 XL results with and without references on the WMT 2020 En->DE and 2019 En->Ru test sets. Lower values are better. Results appear in line with XCOMET-XL and -XXL and reinforce our previous results.

ble, we shared hyperparameters across all runs and datasets. For example, both IPO and dCPO have $\beta$ set to 0.1. Full hyper-parameters can be found in the Appendix A. We used reference-free XCOMET-XL as an early stopping criterion, which was run at the end of each epoch. During generation, we used greedy decoding.

Because PO techniques requires seeing both the preferred and the dis-preferred sequence during the same optimization step, we concatenate them together along the batch dimension so that both se-quences are processed in the same forward pass. This doubling of sequences in each batch requires that the number of training examples per batch be halved and the number of gradient accumulation steps be doubled in order to have the same effective batch size. This incurs no additional memory penalty but doubles the time to see the same number of training examples.

Our initial experiments showed that string-based metrics actually decrease when using PO techniques but we did not observe a discernible

quality difference. This is in line with the observations reported by (Xu et al., 2024b). Therefore, we evaluate with neural metrics so that evaluation could not be biased towards models that produce superficially similar translations. We evaluate with XCOMET-XL and -XXL (Guerreiro et al., 2023) and MetricX 23 (Juraska et al., 2023), both with and without references.

## 5 Results

Our XCOMET metric results are shown in Table 2. The evaluation shows that Tower Base is already competent at performing zero-shot translation for English to German, achieving reference-free XCOMET-XL and -XXL scores that are above the MT hypotheses contained in our dataset; 94.33 and 94.75, respectively. This is in spite of the fact that the model has not yet been instruction fine-tuned to perform zero-shot translation.

The lack of instruction fine-tuning is made obvious in the English to Russian results, where the model is unable to translate well prior to fine-tuning. Specifically, the Tower Base model frequently translated its instructions to Russian and ignored the source text, yielding lower scores. XCOMET-XL and -XXL seem to react differently to these non-translations with XCOMET-XXL punishing them more severely than XCOMET-XL.

Supervised fine-tuning is able to reach the level of the post-edits contained in the APE datasets when evaluating with reference-free evaluation. SFT surpasses the post-edits only with XCOMET-XL on the En→De data but this improvement is not significant. Here, we are evaluating the post-edits included in the dataset as if they were hypotheses for the source sentences.

IPO and dCPO are able to improve XCOMET-XL and -XXL scores for En→De above what the post-edits achieve, but only for -XXL is this improvement significant, as evaluated by pairwise bootstrap resampling implemented in the COMET package. For En→Ru, only dCPO is able to surpass post-edits and even then only for XCOMET-XL.

However, once we initialize the PO methods with the SFT model, we find our best results. SFT→dCPO is significantly better than both the MT and PE data from the dataset, the Tower Base model, and the SFT model for both En→De and Ru; while for just En→Ru, it is better than all other

| Model | PE − MT | | |
|---|---|---|---|
| | Train | Dev | Test |
| Base | 0.038 | 0.048 | 0.049 |
| SFT | 0.060 | 0.070 | 0.073 |
| IPO | 0.120 | 0.124 | 0.134 |
| dCPO | 0.110 | 0.115 | 0.124 |
| SFT→IPO | 0.144 | 0.144 | 0.157 |
| SFT→dCPO | 0.138 | 0.138 | 0.150 |

Table 4: This table shows the average values of the post-edit log-probabilities minus the machine translation log-probabilities for the English→German data. We see that the gap between PE and MT increases more with PO than it does with SFT.

systems.

Results with references do not differ drastically and can also be found in Table 2. We also evaluate with MetricX 23 XL (Juraska et al., 2023) and show our results in Table 3. The relations follow those of XCOMET and reinforce our conclusions.

## 6 Analysis

In addition to evaluating the fine-tuned models with neural metrics, we analyze the behavior of the models after training to see how the log-probabilities of the two sequences change compared to the baseline model. Additionally, we use the log-probabilities as a measure for model preferences. If one sequence is more probable, it is preferred by the model.

In our analysis, we remove machine translation and post-edit pairs where the post-edit remains unedited. This is because we are looking for differences in model behavior between machine translations and post-edits, which can not be done when they are the same sequence.

### 6.1 Log Probability Changes

Figures 2 and 3 are split violin plots showing the difference between log-probabilities before and after training, for German and Russian respectively. The left side of each violin shows the post-edit sequences' change from the baseline model's while the right side shows the machine translations' difference. This way we can examine how each training method affects the two sequence types individually. Additionally, we also measure the difference between the post-edits and machine translations after training in Tables 4 and 5 for German and Russian, respectively.
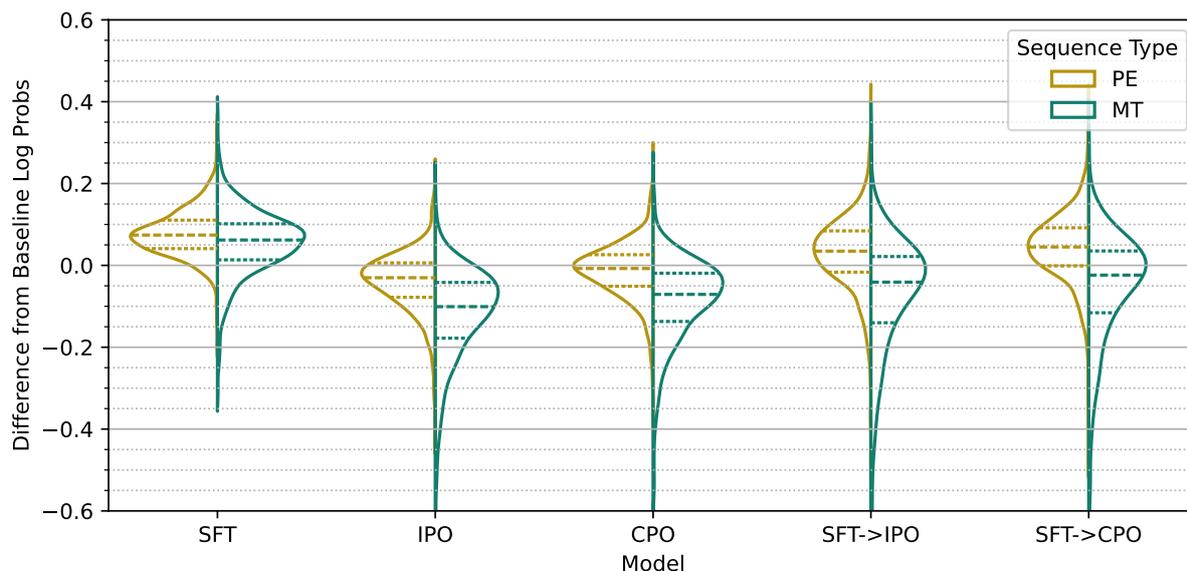
Figure 2: The difference of the models' averaged sequence log-probabilities from the baseline model's on the WMT 2020 En→De test data. Zero for PE is an average log-probability of −0.516 while for MT it is −0.565. This violin plot then shows displacement from these baseline values. Dashed horizontal lines indicate quartiles.
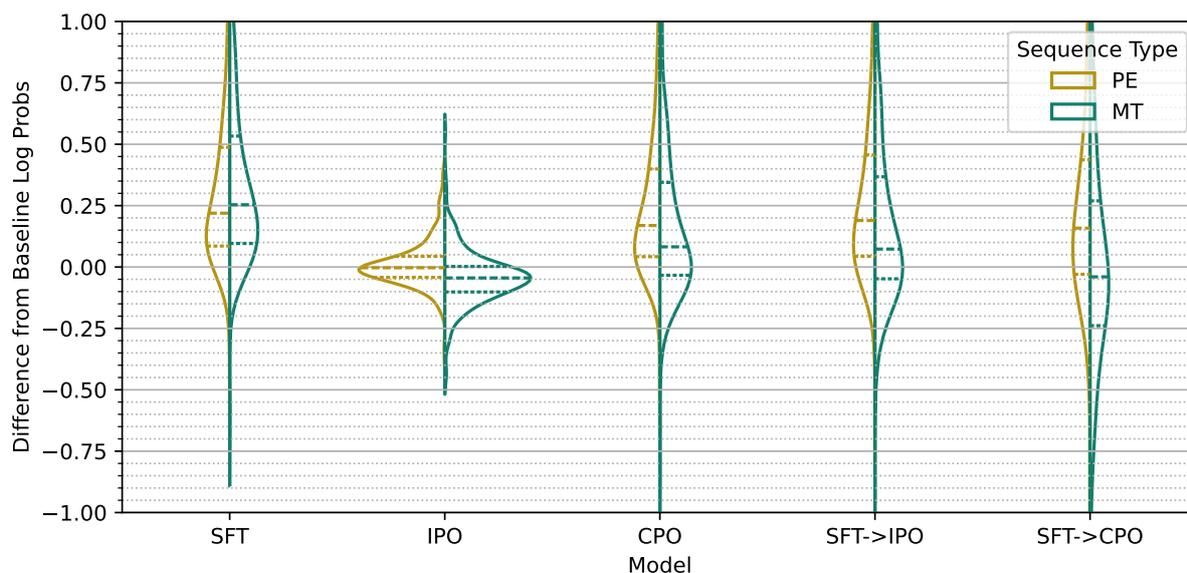


Figure 3: The difference of the models' averaged sequence log-probabilities from the baseline model's on the WMT 2019 En→Ru test data. Zero for PE is an average log-probability of −1.099 while for MT it is −1.260. This violin plot then shows displacement from these baseline values. Dashed horizontal lines indicate quartiles.

As we see in Figure 2, if we perform SFT on post-edits, as would typically be done when treating post-edits as new references, *both* the post-edits and the MT outputs become more likely under our fine-tuned model. Because the post-edits and MT outputs are highly correlated, they likely reside very close to each other in the model's hidden representation. This means, that with a smooth mapping from hidden representations to outputs, increasing the probability of the PE will

also increase the probability of the MT sequence.

For the En→De IPO and dCPO runs, we see the post-edits stay close to the baseline while the MT is pushed further down in log-probability space. Additionally, the distance between the two sequences increases under PO compared to SFT. As shown in Table 4, the average distance that PEs are above MT outputs doubles after PO compared to SFT.

After the IPO training, both sequences become

1295

less likely as seen in the split violin plot for IPO in Figure 2. This method does not have the upwards pressure on the preferred sequences that SFT or dCPO does, so we hypothesize that the downward pressure on the MT output also drags the PE sequence down as well; similar to how SFT increases the probability of MT without training on it. Alternatively, it could be that in order to establish a greater distance between the sequences, probability mass has to be re-allocated to other possible sequences.

With the En→Ru data, we see that the MT sequences benefit *more* from training than the PE sequences do, even though they remain unseen, as shown by the violin plot for SFT in Figure 3. This corresponds to the smaller difference between PE and MT that we see for SFT when compared to Tower Base in Table 5. The need for more fine-tuning of the Tower Base model is also visible in the SFT, SFT-initialized, and dCPO models' larger displacement from the Tower Base log-probabilities. IPO remains close to the 0 point for both sequences because the only pressure for each sequence is for them to move further apart; which is more difficult with the large overlap between the machine-translations and post-edits.

En→Ru appears similar for IPO, where both sequences are moved down in log-probability space, however the violin plot for dCPO and the SFT initialized models have a displacement from the baseline similar to SFT. This is because the baseline model was unable to perform zero-shot translation for En→Ru and, since the dCPO loss includes an SFT term, it learned how to translate which moved all sequences upwards. Unlike SFT, post-edits benefit more than machine translations after dCPO training.

Finally for the En→De SFT initialized models, we see in Figure 2 that post-edits increase in probability over the baseline while machine translation outputs are held close to or below the baseline. The difference between PE and MT is increased here compared to the PO only conditions.

We find that this behavior generalizes also to the development and test sets as shown in Table 4. For En→Ru, the SFT→IPO model and the dCPO model both have post-edits and machine translation increase in likelihood compared to the baseline. This is again due to the baseline model being unable to perform zero-shot translation and both sequences become more likely after it is able to do so. SFT→dCPO appears similar but far more

| Model | PE − MT | | |
| | Train | Dev | Test |
| --- | --- | --- | --- |
| Base | 0.039 | 0.078 | 0.161 |
| SFT | 0.025 | 0.062 | 0.133 |
| IPO | 0.085 | 0.125 | 0.217 |
| dCPO | 0.101 | 0.140 | 0.240 |
| SFT→IPO | 0.110 | 0.151 | 0.263 |
| SFT→dCPO | 0.192 | 0.242 | 0.419 |

Table 5: This table shows the average values of the post-edit log-probabilities minus the machine translation log-probabilities for the En→Ru data. We see that the gap between PE and MT increases more with PO than it does with SFT.

stretched out and with MT moved below the baseline. This model trained for much longer before reaching its early stopping criterion (SFT→dCPO stopped after 10 epochs, compared to SFT→IPO stopping after 2).

The largest improvements in our XCOMET-XL and -XXL scores coincide with training methods that both move the post-edit up in log-probability space while also ensuring that the machine translations are less likely by enough of a margin. SFT on its own also increases the probability of machine translations and does not work to establish a margin between the two sequences. Additionally, this shows us that PO successfully moves the model towards generating post-edit-like translations rather than those like the machine translations.

## 6.2 Preference Changes

Changes in log-probabilities from the baseline model do not necessarily indicate whether the models' preferences have changed. It could be that, in $(mt, pe)$ pairs where it is already the case that $pe > mt$, the distances between $pe$ and $mt$ increased, but examples where $mt > pe$ did not have their ordering changed. To that end, we also examine the baseline model's preference in terms of sequence probability—if a sequence's average log probability is strictly greater than that of another sequence, it is preferred. We plot preferences across all data splits for En→De in Figure 4 and for En→Ru in Figure 5. The exact values with corresponding confidence intervals are in Tables 6 and 7, respectively.

For both language pairs, we find that the Tower Base model does not have strong preferences. On the En→De data set, it prefers post-edits to ma-
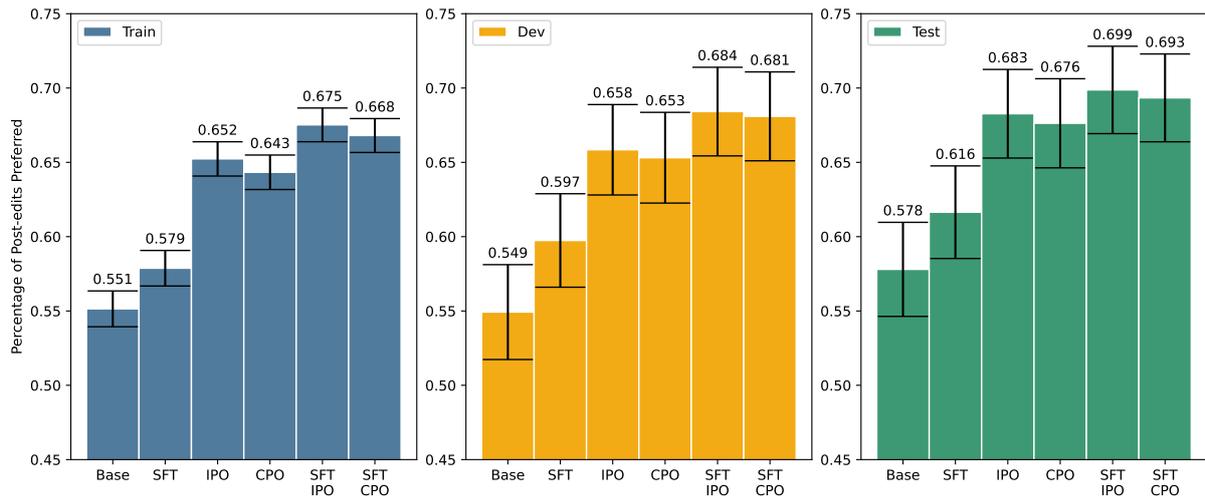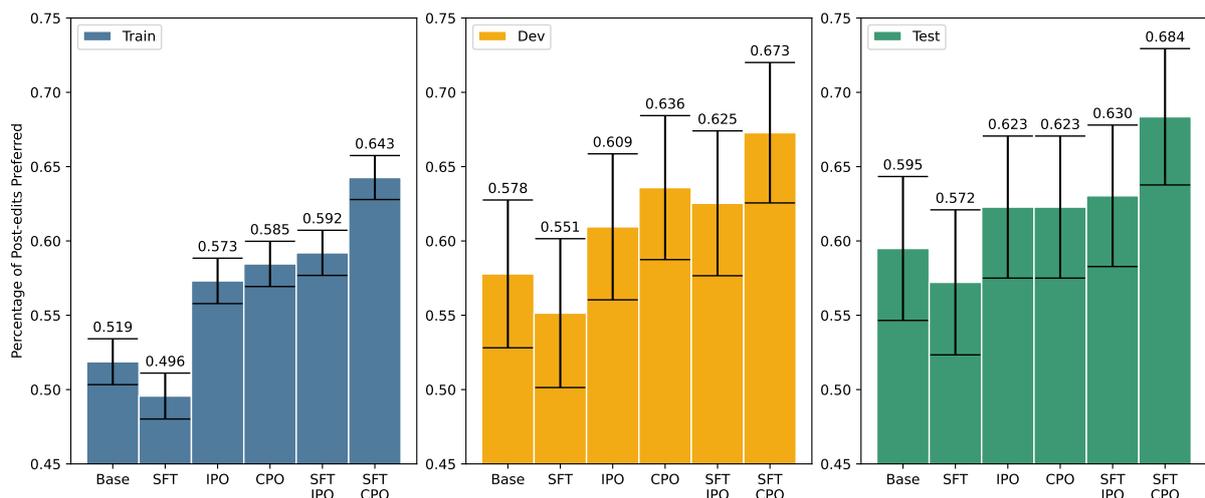
Figure 4: Here we show the percentage of training examples where the post-edit sequence is preferred in terms of average log-probability over the machine translation for the WMT En→De dataset. The black lines indicate the 95% confidence intervals for binomial distributed data—non-overlapping confidence intervals indicate a significant difference.

| Method | Train | Dev | Test |
|---|---|---|---|
| Base | 55.14% (53.94%, 56.35%) | 54.93% (51.73%, 58.12%) | 57.80% (54.64%, 60.96%) |
| SFT | 57.88% (56.68%, 59.07%) | 59.74% (56.60%, 62.89%) | 61.65% (58.53%, 64.76%) |
| IPO | 65.23% (64.08%, 66.39%) | 65.85% (62.80%, 68.89%) | 68.27% (65.29%, 71.25%) |
| dCPO | 64.33% (63.17%, 65.49%) | 65.31% (62.26%, 68.36%) | 67.63% (64.63%, 70.63%) |
| SFT→IPO | 67.52% (66.39%, 68.66%) | 68.42% (65.43%, 71.40%) | 69.87% (66.93%, 72.81%) |
| SFT→dCPO | 66.80% (65.66%, 67.94%) | 68.09% (65.10%, 71.08%) | 69.34% (66.38%, 72.29%) |

Table 6: Percentage of instances where post-edits are preferred over machine translations and their corresponding 95% confidence intervals for Train, Dev, and Test Splits for the WMT En→De 2020 APE Dataset. Non-overlapping confidence intervals correspond to statistically significant differences with $\alpha < 0.05$.



Figure 5: Here we show the percentage of training examples where the post-edit sequence is preferred in terms of average log-probability over the machine translation for the WMT En→Ru dataset.. The black lines indicate the 95% confidence intervals for binomial distributed data—non-overlapping confidence intervals indicate a significant difference.

chine translation 57.80% of the time on the test set    while for En→Ru this preference occurs 59.49%

| Method | Train | Dev | Test |
|---|---|---|---|
| Base | 51.87% (50.33%, 53.42%) | 57.78% (52.81%, 62.76%) | 59.49% (54.65%, 64.33%) |
| SFT | 49.57% (48.02%, 51.11%) | 55.15% (50.14%, 60.15%) | 57.22% (52.34%, 62.09%) |
| IPO | 57.31% (55.79%, 58.84%) | 60.95% (56.04%, 65.86%) | 62.28% (57.50%, 67.06%) |
| dCPO | 58.46% (56.93%, 59.98%) | 63.59% (58.74%, 68.43%) | 62.28% (57.50%, 67.06%) |
| SFT→IPO | 59.20% (57.68%, 60.72%) | 62.53% (57.66%, 67.41%) | 63.04% (58.28%, 67.80%) |
| SFT→dCPO | 64.27% (62.79%, 65.75%) | 67.28% (62.56%, 72.01%) | 68.35% (63.77%, 72.94%) |

Table 7: Percentage of instances where post-edits are preferred over machine translations and their corresponding 95% confidence intervals for Train, Dev, and Test Splits for the WMT En→Ru 2019 APE Dataset. Non-overlapping confidence intervals correspond to statistically significant differences with $\alpha < 0.05$.

of the time. For En→De, SFT significantly improves this preference on the training data but not on the development or test data. SFT actually seems to change the preferences in favor of machine translations on the En→Ru data; which also coincides with a decrease in the average distance between sequences and machine translations increasing in probability more.

When we train with IPO and dCPO on En→De, we find that both improve the preference for post-edits up to 68.27% on test data. The improvements above SFT are significant for both models on the train set while for dev, the confidence intervals overlap, and for test only IPO is significantly better. On En→Ru, we see a similar improvement in preferences but only on the training set are they significant.

Initializing with SFT and then training with PO on En→De yields the best improvements with 69.87% on test. Both SFT→IPO and SFT→dCPO are significantly better than SFT across all data splits. Again, En→Ru shows similar behavior with only the change on the training set being significant.

Across all data splits on En→De, IPO methods seem to establish a slightly stronger preference for post-edits which seems to be accounted for by increase in difference between the two sequence types as shown in Table 4. For En→Ru, dCPO is better at establishing this preference which also coincides with the increase in differences from Table 5.

## 7 Conclusion

Post-editing is part of common translation workflows before publishing to clean up raw-MT outputs. If the post-edits are used for training purposes, they are treated simply as new references and the MT output is treated as a by-product. Post-edits are created with an implicit preference in mind, that the PE should be better than the MT. We find that keeping both the PE and MT allows us to perform preference optimization techniques and improve translation quality with data that would otherwise be discarded.

We find that performing supervised fine-tuning using post-edits as references also increases the likelihood of the machine translations which remained unseen by the system. However, because the original machine translations were erroneous (in order to need correction), it is disadvantageous to increase their likelihood as well. Using PO techniques allows the model to establish a larger margin between the post-edit sequence and the machine translation sequence in log-probability space.

Increasing this margin coincides with significant improvements in neural translation metrics. We additionally find that we can measure the models' preferences in terms of sequence probability—if one sequence is more likely it is preferred. Models trained with SFT do not have a significant change in preferences compared to the baseline models but using PO teaches the model to prefer the post-edits above the machine translations.

In future work, we would like to examine the effect of the distance between post-edit and machine translation sequence probabilities. Currently, IPO sets a single distance for all sequence pairs but this may be sub-optimal when the sequences are correlated to different degrees. For example, if a post-edit and machine translation share a large prefix, the rest of the tokens in the sequences must account for the distance, while for non-overlapping sequences all tokens contribute to the distance between the log-probabilities.

## References

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *Preprint*, arXiv:2402.17733.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *Preprint*, arXiv:2204.05862.

Nathaniel Berger, Miriam Exel, Matthias Huck, and Stefan Riezler. 2023. Enhancing supervised learning with contrastive markings in neural machine translation training. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 69–78, Tampere, Finland. European Association for Machine Translation.

Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.

Rajen Chatterjee, Christian Federmann, Matteo Negri, and Marco Turchi. 2019. Findings of the WMT 2019 shared task on automatic post-editing. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 11–28, Florence, Italy. Association for Computational Linguistics.

Rajen Chatterjee, Markus Freitag, Matteo Negri, and Marco Turchi. 2020. Findings of the WMT 2020 shared task on automatic post-editing. In *Proceedings of the Fifth Conference on Machine Translation*, pages 646–659, Online. Association for Computational Linguistics.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. 2024. A general theoretical paradigm to understand learning from human preferences. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 4447–4455. PMLR.

Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *Preprint*, arXiv:2310.10482.

Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, Marc Sun, and Benjamin Bossan. 2022. Accelerate: Training and inference at scale made simple, efficient and adaptable. https://github.com/huggingface/accelerate.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google Submission to the WMT 2023 Metrics Shared Task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Julia Kreutzer, Joshua Uyheng, and Stefan Riezler. 2018. Reliability and learnability of human bandit feedback for sequence-to-sequence reinforcement learning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1777–1788, Melbourne, Australia. Association for Computational Linguistics.

Arle Richard Lommel, Aljoscha Burchardt, and Hans Uszkoreit. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.

OpenAI, Josh Achiam, et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics (ACL)*, Philadelphia, PA.

Maja Popović. 2015. chrF: character n-gram f-score for automatic mt evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation (WMT)*, Lisbon, Portugal.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn.

2023. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems*, volume 36, pages 53728–53741. Curran Associates, Inc.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas (AMTA'06)*, Cambridge, MA.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024a. A paradigm shift in machine translation: Boosting translation performance of large language models. *Preprint*, arXiv:2309.11674.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024b. Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. In *Forty-first International Conference on Machine Learning*.

Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul F. Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *arXiv*, abs/1909.08593.

## A  Hyperparameters

The hyperparameters for all runs are shown in Table 8. All hyperparameters are shared between runs.

| Hyperparameter | Value |
| --- | --- |
| Max Epochs | 20 |
| Learning Rate | 2e-6 |
| Optimizer | AdamW |
| Learning Rate Scheduler | Cosine |
| Warm-up Ratio | 0.1 |
| Effective Batch Size | 256 |
| Max Gradient Norm | 10.0 |
| Mixed Precision | bfloat16 |
| Early Stopping Criterion | XCOMET-XL w/o Refs |
| Early Stopping Patience | 3 |
| Early Stopping Epsilon | 0.00001 |
| Evaluation Frequency | Epoch |
| Max New Tokens | 64 |
| $\beta^*$ | 0.1 |
| Average Log-Probabilities | True |
| Normalize Loss | True |

Table 8: Hyperparameters for all training runs. * indicates that this parameter only affects the preference optimization techniques

# Translating Step-by-Step: Decomposing the Translation Process for Improved Translation Quality of Long-Form Texts

**Eleftheria Briakou, Jiaming Luo, Colin Cherry, Markus Freitag**
Google
{ebriakou,jmluo,colincherry,freitag}@google.com

## Abstract

In this paper we present a step-by-step approach to long-form text translation, drawing on established processes in translation studies. Instead of viewing machine translation as a single, monolithic task, we propose a framework that engages language models in a multi-turn interaction, encompassing pre-translation research, drafting, refining, and proofreading, resulting in progressively improved translations. Extensive automatic evaluations using Gemini 1.5 Pro across ten language pairs show that translating step-by-step yields large translation quality improvements over conventional zero-shot prompting approaches and earlier human-like baseline strategies, resulting in state-of-the-art results on WMT 2024.

## 1 Introduction

Machine Translation (MT) has been traditionally seen as a sequence transduction task that maps a source text from one language to an equivalent translation in another language. While this simplified definition of the task served the modeling capabilities of statistical and neural machine translation systems for many years, recent advancements in large language modeling offer promise for redefining MT to align more closely with human translation processes. This shift prompts us back to a fundamental question: *what does a good translation process look like*?

Thankfully, this question has been a long-debated topic in the field of translation studies. Despite the lack of consensus around the nature of cognitive steps involved when humans translate, a common thread is apparent, i.e., translation is a multi-faceted *process* encompassing several sub-tasks that navigate a bilingual landscape. This view of translation finds a parallel in the rise of the "chain-of-thought" paradigm popularized by large language models (LLM) (Wei et al., 2022). That is, instead of attempting to generate the response to a



Figure 1: MetricX-23 quality improvements (where lower scores indicate better translation quality) on document-level translation on the WMT24 test set. Translate step-by-step with Gemini 1.5 Pro consistently outperforms zero-shot translation.

complex task directly, LLMs are prompted to derive their final answer by decomposing the original task into several simpler sub-tasks.

But, what form would chain-of-thought take in the context of MT? While initial attempts to model the entire translation process using complex multi-stage processes has shown mixed results (Wu et al., 2024), explicitly modeling certain pre-translation or post-translation processes has led to more consistent gains in translation quality. On the pre-translation side, He et al. (2023) proposes to generate multiple translation candidates conditioned on self-generated translation-related knowledge. On the post-translation side, recent research threads prompt LLMs for refinement with (Feng et al., 2024; Xu et al., 2023b; Ki and Carpuat, 2024) or without (Chen et al., 2023) external quality estimation feedback.

Despite the promising results reported by prior work on decomposing and re-ranking MT with LLMs, it still remains unclear whether LLMs can

1301

benefit from modeling the *entire spectrum of translation processes*. In this work, drawing on literature from translation studies, we view MT as a complex and iterative task adhering to distinct steps, i.e., pre-translation research, drafting, refining, and proofreading. Based on this framework, we ask: *How well can LLMs translate in a step-by-step manner that draws from translation processes?*

Taking Gemini 1.5 Pro (Reid et al., 2024) as a case study, we start by designing instruction prompts for various translation subtasks. Concretely, our framework implements a **multi-turn** interaction with Gemini that breaks down the translation process into four distinct stages. It begins by prompting the model to conduct background research that identifies potential challenges in translating the source text (**research** phase). The next interaction focuses on drafting an initial translation prioritizing faithfulness to the source text (**drafting** phase). This draft is then revised in subsequent turns, ensuring a polished final translation (**refinement** and **proofreading** phases).

To align better with human translation processes, we test the *translate step-by-step* framework on long-form documents derived from the general MT shared tasks for WMT 2023 (Kocmi et al., 2023) and WMT 2024. We evaluate out-of-English translation for ten languages, namely Chinese (ZH), Ukrainian (UK), Russian (RU), Japanese (JA), Hebrew (HE), Czech (CS), German (DE), Hindi (HI), Icelandic (IS), and Spanish (ES). Extensive automatic evaluation according to both reference-based and QE-based versions of MetricX-23 (Juraska et al., 2023) show that translating step-by-step yields strong translation quality improvements across all languages and test sets studied (see Figure 1).

## 2  Background

With the recent rise of LLMs, machine translation is going through a gradual but significant paradigm shift. While much research is focusing on how LLMs' training data are improving their MT capabilities (Xu et al., 2023a; Alves et al., 2024), there are also many opportunities to improve how existing LLMs can be best used for translation. This becomes evident in recent research that explores ways to augment and refine MT to align better with human translation processes. To navigate the diverse landscape of LLM-driven research, we summarize key studies in Table 1 along their four most distinct dimensions:

| PAPER | PRE-TR. | POST-TR. | DEV. | PARAM. | STEPS |
|---|---|---|---|---|---|
| He et al. (2023) | ✔ | ✗ | ✗ | ✗ | 3-4 |
| Xu et al. (2023b) | ✗ | ✔ | ✔ | ✗ | Iterative |
| Feng et al. (2024) | ✗ | ✔ | ✔ | ✗ | 3 |
| Huang et al. (2024) | ✗ | ✔ | ✔ | ✗ | 3 |
| Li et al. (2024) | ✔ | ✗ | ✗ | ✗ | 1 |
| Chen et al. (2023) | ✗ | ✔ | ✗ | ✔ | Iterative |
| Ki and Carpuat (2024) | ✗ | ✔ | ✗ | ✗ | 1 |
| Wu et al. (2024) | ✔ | ✔ | ✗ | ✔ | Iterative |
| Step-by-Step (ours) | ✔ | ✔ | ✔ | ✔ | 4 |

Table 1: List of prior work leveraging LLMs to improve translation quality by modeling either pre- or post-translation processes (PRE-TR. or POST-TR.). For each study we also note key aspects of their methodology: whether prompting strategies are developed on a separate development set (DEV.), whether the approach relies solely on the LLM's parametric knowledge (PARAM.), and the number of steps in the pipeline.

- **Temporal Focus**: This differentiating factor is based on whether an LLM is engaged in the translation process before (pre-translation) or after (post-translation) an initial translation is produced (whether by the same LLM or a different system).

- **Parametric vs. External Knowledge**: This dimension focuses on whether LLMs rely solely on their internal, learned knowledge (encoded in their parameters) or whether they use external resources, i.e., dictionaries, knowledge bases, retrieval engines or QE-based metrics (Mallen et al., 2023).

- **Reported Prompt Development**: This dimension considers whether the prompting strategies are clearly developed on separate development sets, as reported in papers.[1]

- **Number of Steps**: This dimension counts the number of distinct steps that are used in multi-turn interactions with the LLM.

Table 1 shows a clear trend: most studies focus on post-translation refinement. These approaches predominantly rely on external feedback to identify and correct errors, using either automatic met-

---

[1]We include this column not to cast aspersions on previous work, but to encourage a culture moving forward where prompt-based research uses and reports a development set. From personal communication, some of the works receiving an "✗" here underwent little to no prompt optimization.

## Pre-drafting Research

You will be asked to translate a piece of text from English into Chinese following the five stages of the translation process. Here is the context in which the text appears:

Context: However, he said the most "responsible way" to cut spending would be to pass all 12 bills. With many Republicans reticent to see a shutdown and a 1 percent cut to defense spending, the urgency to pass the bills could see enough Republicans ally with Democrats to extend the deadline on the four spending bills and reach an agreement to fund the government, even if doing so puts Johnson in hot water with some of those on the hard right.

To start, let's do some pre-drafting research on the above context:

**Research:**

During this phase, thorough research is essential to address components of the context text that pose translation challenges. The goal is to establish a comprehensive translation plan that covers the following category:
- Idiomatic Expressions:
  - Identify idiomatic expressions that cannot be directly translated word-for-word into Chinese.

**Idiomatic Expressions in the Context:**

Here are some idiomatic expressions in the context that require careful consideration for translation into Chinese:
- **"in hot water":** This idiom means to be in trouble or facing criticism. A direct translation wouldn't convey the meaning.
- **"hard right":** This refers to the most conservative faction within a political party. A literal translation wouldn't be accurate.
- **"reticent to see":** This implies reluctance or unwillingness. A more nuanced translation is needed.

然而，他表示，削减开支最"负责任的方式"是通过所有 12 项法案。由于许多共和党人不愿看到政府停摆和国防开支削减 1%，因此通过这些法案的紧迫性可能会促使足够多的共和党人与民主党人结盟，延长四项支出法案的截止日期，并达成一项为政府提供资金的协议，即便这样做会让约翰逊与一些极右翼人士产生摩擦。

## Drafting

Now, let's move on to the drafting stage.

**Draft Translation:**

In this phase, your primary objective is to create a draft translation that accurately conveys the meaning of the source text presented below. At this stage, it is crucial to focus on adequacy, ensuring that your translation closely adheres to the source text. Your response should conclude with the draft translation. If context is missing, generate a general translation that is adaptable to various contexts. Avoid adding any additional information not present in the source text. All elements of the source text should be present in the translation.

Give your best one translation for the following piece of text based on the pre-drafting analysis without providing alternatives:

English: However, he said the most "responsible way" to cut spending would be to pass all 12 bills. With many Republicans reticent to see a shutdown and a 1 percent cut to defense spending, the urgency to pass the bills could see enough Republicans ally with Democrats to extend the deadline on the four...

然而，他表示，削减开支最"负责任的方式"是通过所有 12 项法案。鉴于许多共和党人不愿看到政府停摆和国防开支削减 1%，通过法案的紧迫性可能会促使足够多的共和党人与民主党人结盟，延长四项支出法案的截止日期，并达成一项为政府提供资金的协议，即使这样做会让约翰逊与一些极硬的右翼人士产生矛盾。

## Refinement

Now let's move to the next stage.

**Post-editing with local refinement**

In this stage, the primary aim is to refine the draft translation by making micro-level improvements that improve the draft's fluency.

Provide only one refined translation and do not output anything else after that.

然而，他表示，削减开支最"负责任的方式"是通过全部 12 项法案。鉴于许多共和党人不愿看到政府停摆和国防开支削减 1%，通过法案的紧迫性可能会促使足够多的共和党人与民主党人合作，延长四项支出法案的截止日期，并达成一项为政府提供资金的协议，即使这样做会让约翰逊与一些强硬的右翼人士产生矛盾。

## Proofreading

You are tasked with proofreading a translation that has been revised for improved fluency. The refined translation has been generated by editing the draft translation.

**Proofreading and Final Editing**

The goal is to provide a polished final translation of the source text. For you reference, below are the source text, the draft, and refined translations.

Source Text
However, he said t....

Draft Translation
然而，他表示....

Refined Translation
然而，他表示....

Please proofread the refined text for grammar, spelling, punctuation, terminology, and overall fluency. Ensure the translation accurately reflects the original meaning and style. Provide only the final, polished translation.

然而，他表示，削减开支最"负责任的方式"是通过全部12项法案。鉴于许多共和党人不愿看到政府停摆和国防开支削减1%，通过法案的紧迫性可能会促使足够多的共和党人与民主党人合作，延长四项支出法案的截止日期，并达成一项为政府提供资金的协议，即便这样做会让约翰逊与一些强硬的右翼人士产生摩擦。

Figure 2: Translate Step-by-Step prompting framework. User prompts (top) and Gemini's responses (bottom) for the translation of an English document into Chinese. The full prompts for each step also appear in §A.3.

rics (Feng et al., 2024; Xu et al., 2023b; Huang et al., 2024) or human annotations of translation errors (Ki and Carpuat, 2024). A notable exception is the study of Chen et al. (2023), which shows that LLMs can iteratively refine their own outputs using only their parametric knowledge.

Comparatively fewer studies explore the pre-translation stage, investigating how LLMs can utilize background information to enhance their translation quality. He et al. (2023) explores this by prompting LLMs for different types of background information (similar examples, topics and keywords) related to the source text. However, they find that this knowledge alone is insufficient to improve the model's translation quality, and ultimately rely on external QE feedback for selection. In contrast, Li et al. (2024) operationalizes background research by incorporating idiom definitions retrieved from an external knowledge base.

A notable exception to the above is the recent work of Wu et al. (2024) which, similar to our approach, explores modeling the entire spectrum of translation processes. While conceptually aligned with our step-by-step approach, their framework is significantly more complex, with 30 distinct LLM roles interacting iteratively. Their use of non-standard metrics makes it difficult to gauge the method's success: the human evaluation does not give annotators source or reference texts, while the

bilingual automatic evaluation collects only preference decisions using the same model family as the method being tested.

Overall, in contrast to prior work, which often relies on complex multi-stage processes and external resources, our goal is to streamline the translation process, *unifying pre- and post-translation stages within one framework, by accessing only the* LLM*'s parametric knowledge throughout*. We emphasize the methodological soundness of our pipeline by developing it on a separate development set, a practice not yet standardized in this area.

## 3 Translate Step-by-Step

Drawing on existing literature on translation studies (Borg, 2018), we design a series of staged prompts that attempt to map the translation process to instructions. This approach views translation as a multi-turn interaction with an LLM where each prompt guides the model's next action. Below, we describe those stages, along with what their function in the translation process is and how they are operationalized as instruction-following tasks. These stages are further illustrated in Figure 2.

**Pre-translation Research** Mirroring the human translation processes, our framework incorporates a pre-translation research stage. This stage primarily

focuses on using the source text (Mossop, 2000) to identify potential translation challenges drawing on real-world knowledge and knowledge of the target language (Dimitrova, 2005). We model this stage by prompting the LLM to identify and explain phrases of the source text that cannot be translated word-for-word into the target language.

**Drafting**   Following the pre-translation research, the next stage aims at producing a draft translation, i.e., "the first stab at the rewriting" (Bassnett and Bush., 2016). This stage represents an initial attempt at rendering the source text into the target language. To that end, we initiate a subsequent interaction and prompt the model to focus on adequacy at this stage, ensuring the draft faithfully captures the meaning of the source.

**Refinement**   The *post-drafting* stages are defined as editing tasks, with the goal of improving the overall quality of the draft translation. We define the first post-drafting stage as a subsequent interaction where the LLM is prompted to improve the draft's fluency such that the text works on its own (Borg, 2018).

**Proofreading**   At the final, post-drafting stage, we task the LLM with the role of proofreading the refined translation to ensure it delivers a polished translation. We model this stage as a new conversation with the LLM, rather than a subsequent interaction, drawing inspiration from human studies suggesting that proofreading requires a new perspective after a break from revising (Shih, 2013).

### 3.1   Lessons During Development

While developing the above method, we found two factors to be important for the success of this approach: working at the document level and representing multi-step interactions as conversations.

**Working at the Document Level**   Our multi-step process became more effective as we moved from the segments provided by WMT to working on multi-segment documents (see §4 for details on the setup). This had a large effect on the pre-translation research step, changing it in two ways. First, some phrases that appeared idiomatic or difficult at the segment level disappeared, as their translations became clear with context. Second, the LLM began identifying larger phrases. The refinement step also improved according to automatic metrics. We verified that our shift to the document level was either neutral or an improvement for our baselines (§5.3).

| Domain | Literary | News | Social | Speech |
|---|---|---|---|---|
| # Docs. | 40 | 43 | 48 | 111 |
| Avg. Length | 192 | 184 | 164 | 73 |

Table 2: Per-domain statistics for WMT 2024.

**Multi-step Interactions as Conversations** Modern LLMs use special markers to indicate human versus assistant turns in multi-turn interactions. When building an automated process like translate step-by-step, for each step, one has the option to either use previous outputs to build a completely new query that summarizes all previous interactions, or to continue the conversation, allowing the LLM to see all previous steps with its own outputs clearly marked. With the exception of the proofreading step, we found that continuing the conversation improved performance. Also, breaking the conversation into smaller turns helps with modularity for ablations.

## 4   Experimental Setting

We start by evaluating the translate step-by-step approach on the task of document-level translation. The experimental setting is described below.

**Model Settings**   Throughout our experiments we use Gemini 1.5 Pro. All model outputs are generated with greedy decoding. All model prompts are provided in Appendix A.3. In zero-shot mode, the model is instructed to translate the source text directly, without providing any explanations.

To effectively isolate the artifacts from pre-translation research, we employ a secondary model call. This call restructures the natural language output into a JSON object, simplifying the parsing process for extracting artifacts.

**Evaluation Sets**   We use WMT 2023 as our ***development*** set. Any prompt development and stage ablation experiments are conducted on this dataset. For our final ***test*** set, we use the WMT 2024 datasets. Each of these datasets was built by translating a set of English documents into multiple languages.

Both datasets are segmented for sentence- or paragraph-level evaluation, but our approach focuses on translating with as much context as possible. Therefore, we use meta-data to merge the original segments into larger ones. Ideally, this would result in complete documents, but current neural metrics have token-count limits beyond which they truncate their inputs. To accommodate neural evaluation, we set a maximum length of 250 (English

| Research | Draft | Refinement | Proofreading | ZH | UK | RU | JA | HE | CS | DE | AVERAGE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| *Ref-based* | | | | | | | | | | | |
| 1. ○ | ○ | ○ | ○ | 3.64 | 4.18 | 3.32 | 2.59 | 4.36 | 2.82 | 1.82 | 3.25 |
| 2. ○ | ● | ○ | ○ | 3.48 ↓0.16 | 4.16 ↓0.02 | 3.32 ↓0.00 | 2.47 ↓0.12 | 4.54 ↑0.18 | 2.67 ↓0.15 | 1.92 ↑0.10 | 3.22 |
| 3. ○ | ○ | ● | ○ | 2.92 ↓0.72 | 3.32 ↓0.86 | 2.43 ↓0.89 | 2.19 ↓0.40 | 3.24 ↓1.12 | 2.35 ↓0.47 | 1.31 ↓0.51 | 2.54 |
| 4. ○ | ● | ● | ○ | 2.85 ↓0.79 | 3.06 ↓1.12 | 2.54 ↓0.78 | 2.09 ↓0.50 | 3.18 ↓1.18 | 2.22 ↓0.60 | 1.37 ↓0.45 | 2.47 |
| 5. ● | ● | ○ | ○ | 3.00 ↓0.64 | 3.46 ↓0.72 | 2.56 ↓0.76 | 2.05 ↓0.53 | 3.89 ↓0.47 | 1.97 ↓0.85 | 1.56 ↓0.26 | 2.64 |
| 6. ● | ● | ● | ○ | 2.63 ↓1.01 | 2.70 ↓1.47 | 2.13 ↓1.19 | 1.73 ↓0.86 | 2.88 ↓1.48 | 1.85 ↓0.96 | 1.17 ↓0.65 | 2.16 |
| 7. ● | ● | ● | ● | 2.67 ↓0.97 | 2.38 ↓1.80 | 2.16 ↓1.16 | 1.70 ↓0.89 | 2.75 ↓1.61 | 1.71 ↓1.10 | 1.07 ↓0.75 | 2.06 |
| *QE-based* | | | | | | | | | | | |
| 8. ○ | ○ | ○ | ○ | 2.64 | 4.87 | 4.16 | 1.73 | 5.55 | 5.39 | 3.96 | 4.04 |
| 9. ○ | ● | ○ | ○ | 2.71 ↑0.07 | 4.78 ↓0.09 | 4.05 ↓0.11 | 1.65 ↓0.07 | 5.22 ↓0.33 | 5.14 ↓0.25 | 4.03 ↑0.08 | 3.94 |
| 10. ○ | ○ | ● | ○ | 2.11 ↓0.52 | 4.33 ↓0.54 | 2.82 ↓1.34 | 1.30 ↓0.43 | 4.49 ↓1.06 | 4.31 ↓1.08 | 2.89 ↓1.07 | 3.18 |
| 11. ○ | ● | ● | ○ | 2.04 ↓0.59 | 4.12 ↓0.75 | 3.31 ↓0.85 | 1.19 ↓0.54 | 4.30 ↓1.25 | 4.40 ↓0.99 | 3.36 ↓0.60 | 3.25 |
| 12. ● | ● | ○ | ○ | 2.26 ↓0.38 | 4.18 ↓0.69 | 3.50 ↓0.66 | 1.54 ↓0.19 | 4.60 ↓0.95 | 4.62 ↓0.77 | 3.73 ↓0.23 | 3.49 |
| 13. ● | ● | ● | ○ | 1.90 ↓0.73 | 3.39 ↓1.48 | 2.76 ↓1.40 | 1.23 ↓0.49 | 4.17 ↓1.38 | 4.12 ↓1.28 | 2.97 ↓0.99 | 2.93 |
| 14. ● | ● | ● | ● | 1.82 ↓0.81 | 3.43 ↓1.44 | 3.11 ↓1.05 | 1.25 ↓0.48 | 4.01 ↓1.54 | 3.56 ↓1.83 | 2.63 ↓1.33 | 2.83 |

Table 3: MetricX-23 evaluation results of translate step-by-step and its ablation variants on the WMT 2023 development datasets. We report both the reference-based and QE-based metric variants. Filled dots indicate active steps in the pipeline, while unfilled dots represent ablated steps. When all steps are ablated, the system defaults to zero-shot translation. Colored boxes highlight performance differences compared to zero-shot: blue shades indicate significant improvements at $p < 0.001$, green shades indicate significant improvements at $p < 0.05$, yellow shades indicate non-significant improvements ($p \geq 0.05$), while red shades indicate non-significant regressions ($p \geq 0.05$) against zero-shot. *Translate step-by-step surpasses zero-shot across the board, with each step incrementally improving translation quality.*

white-space separated) tokens each.[2] The resulting datasets consist of 192 documents of average token length 178 for WMT 2023 and, 243 documents of average token length 130 for WMT 2024, respectively. For WMT 2024 we also report per-domain results. Per-domain document counts and average lengths, as measured in English white-space separated tokens, are presented in Table 2.

**Evaluation Metrics** We evaluate our approach using MetricX-XXL-23 (Juraska et al., 2023), the metric adopted in the most recent WMT 2024 automatic evaluations. We report results on both the reference-based and the QE-based metric variants. Despite being trained at the sentence level, Deutsch et al. (2023) show that MetricX can effectively evaluate multi-sentence sequences, capped at its maximum window length. We note that MetricX is powered by mT5 (Xue et al., 2021), which minimizes the potential bias in favor of Gemini-generated translations.[3] We employ paired permutation tests to determine if the observed improvements across system pairs are statistically significant.[4]

---

[2] We also present results for a shorter set of documents, with a maximum length of 150 tokens in Appendix A.1.

[3] We also report ChrF (Popović, 2015) in Appendix A.2.

[4] https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.permutation_test.html

## 5 Quantitative Results

We start by analyzing the importance of each step in the translate step-by-step pipeline. Ablation results on the WMT 2023 development sets are presented in §5.1. Next, the generalizability of our final step-by-step recipe is evaluated on the WMT 2024 test sets in §5.2, with comparison to prior work in §5.3.

### 5.1 Analyzing Step Importance

Automatic evaluation results on our development sets are presented in Table 3. Overall, translation artifacts extracted through the step-by-step process yield consistently better document translations compared to the zero-shot mode according to both reference- (lines 3–7 vs. 1) and QE-based (lines 10–14 vs. 8) versions of MetricX. Ablating the various steps from the pipeline gives insights into how each step contributes to the overall quality improvements. We describe those below.

**Importance of Pre-translation Research** Modelling pre-translation processes is crucial for achieving higher quality translations compared to the zero-shot. Simply prompting for a draft translation without asking for pre-translation research yields only small and non-significant improvements or even regressions over the zero-shot (lines 2 vs. 1

1305

and 9 vs. 8). This result rules out the possibility that any observed improvements are solely due to a better prompt for the draft translation, which was modified to emphasize faithfulness to the source (§3). However, combining the research and draft steps achieves consistently higher quality translations compared to zero-shot (lines 5 vs. 1 and 12 vs. 8). Importantly, those improvements are consistently statistical significant ($p < 0.0001$) across languages (measured by reference-based metrics), except for Hebrew, which shows non-significant improvements compared to zero-shot ($p \geq 0.05$).

**Importance of Refinement** Moving to the evaluation of the *refined* document translations, we notice an interesting trend. The refinement step consistently improves the translation quality, regardless of the initial translation it processes, i.e., the zero-shot (lines 3 vs. 1 and 10 vs. 8), the single-turn draft (lines 4 vs. 2 and 11 vs. 9), and the research-informed draft (lines 6 vs. 5 and 13 vs. 12). This demonstrates that the effectiveness of the refinement stage is not conditioned on the initial translation. However, the strongest quality improvements—reaching consistently high levels of statistical significance ($p < 0.001$) over the zero-shot translations—are observed when the refinement stage is combined with the pre-translation research (lines 6 vs. 1 and 13 vs. 8), highlighting that those stages bring complimentary benefits.

**Importance of Proofreading** Finally, the evaluation of the *proofreading* document translations, indicate that this stage contributes modest average improvements (lines 7 vs. 6 and 14 vs. 13). Unlike previous stages, the impact of proofreading appears to be more language dependent. Ukrainian stands out as the only language that clearly benefits from a proofreading stage, while others show only minor differences in quality compared to their refined translations.

### 5.2 Generalizability of Step-by-Step

Table 4 presents results on the WMT 2024 test set. Across the board, translating step-by-step exhibits the same trends noticed on our development set (as discussed in §5.1). This confirms the generalizability of our proposed approach, crucially, on a wider range of languages. Concretely, the draft translations outperform the zero-shot translations. The refined stages bring additional quality improvements across the board, with the proofreading stage contributing small improvements for most languages.



Figure 3: Domain-level comparison between zero-shot and step-by-step translations on WMT 2024 using reference-based MetricX-23. Each data point represents the delta from zero-shot (dotted horizontal line). The steps are denoted as follows: 0 (zero-shot), D (draft after research), R (refinement), and P (proofreading).

To better understand the robustness of our approach we present a per-domain analysis in Figure 3. As shown, translation quality improvements of step-by-step translations over zero-shot are observed across all domains, with speech showing the least and social the most significant gains.

### 5.3 Contextualizing Step-by-Step Gains

Having demonstrated how translate step-by-step improves long-form translation with LLMs over zero-shot translation, we now contextualize these gains by comparing our approach to two repre-

| | DE | ES | ZH | RU | UK | JA | HI | IS | CS | AVERAGE |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | *Ref-based* |
| *Zero-shot* | 1.90 | 3.23 | 3.48 | 3.02 | 3.15 | 2.29 | 3.65 | 4.01 | 2.65 | 3.04 |
| SBYS: *Research & Drafting* | 1.68 ↓0.22 | 2.69 ↓0.54 | 2.99 ↓0.49 | 2.53 ↓0.49 | 2.81 ↓0.35 | 1.92 ↓0.37 | 2.52 ↓1.13 | 3.77 ↓0.24 | 2.30 ↓0.35 | 2.58 |
| SBYS: *Refinement* | 1.45 ↓0.45 | 2.29 ↓0.94 | 2.45 ↓1.03 | 2.21 ↓0.81 | 2.58 ↓0.57 | 1.64 ↓0.66 | 2.31 ↓1.35 | 3.14 ↓0.87 | 2.10 ↓0.55 | 2.24 |
| SBYS: *Proofreading* | 1.35 ↓0.54 | 2.27 ↓0.96 | 2.42 ↓1.06 | 2.21 ↓0.81 | 2.49 ↓0.66 | 1.67 ↓0.62 | 2.09 ↓1.56 | 3.15 ↓0.86 | 2.14 ↓0.51 | 2.20 |
| | | | | | | | | | | *QE-based* |
| *Zero-shot* | 1.97 | 2.59 | 2.23 | 1.87 | 2.23 | 1.32 | 4.81 | 3.47 | 2.08 | 2.51 |
| SBYS: *Research & Drafting* | 1.72 ↓0.25 | 2.23 ↓0.36 | 2.08 ↓0.15 | 1.54 ↓0.33 | 1.81 ↓0.41 | 1.19 ↓0.13 | 4.12 ↓0.69 | 3.43 ↓0.04 | 1.97 ↓0.11 | 2.23 |
| SBYS: *Refinement* | 1.38 ↓0.59 | 1.78 ↓0.81 | 1.71 ↓0.52 | 1.21 ↓0.66 | 1.34 ↓0.89 | 0.95 ↓0.37 | 3.47 ↓1.34 | 2.79 ↓0.68 | 1.51 ↓0.56 | 1.79 |
| SBYS: *Proofreading* | 1.25 ↓0.72 | 1.74 ↓0.84 | 1.63 ↓0.60 | 1.14 ↓0.73 | 1.32 ↓0.91 | 0.93 ↓0.40 | 3.35 ↓1.46 | 2.65 ↓0.82 | 1.45 ↓0.63 | 1.72 |

Table 4: MetricX-23 results comparing step-by-step (SBYS) with zero-shot on the WMT 2024 test datasets. When all steps are ablated, the system defaults to zero-shot translation. Colored boxes highlight performance differences compared to zero-shot: blue shades indicate significant improvements at $p < 0.001$, green shades indicate significant improvements at $p < 0.05$, while yellow shades indicate non-significant improvements ($p \geq 0.05$). *Translate step-by-step surpasses zero-shot, with each step incrementally improving translation quality.*

sentative baselines: a) methods that leverage non-parametric knowledge for best translation selection, and b) segment-level baselines that translate documents using the pre-defined segmentation provided in WMT 2024 test sets.

**Conditions** As a representative of the first class, we compare against MAPS (He et al., 2023). This baseline employs an LLM to analyze the source text for topic, keywords, and similar examples, generating three candidate translations conditioned on each knowledge type. Then, a QE metric selects the best candidate. To create a fair comparison, we re-implement their method using Gemini 1.5 Pro, using the prompts provided in their released code. To create an even stronger baseline, we perform candidate selection with the QE variant of MetricX-23, which we know correlates well with the final reference-based MetricX-23, creating an advantage for MAPS.

For the second class of baselines, we consider two approaches: a) zero-shot translation applied to each segment individually using Gemini 1.5 Pro, both with (ZERO-SHOT IN CONTEXT) and without (ZERO-SHOT) access to the full document in the input prompt, and b) a comparison with the segment-level translations from Unbabel-Tower70B, the top-performing system of WMT 2024 based on early automatic evaluations (Kocmi et al., 2024). To get comparable document-level metrics, before evaluation, we concatenate the segment-level translation back into the mini-documents, as described in §4.

We focus our comparisons on EN-DE, EN-JA, EN-ZH, as MAPS requires in-context demonstrations that were made available only for those languages by the original authors. For a fair comparison with Unbabel-Tower70B, we exclude the speech domain from our comparison, as WMT 2024

| METHOD | DOC. | EN-DE | | EN-ZH | | EN-JA | |
|---|---|---|---|---|---|---|---|
| UNBABEL-TOWER70B | ✗ | [1] | 1.42 | [1] | 2.77 | [2] | 2.16 |
| ZERO-SHOT | ✗ | [2] | 1.98 | [3] | 3.65 | [3] | 2.60 |
| ZERO-SHOT IN CONTEXT | ✗ | [2] | 1.86 | [2] | 3.33 | [2] | 2.19 |
| ZERO-SHOT | ✔ | [3] | 2.02 | [3] | 3.91 | [3] | 2.47 |
| MAPS | ✔ | [2] | 1.91 | [2] | 3.25 | [2] | 2.19 |
| SBYS: *Research & Drafting* | ✔ | [2] | 1.75 | [2] | 3.32 | [2] | 1.94 |
| SBYS: *Refinement* | ✔ | [1] | 1.41 | [1] | 2.73 | [1] | 1.58 |
| SBYS: *Proofreading* | ✔ | [1] | 1.27 | [1] | 2.75 | [1] | 1.73 |

Table 5: Comparison of step-by-step (SBYS) with representative baselines (lower scores are better) on WMT 2024 according to Metric-X (reference-based). The second column indicates whether translation is performed on the entire document or by merging segment-level translations. Numbered squares represent significance clusters (Freitag et al., 2023) at $p = 0.05$. *Translate step-by-step matches or exceeds all compared baselines, crucially, without accessing external resources.*

submissions were given ASR transcripts instead of human-sourced transcripts in this domain.

**Results** Table 5 compares step-by-step against various baselines. Notably, even the initial stage, where the draft translation is conditioned on cross-lingual research (SBYS: *Research & Drafting*) demonstrates competitive performance against MAPS, falling within the same statistical significance cluster. This highlights the effectiveness of our pre-translation strategy compared to the background information used by MAPS. Comparing the final, proofreading stage of step-by-step (SBYS: *Proofreading*) with MAPS reveals significant translation quality gains: 0.64 improvement for DE, 0.50 for ZH, and 0.46 for JA. Notably, these improvements are achieved even though MAPS uses the same QE model family as MetricX for final candidate selection, giving it an inherent advantage. In contrast, SBYS relies solely on the model's internal, parametric knowledge throughout the entire translation process.

Comparing the final, proofreading stage of step-

by-step with the segment-level baselines helps put the improvements in perspective. Concretely, the segment-level zero-shot baselines (second and third lines in Table 5) fall significantly behind the step-by-step final translations (SBYS: *Proofreading*) across all languages by more than 0.7 and 0.4 MetricX points when compared to the out-of- and in-context variants, respectively. This demonstrates that simply translating documents at a finer granularity is not sufficient for boosting the LLM's translation quality.

Finally, comparing the final, proofreading stage of our approach with the merged translations from Unbabel-Tower70B, reveals that our approach achieves statistically comparable performance for Chinese and German (0.02 and 0.15 improvements respectively) and significantly better performance for Japanese (0.43 improvement). These improvements over the top-performing WMT 2024 system demonstrate the competitiveness of the step-by-step approach, especially given that the competing system relies on external QE metrics and computationally expensive decoding strategies to improve translation quality.

## 6 Qualitative Analysis

We conduct a qualitative analysis on a small subset of model outputs from all stages to understand the strengths and weaknesses of our step-by-step approach. To this end, we first compute the score difference between the final translation and the zero-shot output on WMT 2024 English to Chinese, and then randomly sample up to 5 examples from either end (i.e., examples for which the final translation quality either substantially improves or degrades over the zero-shot baseline).[5] One of the authors (native speaker of Chinese) manually inspected the sampled outputs and took notes on the salient properties of the pre-translation artifacts and the incremental changes from the different stages of the step-by-step process.

**Pre-drafting** For pre-drafting research, we observe that the LLM is highly capable of understanding the source in a wide variety of contexts. As showcased in Table 6, the LLM is able to correctly interpret slang (example 1: *cheeked up* in the context of making miniatures), recognize figurative

usage (example 2: *the weather didn't cooperate* in the context of flying a plane), and detect humorous expressions (example 3). This strength is especially pronounced when even the references show clear signs of human translators misinterpreting the source (see the next subsection for full examples).

On the other hand, the LLM is also prone to over-generate and seems too eager to confirm with the given instruction to find instances of indirect translation. This resulted in false positives where a direct and literal translation is already adequate (example 4: *a bit dazed* can be directly translated into Chinese), and in some cases bizarre cultural commentaries (example 5 for asking to contextualize the texture of bubble gum).

**Translations** The observed understanding of the source texts seems to directly contribute to more fluent and context-appropriate translations. Table 13 in §A.4 shows several interesting examples. There are quite a few instances where the step-by-step approach produces the correct translation even when the reference fails to recognize the context the phrase appears in. For example, our method correctly interprets the meaning of *cheeked up* in the first example to be "having a full cheek" when the zero-shot translates it to *blushing* and the reference translates it to *talking nonsense*. Similarly for the second example, the term *threading* is correctly understood as a thread of posts on a social media platform by all step-by-step translations, whereas the zero-shot interpretes it as a computing terminology (as in *multithreading*) and the reference interprets it as *study*.

It is also evident that the refinement improves the fluency significantly. The third example shows that both the zero-shot and the draft translates the source literally. This entails preserving the original source structure and translating the source word *ridiculous* directly. The result is somewhat awkward and sounds like translationese, which is particularly jarring considering the social media domain of the source text. This issue is rectified by both the refined translation and the reference. On the other hand, the refinement process is not perfect and does regress to less fluent outputs at times (example 4). Another prominent failure mode typical of all LLM outputs is hallucinations. In the fifth example, the refinement adds "touching the bruise", which is not present in the source.

---

[5]The exact sample ranges of the score difference are [-6, -2] and [1, 6]. Examples from beyond these ranges typically demonstrate clear signs of model degeneration and are therefore excluded from this analysis.

| | SOURCE SNIPPET | PRE-TRANSLATION RESEARCH | NOTES |
|---|---|---|---|
| 1. ✔ | trying my hand at miniatures redid the head! *[content omitted]* @user26 cheeked up. I love HER* | **"cheeked up"**: This slang term likely refers to making the miniature's cheeks fuller or more prominent. A direct translation wouldn't be accurate. | *Understanding source correctly in the context of making miniatures.* |
| 2. ✔ | #TeaganAir #pilot #flying #pa28 #PilotTraining *[content omitted]* Last week the weather didn't cooperate so we're back at it today. | **"The weather didn't cooperate"**: This is a figurative way of saying the weather was not suitable for flying. | *Good example of indirect translation.* |
| 3. ✔ | It is nice to have art and entertainment to help us make sense of our life. @user32 I could perform a liver transplant thanks to Hospital Playlist. | **"I could perform a liver transplant thanks to Hospital Playlist"**: This is a humorous exaggeration, implying that the show provided a lot of information about liver transplants. | *Good humor detection.* |
| 4. ✗ | He was covered in cuts and scratches and seemed a bit dazed, but he didn't seem too badly injured. | **"a bit dazed"**: This describes a state of confusion or disorientation. A literal translation might not convey the same meaning. | *Literal translation would do just fine.* |
| 5. ✗ | been really enjoying kneadatite (green stuff), it's like sculpting with bubble gum but you get used to that and it's tough and flexible when it cures | **"like sculpting with bubble gum"**: : This is a simile comparing the texture of Kneadatite to bubble gum. It needs to be translated in a way that makes sense to a Chinese audience who might not be familiar with the specific texture of bubble gum. | *Strange commentary on contextualizing the texture of bubble gum to a Chinese audience.* |

Table 6: Samples of pre-translation research outputs along with notes from the author-annotator.

## 7 Conclusion

We introduce a step-by-step approach to long-form text translation using LLMs. Inspired by literature on translation studies, we decompose the translation process into distinct stages, modeling pre-translation research, drafting, refinement, and proofreading though a multi-turn interaction with Gemini 1.5 Pro. Extensive automatic evaluations on WMT 2023 and WMT 2024 tasks in ten languages demonstrate that our approach improves translation quality over directly translating the entire document with a single prompt.

Furthermore, comparison with competitive baselines, including similar human-like LLM-driven approaches and top-performing systems that employ segment-by-segment translation of a document, reveals the strong performance of our approach. Our findings highlight the potential of LLMs to progressively improve their translations, moving beyond the traditional view of machine translation as a monolithic sequence mapping task.

## Limitations

While our study reveals promising step-by-step improvements across various languages and domains, we acknowledge the limitations of solely relying on automatic metrics for evaluation. While metric improvements give us a consistent signal, human evaluation is needed to further validate the effec-

tiveness of the approach and reveal a more nuanced understanding of the translation properties introduced at each step. We also acknowledge that our analysis is based solely on one family of metrics, due to context window limitations of other neural metrics in evaluating longer texts.

Finally, our pipeline is developed and tested solely on Gemini. Since different LLMs might exhibit different instruction-following capabilities across languages, the generalizability of this approach to other LLMs requires further investigation.

## Ethics Statement

This paper explores the use of LLMs to improve translation quality. In doing so, our approach starts from an initial translation that prioritizes faithfulness to the source text. Subsequent stages focus on improving fluency which, as they deviate more from the source, increase the risk of hallucinations (Guerreiro et al., 2023)—a critical issue in machine translation, potentially leading to misleading translations.

Moreover, the increasing fluency of machine translations presents new challenges when prioritized over adequacy (Martindale and Carpuat, 2018), as users might trust their outputs blindly, even when incorrect. This highlights the need for careful adoption of those translation systems and the developing of strategies that help users calibrate their trust appropriately.

# References

Duarte M. Alves, José P. Pombal, Nuno M. Guerreiro, Pedro H. Martins, Joao Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, Jos'e G. C. de Souza, and André Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *ArXiv*, abs/2402.17733.

Susan Bassnett and Peter R. Bush. 2016. The translator as writer. london: Continuum.

Claudine Borg. 2018. The phases of the translation process: are they always three?

Pinzhen Chen, Zhicheng Guo, Barry Haddow, and Kenneth Heafield. 2023. Iterative translation refinement with large language models. *ArXiv*, abs/2306.03856.

Daniel Deutsch, Juraj Juraska, Mara Finkelstein, and Markus Freitag. 2023. Training and meta-evaluating machine translation evaluation metrics at the paragraph level. In *Proceedings of the Eighth Conference on Machine Translation*, pages 996–1013, Singapore. Association for Computational Linguistics.

Birgitta Englund Dimitrova. 2005. Expertise and explicitation in the translation process.

Zhaopeng Feng, Yan Zhang, Hao Li, Wenqiang Liu, Jun Lang, Yang Feng, Jian Wu, and Zuozhu Liu. 2024. Improving llm-based machine translation with systematic self-correction.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Nuno M. Guerreiro, Duarte M. Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André F. T. Martins. 2023. Hallucinations in large multilingual translation models. *Transactions of the Association for Computational Linguistics*, 11:1500–1517.

Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shuming Shi, and Xing Wang. 2023. Exploring human-like translation strategy with large language models. *Transactions of the Association for Computational Linguistics*, 12:229–246.

Yi-Chong Huang, Xiaocheng Feng, Baohang Li, Chengpeng Fu, Wenshuai Huo, Ting Liu, and Bing Qin. 2024. Aligning translation-specific understanding to general understanding in large language models. *ArXiv*, abs/2401.05072.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Dayeon Ki and Marine Carpuat. 2024. Guiding large language models to post-edit machine translation with error annotations. *ArXiv*, abs/2404.07851.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinth'or Steingr'imsson, and Vil'em Zouhar. 2024. Preliminary wmt24 ranking of general mt systems and llms.

Shuang Li, Jiangjie Chen, Siyu Yuan, Xinyi Wu, Hao Yang, Shimin Tao, and Yanghua Xiao. 2024. Translate meanings, not just words: Idiomkb's role in optimizing idiomatic translation with language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):18554–18563.

Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.

Marianna Martindale and Marine Carpuat. 2018. Fluency over adequacy: A pilot study in measuring user trust in imperfect MT. In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 13–25, Boston, MA. Association for Machine Translation in the Americas.

Brian Mossop. 2000. The workplace procedures of professional translators.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the*

*Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.

Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem W. Ayoub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren Sezener, Luke Vilnis, Oscar Chang, Nobuyuki Morioka, George Tucker, Ce Zheng, Oliver Woodman, Nithya Attaluri, Tomás Kociský, Evgenii Eltyshev, Xi Chen, Timothy Chung, Vittorio Selo, Siddhartha Brahma, Petko Georgiev, Ambrose Slone, Zhenkai Zhu, James Lottes, Siyuan Qiao, Ben Caine, Sebastian Riedel, Alex Tomala, Martin Chadwick, J Christopher Love, Peter Choy, Sid Mittal, Neil Houlsby, Yunhao Tang, Matthew Lamm, Libin Bai, Qiao Zhang, Luheng He, Yong Cheng, Peter Humphreys, Yujia Li, Sergey Brin, Albin Cassirer, Ying-Qi Miao, Lukás Zilka, Taylor Tobin, Kelvin Xu, Lev Proleev, Daniel Sohn, Alberto Magni, Lisa Anne Hendricks, Isabel Gao, Santiago Ontan'on, Oskar Bunyan, Nathan Byrd, Abhanshu Sharma, Biao Zhang, Mario Pinto, Rishika Sinha, Harsh Mehta, Dawei Jia, Sergi Caelles, Albert Webson, Alex Morris, Becca Roelofs, Yifan Ding, Robin Strudel, Xuehan Xiong, Marvin Ritter, Mostafa Dehghani, Rahma Chaabouni, Abhijit Karmarkar, Guangda Lai, Fabian Mentzer, Bibo Xu, YaGuang Li, Yujing Zhang, Tom Le Paine, Alex Goldin, Behnam Neyshabur, Kate Baumli, Anselm Levskaya, Michael Laskin, Wenhao Jia, Jack W. Rae, Kefan Xiao, Antoine He, Skye Giordano, Lakshman Yagati, Jean-Baptiste Lespiau, Paul Natsev, Sanjay Ganapathy, Fangyu Liu, Danilo Martins, Nanxin Chen, Yunhan Xu, Megan Barnes, Rhys May, Arpi Vezer, Junhyuk Oh, Ken Franko, Sophie Bridgers, Ruizhe Zhao, Boxi Wu, Basil Mustafa, Sean Sechrist, Emilio Parisotto, Thanumalayan Sankaranarayana Pillai, Chris Larkin, Chenjie Gu, Christina Sorokin, Maxim Krikun, Alexey Guseynov, Jessica Landon, Romina Datta, Alexander Pritzel, Phoebe Thacker, Fan Yang, Kevin Hui, A.E. Hauth, Chih-Kuan Yeh, David Barker, Justin Mao-Jones, Sophia Austin, Hannah Sheahan, Parker Schuh, James Svensson, Rohan Jain, Vinay Venkatesh Ramasesh, Anton Briukhov, Da-Woon Chung, Tamara von Glehn, Christina Butterfield, Priya Jhakra, Matt Wiethoff, Justin Frye, Jordan Grimstad, Beer Changpinyo, Charline Le Lan, Anna Bortsova, Yonghui Wu, Paul Voigtlaender, Tara N. Sainath, Charlotte Smith, Will Hawkins, Kris Cao, James Besley, Srivatsan Srinivasan, Mark Omernick, Colin Gaffney, Gabriela de Castro Surita, Ryan Burnell, Bogdan Damoc, Junwhan Ahn, Andrew Brock, Mantas Pajarskas, Anastasia Petrushkina, Seb Noury, Lorenzo Blanco, Kevin Swersky, Arun Ahuja, Thi Avrahami, Vedant Misra, Raoul de Liedekerke, Mariko Iinuma, Alex Polozov, Sarah York, George van den Driessche, Paul Michel, Justin Chiu, Rory Blevins, Zach Gleicher, Adrià Recasens, Alban Rrustemi, Elena Gribovskaya, Aurko Roy, Wiktor Gworek, S'ebastien M. R. Arnold, Lisa Lee, James Lee-Thorp, Marcello Maggioni, Enrique Piqueras, Kartikeya Badola, Sharad Vikram, Lucas Gonzalez, Anirudh Baddepudi, Evan Senter, Jacob Devlin, James Qin, Michael Azzam, Maja Trebacz, Martin Polacek, Kashyap Krishnakumar, Shuo yiin Chang, Matthew Tung, Ivo Penchev, Rishabh Joshi, Kate Olszewska, Carrie Muir, Mateo Wirth, Ale Jakse Hartman, Joshua Newlan, Sheleem Kashem, Vijay Bolina, Elahe Dabir, Joost R. van Amersfoort, Zafarali Ahmed, James Cobon-Kerr, Aishwarya B Kamath, Arnar Mar Hrafnkelsson, Le Hou, Ian Mackinnon, Alexandre Frechette, Eric Noland, Xiance Si, Emanuel Taropa, Dong Li, Phil Crone, Anmol Gulati, S'ebastien Cevey, Jonas Adler, Ada Ma, David Silver, Simon Tokumine, Richard Powell, Stephan Lee, Michael B. Chang, Samer Hassan, Diana Mincu, Antoine Yang, Nir Levine, Jenny Brennan, Mingqiu Wang, Sarah Hodkinson, Jeffrey Zhao, Josh Lipschultz, Aedan Pope, Michael B. Chang, Cheng Li, Laurent El Shafey, Michela Paganini, Sholto Douglas, Bernd Bohnet, Fabio Pardo, Seth Odoom, Mihaela Rosca, Cicero Nogueira dos Santos, Kedar Soparkar, Arthur Guez, Tom Hudson, Steven Hansen, Chulayuth Asawaroengchai, Ravichandra Addanki, Tianhe Yu, Wojciech Stokowiec, Mina Khan, Justin Gilmer, Jaehoon Lee, Carrie Grimes Bostock, Keran Rong, Jonathan Caton, Pedram Pejman, Filip Pavetic, Geoff Brown, Vivek Sharma, Mario Luvci'c, Rajkumar Samuel, Josip Djolonga, Amol Mandhane, Lars Lowe Sjosund, Elena Buchatskaya, Elspeth White, Natalie Clay, Jiepu Jiang, Hyeontaek Lim, Ross Hemsley, Jane Labanowski, Nicola De Cao, David Steiner, Sayed Hadi Hashemi, Jacob Austin, Anita Gergely, Tim Blyth, Joe Stanton, Kaushik Shivakumar, Aditya Siddhant, Anders Andreassen, Carlos L. Araya, Nikhil Sethi, Rakesh Shivanna, Steven Hand, Ankur Bapna, Ali Khodaei, Antoine Miech, Garrett Tanzer, Andy Swing, Shantanu Thakoor, Zhufeng Pan, Zachary Nado, Stephanie Winkler, Dian Yu, Mohammad Saleh, Lorenzo Maggiore, Iain Barr, Minh Giang, Thais Kagohara, Ivo Danihelka, Amit Marathe, Vladimir Feinberg, Mohamed Elhawaty, Nimesh Ghelani, Dan Horgan, Helen Miller, Lexi Walker, Richard Tanburn, Mukarram Tariq, Disha Shrivastava, Fei Xia, Chung-Cheng Chiu, Zoe C. Ashwood, Khuslen Baatarsukh, Sina Samangooei, Fred Alcober, Axel Stjerngren, Paul Komarek, Katerina Tsihlas, Anudhyan Boral, Ramona Comanescu, Jeremy Chen, Ruibo

1311

Liu, Dawn Bloxwich, Charlie Chen, Yanhua Sun, Fangxiaoyu Feng, Matthew Mauger, Xerxes Dotiwalla, Vincent Hellendoorn, Michael Sharman, Ivy Zheng, Krishna Haridasan, Gabriel Barth-Maron, Craig Swanson, Dominika Rogozi'nska, Alek Andreev, Paul Kishan Rubenstein, Ruoxin Sang, Dan Hurt, Gamaleldin Elsayed, Ren shen Wang, Dave Lacey, Anastasija Ili'c, Yao Zhao, Woohyun Han, Lora Aroyo, Chimezie Iwuanyanwu, Vitaly Nikolaev, Balaji Lakshminarayanan, Sadegh Jazayeri, Raphael Lopez Kaufman, Mani Varadarajan, Chetan Tekur, Doug Fritz, Misha Khalman, David Reitter, Kingshuk Dasgupta, Shourya Sarcar, T. Ornduff, Javier Snaider, Fantine Huot, Johnson Jia, Rupert Kemp, Nejc Trdin, Anitha Vijayakumar, Lucy Kim, Christof Angermueller, Li Lao, Tianqi Liu, Haibin Zhang, David Engel, Somer Greene, Anais White, Jessica Austin, Lilly Taylor, Shereen Ashraf, Dangyi Liu, Maria Georgaki, Irene Cai, Yana Kulizhskaya, Sonam Goenka, Brennan Saeta, Kiran Vodrahalli, Christian Frank, Dario de Cesare, Brona Robenek, Harry Richardson, Mahmoud Alnahlawi, Christopher Yew, Priya Ponnapalli, Marco Tagliasacchi, Alex Korchemniy, Yelin Kim, Dinghua Li, Bill Rosgen, Kyle Levin, Jeremy Wiesner, Praseem Banzal, Praveen Srinivasan, Hongkun Yu, cCauglar Unlu, David Reid, Zora Tung, Daniel F. Finchelstein, Ravin Kumar, Andre Elisseeff, Jin Huang, Ming Zhang, Rui Zhu, Ricardo Aguilar, Mai Gim'enez, Jiawei Xia, Olivier Dousse, Willi Gierke, Soheil Hassas Yeganeh, Damion Yates, Komal Jalan, Lu Li, Eri Latorre-Chimoto, Duc Dung Nguyen, Ken Durden, Praveen Kallakuri, Yaxin Liu, Matthew Johnson, Tomy Tsai, Alice Talbert, Jasmine Liu, Alexander Neitz, Chen Elkind, Marco Selvi, Mimi Jasarevic, Livio Baldini Soares, Albert Cui, Pidong Wang, Alek Wenjiao Wang, Xinyu Ye, Krystal Kallarackal, Lucia Loher, Hoi Lam, Josef Broder, Daniel Niels Holtmann-Rice, Nina Martin, Bramandia Ramadhana, Daniel Toyama, Mrinal Shukla, Sujoy Basu, Abhi Mohan, Nicholas Fernando, Noah Fiedel, Kim Paterson, Hui Li, Ankush Garg, Jane Park, Donghyun Choi, Diane Wu, Sankalp Singh, Zhishuai Zhang, Amir Globerson, Lily Yu, John Carpenter, Félix de Chaumont Quitry, Carey Radebaugh, Chu-Cheng Lin, Alex Tudor, Prakash Shroff, Drew Garmon, Dayou Du, Neera Vats, Han Lu, Shariq Iqbal, Alexey Yakubovich, Nilesh Tripuraneni, James Manyika, Haroon Qureshi, Nan Hua, Christel Ngani, Maria Abi Raad, Hannah Forbes, Anna Bulanova, Jeff Stanway, Mukund Sundararajan, Victor Ungureanu, Colton Bishop, Yunjie Li, Balaji Venkatraman, Bo Li, Chloe Thornton, Salvatore Scellato, Nishesh Gupta, Yicheng Wang, Ian Tenney, Xihui Wu, Ashish Shenoy, Gabriel Carvajal, Diana Gage Wright, Ben Bariach, Zhuyun Xiao, Peter Hawkins, Sid Dalmia, Cl'ement Farabet, Pedro Valenzuela, Quan Yuan, Christoper A. Welty, Ananth Agarwal, Mianna Chen, Wooyeol Kim, Brice Hulse, Nandita Dukkipati, Adam Paszke, Andrew Bolt, Elnaz Davoodi, Kiam Choo, Jennifer Beattie, Jennifer Prendki, Harsha Vashisht, Rebeca Santamaria-Fernandez, Luis C. Cobo, Jarek

Wilkiewicz, David Madras, Ali Elqursh, Grant Uy, Kevin Ramirez, Matt Harvey, Tyler Liechty, Heiga Zen, Jeff Seibert, Clara Huiyi Hu, A. Ya. Khorlin, Maigo Le, Asaf Aharoni, Megan Li, Lily Wang, Sandeep Kumar, Alejandro Lince, Norman Casagrande, Jay Hoover, Dalia El Badawy, David Soergel, Denis Vnukov, Matt Miecnikowski, Jiri Simsa, Anna Koop, Praveen Kumar, Thibault Sellam, Daniel Vlasic, Samira Daruki, Nir Shabat, John Zhang, Guolong Su, Kalpesh Krishna, Jiageng Zhang, Jeremiah Liu, Yi Sun, Evan Palmer, Alireza Ghaffarkhah, Xi Xiong, Victor Cotruta, Michael Fink, Lucas Dixon, Ashwin Sreevatsa, Adrian Goedeckemeyer, Alek Dimitriev, Mohsen Jafari, Remi Crocker, Nicholas Fitzgerald, Aviral Kumar, Sanjay Ghemawat, Ivan Philips, Frederick Liu, Yannie Liang, Rachel Sterneck, Alena Repina, Marcus Wu, Laura Knight, Marin Georgiev, Hyo Lee, Harry Askham, Abhishek Chakladar, Annie Louis, Carl Crous, Hardie Cate, Dessie Petrova, Michael Quinn, Denese Owusu-Afriyie, Achintya Singhal, Nan Wei, Solomon Kim, Damien Vincent, Milad Nasr, Christopher A. Choquette-Choo, Reiko Tojo, Shawn Lu, Diego de Las Casas, Yuchung Cheng, Tolga Bolukbasi, Katherine Lee, Saaber Fatehi, Rajagopal Ananthanarayanan, Miteyan Patel, Charbel El Kaed, Jing Li, Jakub Sygnowski, Shreyas Rammohan Belle, Zhe Chen, Jaclyn Konzelmann, Siim Poder, Roopal Garg, Vinod Koverkathu, Adam Brown, Chris Dyer, Rosanne Liu, Azade Nova, Jun Xu, Slav Petrov, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *ArXiv*, abs/2403.05530.

Cy Shih. 2013. Translators' end-revision processing patterns and maxims: a think-aloud protocol study. *Arab World English Journal*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NeurIPS 2022, Red Hook, NY, USA. Curran Associates Inc.

Minghao Wu, Yulin Yuan, Gholamreza Haffari, and Longyue Wang. 2024. (perhaps) beyond human translation: Harnessing multi-agent collaboration for translating ultra-long literary texts. *ArXiv*, abs/2405.11804.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023a. A paradigm shift in machine translation: Boosting translation performance of large language models. *ArXiv*, abs/2309.11674.

Wenda Xu, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, Biao Zhang, Zhongtao Liu, William Yang Wang, Lei Li, and Markus Freitag. 2023b. Llmrefine: Pinpointing and refining large language models via fine-grained actionable feedback.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

# A Appendices

## A.1 Results on Shorter Documents

Table 9 presents automatic evaluation results of step-by-step on shorter documents, where segments are grouped together such that they do not exceed a token limit of 150 white-space separated tokens. The dataset statistics are presented in Table 10. We observe the same trends with the ones reported with larger documents in §5.

## A.2 Results on ChrF

Tables 7 and 8 report ChrF scores on WMT 2023 and WMT 2024, respectively. As anticipated with string-based metrics, LLM translations which prioritize fluency receive lower scores compared to those that are by construction instructed to be closer to the source text. This behavior is in line with observations of prior work that employ similar human-like translation strategies with LLMs (Wu et al., 2024).

## A.3 Prompts

Tables 11 and 12 present the complete prompts we used for our translate step-by-step framework and baselines. It has come to our attention that the prompts used in the experiments contain a few typographical errors. Preliminary results using revised prompts show comparable, if not slightly improved results (in the range of $0.1 - 0.2$ MetricX-23 score points), across all steps.

## A.4 More example outputs

Table 13 gives more example outputs to support the discussion in §6.

| | Research | Draft | Refinement | Proofreading | ZH | UK | RU | JA | HE | CS | DE | AVERAGE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1. | ○ | ○ | ○ | ○ | 48.04 | 61.85 | 63.55 | 38.75 | 64.03 | 67.62 | 71.81 | 59.38 |
| 2. | ○ | ● | ○ | ○ | 48.69 ↑0.65 | 61.81 ↓0.04 | 63.93 ↑0.38 | 39.00 ↑0.25 | 64.68 ↑0.65 | 67.63 ↑0.01 | 71.79 ↓0.02 | 59.65 |
| 3. | ○ | ○ | ● | ○ | 41.48 ↓6.56 | 59.44 ↓2.41 | 59.33 ↓4.22 | 36.19 ↓2.56 | 60.26 ↓3.77 | 63.44 ↓4.18 | 66.89 ↓4.92 | 55.29 |
| 4. | ○ | ● | ● | ○ | 43.14 ↓4.90 | 59.58 ↓2.27 | 60.37 ↓3.18 | 37.45 ↓1.30 | 60.92 ↓3.11 | 63.04 ↓4.58 | 68.71 ↓3.10 | 56.17 |
| 5. | ● | ● | ○ | ○ | 45.98 ↓2.06 | 61.51 ↓0.34 | 63.04 ↓0.51 | 39.30 ↑0.55 | 62.89 ↓1.14 | 67.17 ↓0.45 | 71.07 ↓0.74 | 58.71 |
| 6. | ● | ● | ● | ○ | 41.03 ↓7.01 | 58.72 ↓3.13 | 59.44 ↓4.11 | 37.65 ↓1.10 | 59.91 ↓4.12 | 63.02 ↓4.60 | 67.61 ↓4.20 | 55.34 |
| 7. | ● | ● | ● | ● | 40.71 ↓7.33 | 58.78 ↓3.07 | 59.23 ↓4.32 | 37.51 ↓1.24 | 59.65 ↓4.38 | 63.11 ↓4.51 | 67.49 ↓4.32 | 55.21 |

Table 7: ChrF evaluation results of translate step-by-step and its ablation variants on the WMT 2023 development datasets. Filled dots indicate active steps in the pipeline, while unfilled dots represent ablated steps. When all steps are ablated, the system defaults to zero-shot translation

| | DE | ES | ZH | RU | UK | JA | HI | IS | CS | AVERAGE |
|---|---|---|---|---|---|---|---|---|---|---|
| *Zero-shot* | 65.48 | 72.96 | 44.21 | 55.51 | 59.90 | 39.75 | 55.94 | 53.23 | 60.81 | 56.42 |
| *Research & Drafting* | 64.67 ↓0.81 | 72.30 ↓0.66 | 42.73 ↓1.48 | 57.30 ↑1.79 | 60.06 ↑0.16 | 41.19 ↑1.44 | 56.16 ↑0.22 | 53.09 ↓0.14 | 60.31 ↓0.50 | 56.42 |
| *Refinement* | 61.72 ↓3.76 | 69.22 ↓3.74 | 38.26 ↓5.95 | 55.09 ↓0.42 | 57.25 ↓2.65 | 39.15 ↓0.60 | 52.60 ↓3.34 | 52.62 ↓0.61 | 57.29 ↓3.52 | 53.69 |
| *Proofreading* | 61.62 ↓3.86 | 69.04 ↓3.92 | 38.41 ↓5.80 | 54.96 ↓0.55 | 57.14 ↓2.76 | 38.87 ↓0.88 | 53.47 ↓2.47 | 52.32 ↓0.91 | 56.98 ↓3.83 | 53.65 |

Table 8: ChrF results comparing step-by-step with zero-shot performance on the WMT 2024 test datasets.

| | DE | ES | ZH | RU | UK | JA | HI | IS | CS | AVERAGE |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | *Ref-based* |
| *Zero-shot* | 1.89 | 3.10 | 3.25 | 2.90 | 2.99 | 2.31 | 3.03 | 3.79 | 2.37 | 2.85 |
| *Research & Drafting* | 1.67 ↓0.23 | 2.61 ↓0.49 | 2.80 ↓0.45 | 2.53 ↓0.37 | 2.67 ↓0.32 | 1.91 ↓0.40 | 2.00 ↓1.03 | 3.45 ↓0.34 | 2.16 ↓0.21 | 2.42 |
| *Refinement* | 1.44 ↓0.45 | 2.20 ↓0.90 | 2.33 ↓0.92 | 2.17 ↓0.73 | 2.34 ↓0.65 | 1.61 ↓0.70 | 1.62 ↓1.41 | 3.02 ↓0.76 | 2.00 ↓0.37 | 2.08 |
| *Proofreading* | 1.36 ↓0.53 | 2.11 ↓0.99 | 2.28 ↓0.97 | 2.20 ↓0.69 | 2.27 ↓0.72 | 1.65 ↓0.66 | 1.60 ↓1.43 | 3.04 ↓0.75 | 1.98 ↓0.39 | 2.05 |
| | | | | | | | | | | *QE-based* |
| *Zero-shot* | 1.81 | 2.34 | 2.13 | 1.67 | 1.96 | 1.26 | 1.98 | 3.15 | 1.85 | 2.02 |
| *Research & Drafting* | 1.62 ↓0.18 | 2.03 ↓0.31 | 1.85 ↓0.28 | 1.40 ↓0.26 | 1.60 ↓0.36 | 1.10 ↓0.15 | 1.40 ↓0.58 | 2.93 ↓0.21 | 1.73 ↓0.12 | 1.74 |
| *Refinement* | 1.24 ↓0.57 | 1.61 ↓0.73 | 1.51 ↓0.62 | 1.05 ↓0.62 | 1.17 ↓0.79 | 0.91 ↓0.35 | 0.96 ↓1.01 | 2.37 ↓0.78 | 1.31 ↓0.53 | 1.35 |
| *Proofreading* | 1.12 ↓0.68 | 1.54 ↓0.80 | 1.44 ↓0.69 | 0.99 ↓0.67 | 1.10 ↓0.86 | 0.88 ↓0.38 | 0.92 ↓1.06 | 2.24 ↓0.90 | 1.22 ↓0.62 | 1.27 |

Table 9: MetricX-23 evaluation results comparing step-by-step with zero-shot performance on the WMT 2024 test datasets, where each document has a maximum length of 150 tokens. *Translate step-by-step surpasses zero-shot, with each step incrementally improving translation quality.*

| Domain | Literary | News | Social | Speech |
|---|---|---|---|---|
| # Docs. | 66 | 73 | 75 | 112 |
| Avg. Length | 120 | 110 | 105 | 72 |

Table 10: Per-domain statistics for WMT 2024, when blobbing with 150 max for total of 327 docs.

You will be asked to translate a piece of text form English into Chinese following the five stages of the translation process. Here is the context in which the text appears:

Context: *placeholder source text*

To start, let's do some pre-drafting research on the above context:

**Research:**
During this phase, thorough research is essential to address components of the context text that pose translation challenges. The goal is to establish a comprehensive translation plan that covers the following category:

* **Idiomatic Expressions:**

    * Identify idiomatic expressions that cannot be directly translated word-for-word into Chinese.

Now, let's move on to the drafting stage.

**Draft Translation:**
In this phase, your primary objective is to create a draft translation that accurately conveys the meaning of the source text presented below. At this stage, it is crucial to focus on adequacy, ensuring that your translation closely adheres to the source text. Your response should conclude with the draft translation. If context is missing, generate a general translation that is adaptable to various contexts. Avoid adding any additional information not present in the source text. All elements of the source text should be present in the translation.

Give your best one translation for the following piece of text based on the pre-drafting analysis without providing alternatives:

English: *placeholder source text*

Now let's move to the next stage.

**Post-editing with local refinement**
In this stage, the primary aim is to refine the draft translation by making micro-level improvements that improve the draft's fluency.

Provide only one refined translation and do not output anything else after that.

You are tasked with proofreading a translation that has been revised for improved fluency. The refined translation has been generated by editing the draft translation.

**Proofreading and Final Editing**
The goal is to provide a polished final translation of the source text. For you reference, below are the source text, the draft, and refined translations.

**Source Text**
*placeholder source text*

**Draft Translation**
*placeholder draft translation*

**Refined Translation**
*placeholder draft refined translation*

Please proofread the refined text for grammar, spelling, punctuation, terminology, and overall fluency. Ensure the translation accurately reflects the original meaning and style. Provide only the final, polished translation.

Table 11: Complete prompts used by the translate step-by-step pipeline.

You are asked to translate the text below into Chinese. Please output only the translation of the text without any other explanation.

English: *placeholder source text*
Chinese:

You are asked to translate the text below into Chinese. You are also given access to the context it appears.

Context: *placeholder document context*

Please output only the translation of the text without any other explanation.

English: *placeholder source text*
Chinese:

Analyze the previous responses and create a JSON object that organizes the linguistic information they contain. This object should have two sections: "idiomatic_expressions", and "draft_translation":

- **"idiomatic_expressions"**:
  - This section should also be a list of dictionaries.
  - Each dictionary represents a phrase and has the keys: "source_phrase", "description", "translation", and "literal_translation".
  - The "translation" key should hold a list of all provided translations for the phrase.
  - If the response doesn't provide a literal translation, use 'null' for the "literal_translation" value.
  - If the response doesn't identify relevant idiomatic expressions use 'null' for the corresponding value.
  - Don't include transliterated versions of the translations.

- **"draft_translation"**:
  - A string containing the draft translation.
  - If a translation includes multiple alternatives for the same source text, extract the first one.
  - Don't include transliterated versions of the translations.
  - If the translation includes alternatives separated by "/", extract the first one.

Table 12: Prompts used for baseline comparisons (zero-shot and zero-shot in context) and structuring the pre-translation research translation artifacts.

| | |
|---|---|
| **1. ✔Good example** | |
| **SOURCE** | @user26 cheeked up. I love HER* |
| **REFERENCE** | @user26 口出狂言。我喜欢她* |
| | *Talking nonsense. I like her** |
| **ZERO-SHOT** | @user26 脸红了。我爱她* |
| | *Blushing. I love her** |
| **RESEARCH AND DRAFT** | @user26 脸颊做得很饱满。我喜欢她* |
| | *The cheeks are very plump. I like her** |
| **REFINEMENT** | @user26 脸颊做得鼓鼓的，超可爱！我喜欢她* |
| | *The cheeks are very chubby, super cute! I like her** |
| **NOTES** | *Both the reference and the zero-shot misinterpreted the phrase "cheeked up" in the context of making miniatures. Both the draft and the refined version understood it correctly, with the latter being a bit more colloquial.* |
| **2. ✔Good example** | |
| **SOURCE** | I donked up my threading |
| **REFERENCE** | 我完成了这部分学习 |
| | *I finished this part of study.* |
| **ZERO-SHOT** | 我把我的线程搞砸了 |
| | *I messed up my threads [note: a terminology in computing].* |
| **RESEARCH AND DRAFT** | 我搞砸了我的帖子顺序 |
| | *I messed up the order of my posts.* |
| **REFINEMENT** | 我搞砸了帖子顺序 |
| | *I messed up the post order.* |
| **NOTES** | *The reference misinterpreted both "donked up" and "threading" in this context, and the zero-shot translation mistook the latter for the computing terminology.* |
| **2. ✔Good example** | |
| **SOURCE** | It is a ridiculous level. |
| **REFERENCE** | 这关太难了。 |
| | *This level is too hard.* |
| **ZERO-SHOT** | 这是一个荒谬的关卡。 |
| | *This is a ridiculous level.* |
| **RESEARCH AND DRAFT** | 这是一个荒谬的关卡。 |
| | *This is a ridiculous level.* |
| **REFINEMENT** | 这个关卡真是太难了。 |
| | *This level is really too hard.* |
| **NOTES** | *Zero-shot and draft translations are direct, literal, and stylistically awkward in Chinese, especially in the social media domain.* |
| **4. ✗ Bad example** | |
| **SOURCE** | Ivory chuckled from the memory. |
| **REFERENCE** | 艾弗里回忆到这里，笑了起来。 |
| | *Ivory thought of this and laughed* |
| **ZERO-SHOT** | 艾弗里想起这件事，不禁轻笑起来。 |
| | *Ivory thought of this and couldn't help starting to chuckle.* |
| **RESEARCH AND DRAFT** | 艾弗里回想起当时的情景，不禁轻笑出声 |
| | *Ivory thought of the situation at that time and couldn't help letting out a chuckle.* |
| **REFINEMENT** | 艾弗里摸了摸肩膀上上次留下的瘀伤，不禁轻笑出声。 |
| | *Ivory touched the bruise on the shoulder left from the last time and couldn't help letting out a chuckle.* |
| **NOTES** | *All translations did a good job translating the phrase "chuckled from the memory" fluently into Chinese, but the refinement hallucinated by adding "touched the bruise."* |
| **5. ✗ Bad example** | |
| **SOURCE** | one of them ran for something and grabbed it |
| **REFERENCE** | 其中一个人跑过去抓起了什么东西 |
| | *One of them ran over and grabbed something.* |
| **ZERO-SHOT** | 其中一个人跑去拿了什么东西 |
| | *One of them ran and grabbed something* |
| **RESEARCH AND DRAFT** | 其中一个人跑到某个东西那里，抓起它 |
| | *One of them ran to something and grabbed it.* |
| **REFINEMENT** | 其中一个人跑到某个东西那里，抓起它 |
| | *One of them ran to something and grabbed it.* |
| **NOTES** | *Although all translations are adequate and capture the semantic meaning of the source correctly, the draft and the refinement keep the original source structure and the resulting translation sounds like translationese.* |

Table 13: Sample of step-by-step and zero-shot outputs along with notes.

# Scaling Laws of Decoder-Only Models on the Multilingual Machine Translation Task

**Gaëtan Caillaut** and **Raheel Qader** and **Mariam Nakhlé** and **Jingshu Liu** and **Jean-Gabriel Barthélemy**

Lingua Custodia, Paris, France

`firstname.name@linguacustodia.com`

## Abstract

Recent studies have showcased remarkable capabilities of decoder-only models in many NLP tasks, including translation. Yet, the machine translation field has been largely dominated by encoder-decoder models based on the Transformer architecture. As a consequence, scaling laws of encoder-decoder models for neural machine translation have already been well studied, but decoder-only models have received less attention. This work explores the scaling laws of decoder-only models on the multilingual and multidomain translation task. We trained a collection of six decoder-only models, ranging from 70M to 7B parameters, on a sentence-level, multilingual (8 languages) and multidomain (9 domains) dataset. We conducted a series of experiments showing that the loss of decoder-only models can be estimated using a scaling law similar to the one discovered for large language models, but we also show that this scaling law has difficulties to generalize to *too large models* or to a different data distribution. We also study different scaling methods and show that scaling the depth and the width of a model lead to similar test loss improvements, but with different impact on the model's efficiency.

## 1 Introduction

Most modern machine translation systems are based on Transformers (Vaswani et al., 2017), with an encoder-decoder architecture. Despite the tremendous advances made possible with the release of open-source decoder-only Large Language Models (LLMs) (Jiang et al., 2023; Biderman et al., 2023; Touvron et al., 2023), most NLP tasks still rely on encoder-decoder models. Based on the statistics obtained from the WMT23 shared task on general machine translation (Kocmi et al., 2023), 16 out of the 17 participants submitted a system based on an encoder-decoder model. Yet, recent studies show that decoder-only models can achieve comparable results (Gao et al., 2022; Fu et al.,

2023), or even surpass state-of-the-art encoder-decoder systems, when properly finetuned (Xu et al., 2023). Moreover, the decoder-only architecture is easier to train on massive amounts of data as one can simply concatenate documents and feed as much relevant data as possible into the model during training ; while encoder-decoder models requires either to pad the inputs or rely on complex masking strategies (Raffel et al., 2020) to combine multiple inputs in the same sample.

Furthermore, the decoder architecture is much more flexible than the encoder-decoder architecture as decoders treat all tokens similarly, while encoder-decoders make a distinction between input (source) tokens and output (target) tokens, which are processed, respectively, by the encoder and the decoder. As a consequence, it is more tedious to apply complex *self-reasoning* mechanisms, such as chain-of-thought (Wei et al., 2022), or to interface it with external tools (Schick et al., 2024), because the outputs of such method (the *reasoning* process) should, preferably, be treated as inputs of the model. For the same reasons, it is much more computationally expensive to rely on an encoder-decoder for conversational purposes, making this architecture less efficient for modern workflows such as iterative translation. Indeed, at each round (the user's query and the system's answer) should be appended to the input side, and reprocessed by the encoder for the next round. Decoder-only models support it by design, without needing to recompute the representation of the ever-growing inputs. While we do not explore these directions in this work, we do leverage the flexibility of the decoder architecture to include input-or-output parameters. As we are tackling the multilingual and multidomain machine translation task, the model needs input tokens to represent the language direction and the domain. We propose to train the model to predict the source language and the domain so that, during inference, they can be seamlessly predicted or provided by

the user.

Generally speaking, decoder-only models simply expect the input to be the whole discussion and process it in a single forward step. Causal masking enable efficient caching of already computed keys and values so inference is much cheaper. The main downside of decoder-only over encoder-decoder model is the potential inferior quality of the input representation, as input tokens attend only on past tokens. But it should not be a major issue, as generated tokens attend to the whole past sequence, they do have access to the same quantity of information as with an encoder-decoder model. In addition, previous work propose to update the attention mask so that input tokens can attend to all input tokens while generated tokens can attend only on past tokens (Tay et al., 2022; Raffel et al., 2020).

For all these reasons, we would like to embrace the decoder architecture for machine translation, even if it seems to be the exclusive preserve of encoder-decoder models. The flexibility and the simpler training setup of decoders should make them both more suitable and efficient for most real world applications, and the decoder architecture is more appropriate to answer the ever-growing demand for iterative, interactive and machine assisted translation workflow. To this aim, we study the scaling laws of neural machine translation models under different settings. Our contributions are as follow:

- We show that decoder-only models for translation follow the same scaling law as LLM

- Scaling laws do not scale uniformly across directions and domains and do not generalize well to other directions or domains

- Scaling width-wise and depth-wise yield similar improvements, but the former is more efficient

- We discovered a critical issue related to the packing of training samples in batches and propose a solution to fix it

## 2 Background

As the size, data requirement, and training costs of language models rise, it quickly becomes critical to estimate the *right* training configuration for a given training budget — expressed in number of floating point operations (FLOP) — required to train the model. Kaplan et al. (2020) discovered a power law relationship between the loss of a language model and its number of parameters, and that larger models perform better given the same amount of data. Even though most work in this area show that larger models tend to be more powerful, recent studies show that other parameters must be taken into account as well. For instance, the Chinchilla scaling law (Hoffmann et al., 2022) shows that model and dataset sizes are loosely tied and need to be scaled equally. In other words, even if increasing only the model size will most likely improve its performances, the compute-optimal solution often requires to also increase the quantity of training data, while preserving the same training cost. These findings had a great impact on LLM research, as researchers stopped increasing blindly the size of their models, in favor of more data, when it was necessary. For instance, the 176B BLOOM (Le Scao et al., 2022) model would probably have been trained very differently (or not at all) if this study was released sooner. As stated in the paper, "in light of revised scaling laws published during training, we decided to train the large models for an additional 25 billion tokens on repeated data", the authors discovered that training such a big model was sub-optimal given the quantity of data they had. As a consequence, many researchers started to work on the collection of large, high quality datasets (Nguyen et al., 2024; Penedo et al., 2023) or on means to enhance existing datasets (Sorscher et al., 2022; Tirumala et al., 2024).

Most of these scaling laws studies focus exclusively on causal generative language models. While it's likely that many of these findings could apply to translation models, the differences between the two tasks cannot be taken for granted. Translation is a lot stricter than causal language modeling since the model has to take into account each information in the source and precisely generate the target sentence without adding or omitting any information. Hence, many studies have naturally emerged to observe the scaling behavior of translation models (Gordon et al., 2021; Fernandes et al., 2023; Ghorbani et al., 2021). Yet, these works focus on encoder-decoder models. For instance, Gordon et al. (2021); Fernandes et al. (2023) showed that, when the encoder and the decoder are scaled proportionally, the model's loss follow a power-law similar to the observation made on language models. Ghorbani et al. (2021) tackle the

problem in a different setup, and propose to scale the encoder and the decoder individually. They show that encoder-scaling and decoder-scaling affect the model's performances differently, and they propose a new formula describing the scaling behavior of the cross-entropy loss as a bivariate function of encoder and decoder size. They found out that scaling decoder is, according to their experiments, always more beneficial, in terms of cross-entropy loss performance, than scaling the encoder.

Recently, Alves et al. (2024) introduced the TowerInstruct, an LLM based on a decoder architecture (LLama 2 (Touvron et al., 2023)) finetuned to handle several translation tasks. They show that a properly finetuned LLM can perform translation better than state-of-the-art models on high-resource languages. But the most promising aspect of this work is the inherent capacity of LLM to handle different tasks. They finetuned TowerInstruct so it can, for instance, clean source sentences before translating them, follow terminological constraints or respect a given level of language. However, this work is still empirical and we do not know, yet, the limits of such models. Inspired by the performances of TowerInstruct, an LLM finetuned for machine translation tasks, we study, in the following, the scaling behavior of decoder-based machine translation models trained from scratch. To this aim, we fit multiple scaling laws to see if translation models follow the same scaling laws as language modeling models (such as the Chinchilla law) or if they follow their own task-specific law.

## 3 Training methodology

We present in this section all details related to the training of our six models.

### 3.1 Data

To conduct our experiments, we collected many bilingual data from public repositories (CCMatrix (Schwenk et al., 2021b), WikiMatrix (Schwenk et al., 2021a), UN Parallel Corpus (Ziemski et al., 2016), Paracrawl (Bañón et al., 2020) and Europarl (Koehn, 2005)). We also included a subset of an in-house proprietary dataset collected over time, as well as a small portion of financial documents in order to observe the scaling behavior on domain-specific data. An overview of the dataset distribution is given in Table 1. The financial data is divided into 8 sub-domains, which are described in Appendix A. The data is made of bilingual texts

with one sample being one sentence pair.

| Pair | Domain | Sentences | Tokens |
|------|--------|-----------|--------|
| en–de | general | 46.53 M | 2694.16 M |
|       | finance | 1.29 M | 65.93 M |
| en–es | general | 51.88 M | 3525.21 M |
|       | finance | 1.34 M | 71.48 M |
| en–fr | general | 81.39 M | 5430.77 M |
|       | finance | 8.29 M | 494.47 M |
| en–it | general | 26.21 M | 1657.58 M |
|       | finance | 0.73 M | 36.17 M |
| en–nl | general | 42.74 M | 2057.81 M |
|       | finance | 1.36 M | 63.96 M |
| en–pt | general | 42.02 M | 2086.62 M |
|       | finance | 0.61 M | 22.55 M |
| en–sv | general | 46.35 M | 2180.64 M |
|       | finance | 0.24 M | 9.68 M |
| fr–de | general | 23.60 M | 1470.68 M |
|       | finance | 1.46 M | 72.92 M |
| fr–es | general | 32.90 M | 2731.79 M |
|       | finance | 0.48 M | 23.39 M |
| fr–it | general | 28.02 M | 1845.84 M |
|       | finance | 1.10 M | 61.63 M |
| fr–nl | general | 31.94 M | 2034.74 M |
|       | finance | 0.62 M | 29.18 M |
| **Total:** | general | 453.58 M | 27 715.84 M |
|        | finance | 17.53 M | 951.36 M |
|        | all | 471.11 M | 28 667.20 M |

Table 1: Distribution of the training dataset. It covers 8 languages over 11 language pairs and 9 domains (general + 8 financial sub-domains).

We applied temperature sampling ($t = 5$) in order to increase the visibility of under represented pairs. Given a collection $\mathcal{D}$ of datasets, the probability of choosing a sample from a dataset $D_i \in \mathcal{D}$ after temperature sampling is given by $P_t(D_i)$ and is calculated from the original dataset statistical distribution $P(D_i)$.

$$P(D_i) = \frac{N_i}{\sum_{j=0}^{|\mathcal{D}|} N_j}$$

$$T(D_i, t) = P(D_i)^{1.0/t}$$

where $N_i$ is the size of dataset $D_i$ and $T(D_i, t)$ is the factor by which the dataset $D_i$ should be oversampled. The new size $k_i$ of the oversampled dataset $D_i$ is given by:

$$k_i = \left\lfloor \frac{T(D_i, t) \cdot \max_{j=0}^{|\mathcal{D}|}(N_j)}{\max_{j=0}^{|\mathcal{D}|}(T(D_j, t))} \right\rfloor$$

Finally, the probability of picking a sample from dataset $D_i$ after temperature sampling is given by

$$P_t(D_i) = \frac{k_i}{\sum_{j=0}^{|D|} k_j}$$

Since the balance between general and financial is also extremely skewed, we applied the temperature sampling separately on the general and financial domains.

### 3.2 Tokenizer

As we planned to train a multilingual model, we trained a Byte-Level BPE tokenizer (Wei et al., 2021) from scratch because, according to the authors, it is expected to better share the tokens among the multiple languages, resulting in less rare tokens and, hence, better embeddings. The tokenizer has been trained on the whole, non-oversampled, dataset, and we set the vocabulary size to 100 000.

We also reserved a small set of special tokens representing the supported languages and domains. They are inserted inside the input sequence so the model knows this information while generating a translation. For instance, the *English language token* is `<lang_en>` and the *general domain token* is `<dom_general>`.

### 3.3 Data format

Each sample of the datasets has two categories of features: inputs and outputs. Input features are data that will be given during inference, and output features are data that should be predicted by the model. Hence, inputs are the source sentence and the target language (because the model needs to know the desired target language); and outputs are the source language, the domain and the translated sentence.

Predicting the source language is not required, but we decided to include it to give to the model the ability to automatically detect the source language, as it is a very common and handy feature of most commercial translation tools. One could argue that this should be an input parameter, but we decided that the model should be able to classify by itself the language of the source sentence. Yet, the source language token can still be given as input at inference time to force a particular language. This also apply to the domain token.

Since we plan to train a decoder-only model, training samples have been formatted such that the input tokens are first seen by the model, so the model has access to the whole input when generating the first output token. This is why we chose to encode the sentence pairs in the following format:

SOURCE </src> <target lang> <source
 lang> <domain> TARGET <eos>

where `</src>` and `<eos>` are special tokens used to indicate, respectively, the end of the source and target sequences.

This data format gives the possibility to either provide the source language if required, or let the model predict it automatically. For instance, in the real example below, the green part represents the mandatory input (source sentence and target language), the blue part the optional input (source language) and the gray part is the output generated by the model.

The buyer pays at an ATM. </src>
 <lang_fr> <lang_en> <dom_general>
L'acheteur effectue le paiement sur les
 bornes automatiques. <eos>

#### 3.3.1 The <eos> token issue

All the models were trained in the same way LLM are trained. Sentence pairs were packed until the training batch was completely filled. These samples were separated by the usual end-of-sentence token `<eos>`. Ideally, one should also apply proper masking so tokens cannot attend to tokens from past sentence pairs. However, this features is not implemented in flash-attention 2 (Dao et al., 2022), so we trained the models without masks (except the causal mask). We expect the training task to be slightly more complex to solve, as the model now needs to learn to ignore every token before an `<eos>` token, but we decided that the gain in training speed is worthwhile.

| Model | without <eos> | with <eos> |
|---|---|---|
| 70M | 30.80 | 41.11 |
| 160M | 39.12 | 45.13 |
| 410M | 40.85 | 46.82 |

Table 2: BLEU scores of the same models when sources are prefixed with and without the `<eos>` token.

Our experiments showed that the quality of translations generated by the models were far below our expectations. We found that the absence of the `<eos>` token before the source sentence was confusing the model, explaining the drop in translation

quality shown in Table 2. The <eos> token, which was meant to signal the end of the translation, is actually also interpreted as a "start of translation" token. Indeed, during training, all sentence pairs (except the first one) are prefixed with the <eos> token. This phenomenon is clear in the example below, in which three sentence pairs are packed in the same training sample.

```
x₁₋₁ x₁₋₂ x₁₋₃ </src> <tgtlang₁> <srclang₁>
y₁₋₁ y₁₋₂ <eos> x₂₋₁ x₂₋₂ </src> <tgtlang₂>
 <srclang₂> y₂₋₁ y₂₋₂ <eos> x₃₋₁ x₃₋₂ x₃₋₃
</src> <tgtlang₃> <srclang₃> y₃₋₁ y₃₋₂ y₃₋₃
          y₃₋₄ <eos> x₄₋₁ . . .
```

The impact of <eos> absence on the test loss can be seen in Figure 1. The model clearly outputs better translation when the source sentence is prefixed with an <eos> token. This is particularly blatant when comparing the 160M and 410M models, respectively with and without the <eos> token prefix. The 410M model, albeit being more than two times bigger than the 160M model, cannot generate better translations without the <eos> prefix.



Figure 1: Test loss of our three smallest models (70M, 160M and 410M) with and without the <eos> prefix.

This problem should be negligible when training LLM, as documents are usually longer than sentence pairs, so <eos> tokens are scarcer. However, its impact will increase as batch size grows, since more sentence pairs can be packed into the same batch, making even more *obvious* that sentence pairs *should* start with an <eos> token. We experimented with a relatively small input length (512 tokens) and the absence of the <eos> token during inference already lead to significant drop in performance. Generally speaking, this issue should not be ignored when more than one sequence are packed in a single training sample. When possible, one should properly mask previous training sam-

ples. As it is not possible, currently, to leverage the state-of-the-art self-attention algorithms, we recommend to always prefix all source sentences with the same prefix token(s), both during training and inference. An alternative solution might be to prefix all sequences with a <bos> (begin of sentence) token, but we do not think it will solve this particular issue since the model will likely see that most sentence pairs start with the <eos><bos> sequence, which is still not the intended behavior. In the remaining of this paper, we will only consider translations generated with an <eos> prefix.

### 3.4 Training strategy

As we aim to train models dedicated to the translation task, we computed the loss only on target tokens, so the model learns to generate only text given a source sentence. This is different from pretrained language models as there is no notion of source and target sentence. The *target-only* strategy has proven to be effective for training text-to-text models (Touvron et al., 2023), and is also similar to the way loss of encoder-decoder models is calculated, which are commonly used for machine translation (Costa-jussà et al., 2022). Finally, we packed as many sentence pairs that we could in a single batch, in order to increase the training efficiency.

### 3.5 Model architectures

We used almost the same model architectures used in the Pythia suite (Biderman et al., 2023), the only difference being the number of attention head of the 160M model, as flash-attention expects a multiple of 8. We trained the models using the GPT-NeoX library (Andonian et al., 2023). We made a few changes to the data processing scripts in order to ignore source tokens during the loss computation. An overview of the different models we trained is given in Table 3.

All models are trained with a fixed batch size of 262 144 tokens (512 sequences of length 512 tokens) per GPU, on 8 Nvidia A100 GPUs. The models are trained in bfloat16 precision using the Adam optimizer with weight decay set to 0.1, 100 warmup steps and cosine learning rate decay. The maximum learning rate of sub-1B models is set to $1 \times 10^{-3}$, and $1 \times 10^{-4}$ for larger model because of loss instabilities during the training.

The models are trained for 100 000 steps on approximately 210B tokens, although only half of them were actually used to train the model as we

| Model | Non-embedding | Embedding | Layers | Dim | Heads | Max $lr$ |
|---|---|---|---|---|---|---|
| 70M | 70 295 552 | 51 380 224 | 6 | 512 | 8 | $1e^{-3}$ |
| 160M | 162 126 336 | 77 070 336 | 12 | 768 | 16 | $1e^{-3}$ |
| 410M | 405 071 872 | 102 760 448 | 24 | 1024 | 16 | $1e^{-3}$ |
| 610M | 607 448 064 | 154 140 672 | 16 | 1536 | 16 | $1e^{-3}$ |
| 1B | 1 011 257 344 | 205 520 896 | 16 | 2048 | 8 | $1e^{-4}$ |
| 6.9B | 6 855 204 864 | 411 041 792 | 32 | 4096 | 32 | $1e^{-4}$ |

Table 3: Architectures of the trained models. All models are trained with the very same setup (data, random seed, batch size, number of GPU, . . . ). They closely follow the Pythia models but parameters counts do no match because of the bigger vocabulary size, which increases the size of both the embedding and classification layer.

do not take into account source tokens when calculating the loss.

## 4 Experiments and results

In this section, we will study the impact of variations in training data size and parameters count on the test loss, for all our models. We will also verify if these changes correlate with their real translation performances using standard metrics such as BLEU and COMET. We finally explore two different model scaling strategies.

### 4.1 Applying machine translation scaling law



Figure 2: Test loss of all model checkpoints. Each step represents 512 training samples. Larger models always converge faster given the same amount of training data.

All existing scaling-laws studies show that larger models exhibit better generalization capabilities (Gordon et al., 2021; Fernandes et al., 2023; Ghorbani et al., 2021; Rae et al., 2021; Kaplan et al., 2020; Biderman et al., 2023). This study is no exception, as can be seen in Figure 2, larger decoder models always converge faster and require less training data to reach the same loss value.

We first fitted multiple curves following the setting of Ghorbani et al. (2021); Fernandes et al.

(2023), who studied scaling laws for machine translation. The form of the law is given below:

$$L(N) = \alpha N^{-p} + \beta \qquad (1)$$

where $N$ is the number of trainable parameters, and the other variables are fitted by minimizing the huber loss (with a delta value of $0.01$) using the BFGS algorithm from SciPy (Virtanen et al., 2020).

As shown in Figure 3, the test losses of our translation models can be realistically described by the power law fitted on observations made on all our models (the purple dotted line). This suggests that, indeed, performances of translation models follow a scaling law, that can be expressed by the formula above. We also fitted curves on less data points in order to verify if we could estimate the loss of the 6.9B model. Unfortunately, the fitted curves become deviate from the real observations as soon as we remove the data points from the largest model (the 6.9B model). This is extremely problematic, as the main goal of scaling laws is to estimate the performances of not-yet-trained larger models. Yet, we show that it is difficult to find a good estimation of the 6.9B model's performance without actually training it. For instance, the law fitted on the observations made on the subset 70M-160M-410M-610M-1B (in green) cannot give a good approximation of the *unseen* 6.9B model's performance, and the others are even worse. Therefore, we think one might be particularly cautious when applying such scaling laws to estimate larger models behaviors. Even if our law fitted on all data points seems to be a good estimator of the test loss, we think it will deviate from real observations as the model grows in size.

We also fitted scaling laws on a per-domain and per-direction basis, on all available data points. This is particularly interesting as it highlights discrepancies between domains and directions. As

Figure 3: Test losses estimated by power law fitted on different subset of models. Laws fitted on all models and 70M-160M-410M-1B models subset match our observations.



Figure 4: Scaling law fitted on the general domain and some financial subdomains. The law are fitted on the English-French direction only.



Figure 5: Scaling law fitted on the general domain for English-X direction.

shown in Figure 4, it seems to be significantly easier to translate sentences from the *kiid* (Key Investor Information Document) financial domain, but translating general domain sentences is the most difficult, even though the huge majority of our training set is from the general domain. We suspect this curve are, somehow, indicators of the diversity inside each domain. Indeed, *kiid* documents are, by law, all following the same structure and must contain a specific set of information, written in a certain way. On the contrary, *general* domain documents do not follow any rule, making this domain the most heterogeneous one, and thus the most difficult to translate. Other phenomena might explain the differences between these curves. For instance, we also think the presence of many very specific and rare words in the *regulatory* domain explains partly the lower translation quality in this domain.

We also fitted one curve per direction and observed similar phenomena, as shown in Figure 5. For example, our models seem to be better at translating from English to German than from English to French, although our training dataset contains twice as many English-French pairs (before oversampling).

These observations show that the scaling behavior of translation models depends on the training data distribution, and thus scaling laws estimated on a given dataset will not match the real scaling behavior on another one, although they might have the same general shape. For instance, it is not realistic to rely on a scaling law fitted on the en–fr direction to estimate the performances on the en–de direction.

## 4.2 Applying language modeling scaling law

So far, we experimented with a scaling law formula based on the model size only, ignoring the

training dataset size. Even if we just showed that lower perplexity/loss can be obtained with fewer data samples (in the case of the en–fr and en–de directions), larger training datasets still tend to increase the overall models' quality. But, it's also a waste of computing resources to train a model on more data than required, this is why modern language modeling scaling formula take into account both the number of trainable parameter and the training dataset size. Hence, we fitted multiple Chinchilla laws following the setting of Hoffmann et al. (2022), whose form is given below, on various combinations of input data to see if it can be used to reliably predict model performances.

$$L(N, D) = E + \frac{a}{N^\alpha} + \frac{b}{D^\beta} \qquad (2)$$

$E$, $a$, $\alpha$, $b$ and $\beta$ are variables fitted by minimizing the huber loss (with a delta value of $0.01$) using the BFGS algorithm from SciPy (Virtanen et al., 2020) ; $N$ and $D$ are, respectively, the number of non-embedding parameters of the model and the number of training samples. More details are given in the original paper.

As shown in Figure 6, the test loss of our translation models can be realistically described by the power law fitted on observations made on all our models (the purple dotted line). Furthermore, the

Figure 6: Test losses estimated by the Chinchilla law fitted on different model subsets. Curves deviate from the real observations when we remove too many data points to fit the curve.

general shape of the fitted curves is more stable, and thus more trustworthy. Indeed, the curve fitted on all models is very close to the one fitted without the 6.9B model, indicating that behaviors of larger models can be better estimated with this form of scaling law. However, as with the previous scaling law, the curve deviate from real observations when it is fitted on less data points. While it is not a surprising finding, it shows that scaling laws should not be trusted beyond a certain model size. However, we cannot provide a reasonable window in which the estimated loss is realistic.



Figure 7: Test loss of all models, each data point represents 5k training steps, or 2.5M samples. Given a fixed FLOP, it's often more beneficial to increase the dataset size when possible.

These experiments shows two things. First, the test loss of decoder-based translation models follows a scaling law similar to language modeling models, as the curves fitted on all data points match the real observations. The form of the law (a power law) indicates that larger models will always generalize better, until a certain point where the curve will stay mostly flat. The second thing we show is that finding a good and universal estimation for the model's loss is very difficult, as fitted curves do not generalize well beyond an unknown model size.



Figure 8: Estimation of models' test losses if they were trained on more data. According to the Chinchilla law fitted on all available observations, the 70M model should be on-par with the 410M performances with four times more data, and the 610M model should match the 6.9B model with only two times more data.

### 4.3 Correlating scaling law with real translation quality

Let us suppose we know the function modeling the real loss given a model size and an amount of training data. We still do not know if targeting lower loss values will actually improve the quality of the translations generated by the model. We provide in the following an empirical study showing the correlation between the model's loss and its translation performance. We computed BLEU (Papineni et al., 2002), COMET (Rei et al., 2022a) and CometKiwi (Rei et al., 2022b) scores for all six models, and we observed that, indeed, a lower loss does correlate with a performance increase, as shown in Table 4. This trend can be observed on the general domain for all directions, as shown in Appendix C. However, on the financial domain, CometKiwi does not always increase, it reaches a peak on the 610M model, then decreases. We conjecture that CometKiwi cannot correctly evaluate domain specific translations, as it is a reference-free model trained mainly on generalist sentences. We show in Appendix C that BLEU and COMET always increase with models' size, while CometKiwi often decreases at some point.

We also compare our models to well established LLM, and we show that smaller but specialized models clearly outperforms large and generalist LLM, as shown by our 410M model performing on par with `Llama 8B`. Our largest models are also real competitors to `Tower 7B`, even though it has been trained on much more data and specialized for machine translation. `Tower 7B` has the highest CometKiwi score, but as we just showed, it might

| Model | BLEU | COMET | CometKiwi |
|---|---|---|---|
| *General domain* | | | |
| 70M | 29.62 | 81.31 | 80.72 |
| 160M | 32.43 | 84.00 | 83.45 |
| 410M | 33.60 | 84.81 | 84.14 |
| 610M | 34.08 | 85.10 | 84.35 |
| 1B | 34.42 | 85.10 | 84.33 |
| 6.9B | 36.07$^\dagger$ | 85.88 | 84.82 |
| Llama3.1 8B | 30.43 | 84.82 | 84.47 |
| Mistral 7B | 23.26 | 80.08 | 82.29 |
| Tower 7B | 33.50 | 85.91$^\dagger$ | 85.02$^\dagger$ |
| Tower 7B$^*$ | 34.38 | 86.22 | 85.23 |
| *Financial domain* | | | |
| 70M | 44.63 | 86.95 | 80.88 |
| 160M | 49.02 | 88.27 | 81.80 |
| 410M | 50.85 | 88.64 | 81.73 |
| 610M | 52.00 | 88.85 | 81.71 |
| 1B | 53.28 | 89.98$^\ddagger$ | 81.61 |
| 6.9B | 58.34$^\ddagger$ | 89.62 | 81.35 |
| Llama3.1 8B | 34.99 | 84.42 | 81.75 |
| Mistral 7B | 38.93 | 76.52 | 76.17 |
| Tower 7B | 38.93 | 86.49 | 82.66$^\ddagger$ |
| Tower 7B$^*$ | 39.08 | 86.52 | 82.74 |

Table 4: Evaluation of the six models trained during this study on our in-house evaluation dataset. We reports both the scores on the general (G) domain and average over all financial (F) subdomains. We also include best performing LLM. As Tower has not been trained on Swedish, we also evaluate it after removing directions including Swedish (the Tower 7B$^*$ rows). Best scores on the general and financial domains are indicated by $^\dagger$ and $^\ddagger$ respectively.

not be reliable for specialized domains. Our models are obviously performing better on the financial domain, because only our models were finetuned on financial data. We also remark that Mistral's scores are quite low on the general domain, a quick manual inspection revealed that the model often give details and explanations about the produced translation, even when asked not to. As a consequence, we think that Mistral lower score is mostly caused by the model not following rigorously the instructions (see Appendix B).

So, while it certainly boost performances, increasing the model size is often not the optimal solution to improve the model's performance. The training dataset is also extremely important. Indeed, as can be observed in Figure 7, given a fixed FLOP budget, it is often preferable to increase the number of training samples. For instance, the 160M model appears to always be better than the 410M, 610M and 1B models given the same FLOP budget, as indicated by the 160M's curve being below other

models' curves. This observation is also validated by the fitted law, as indicated in Figure 8. Most of the time, and according to the fitted Chinchilla law, it would have been better to just train our models on more data, instead of training larger models. For instance, we estimate that the 160M model would be on-par with the 410M model if trained on approximately twice as many data, which would not exceed the total number of FLOP of our current 410M model.

To conclude with, we find that scaling laws are a powerful tool to have a glimpse of what we can expect from a *relatively larger* model trained on the same dataset, but it will probably fail to predict the performances of *much larger* models, even if trained on a similar data distribution. It has to be kept in mind when using such scaling laws to plan a training budget: **at some point, the fitted law will fail**. Planning a training budget based on observations made on a 10B model might be fine to train a 70B model, but completely wrong for a 500B one. Furthermore, a given scaling law can only estimate the end performances of a model trained on the same data distribution used to fit the scaling law. For instance, we show in Figures 4 and 5 that laws fitted on different language directions or domains are very different, and thus should not be applied to estimate the performances of the model on another direction.

### 4.4 Scaling strategies



Figure 9: In our experiments, we increased the width and the depth of the 70M model so the additional cost in terms of FLOP is similar (left). Scaling the depth or the width can lead to similar performance gains (right). The two figures are similar, except that the loss decrease can be observed either through the FLOP budget prism (left) or throughout training time / size of dataset (right).

We also studied whether one should favor scaling the depth (increasing the number of layers) or the width (increasing the hidden size) of a decoder model. We took the smallest model as a baseline and scaled it depth-wise and width-wise so that the increase in parameters increased the total training FLOP by a similar amount, as illustrated in

| Model | Layers | Dim | Non-embedding | Embedding | FLOP per s. | Samples per s. |
|---|---|---|---|---|---|---|
| 70M | 6 | 512 | 70 295 552 | 51 380 224 | $1.06 \times 10^{14}$ | 1170 |
| 70M+d768 | 6 | 768 | 119 599 104 | 77 070 336 | $1.74 \times 10^{14}$ | 900 |
| 70M+12l | 12 | 512 | 89 209 856 | 51 380 224 | $1.37 \times 10^{14}$ | 760 |
| 70M+d1024 | 6 | 1024 | 178 339 840 | 102 760 448 | $2.43 \times 10^{14}$ | 725 |
| 70M+24l | 24 | 512 | 127 038 464 | 51 380 224 | $1.6 \times 10^{14}$ | 445 |

Table 5: Sizes and architecture of models scaled in depth (70M+12l and 70M+24l) and models scaled in width (70M+d768 and 70M+d1024) compared to the base 70M model. Increasing the depth of the model has limited impact on the total parameters count, but decreases significantly the efficiency (higher FLOP per second but less training samples per second). Scaling the width of the model takes advantage of modern GPU architectures, but adds many trainable parameters.

Figure 9. An overview of the scaled model architectures can be seen in Table 5. Interestingly, we observed that both scaling methods yield the same performance improvement. As shown in Figure 9, given a similar FLOP cost, scaling the depth or the width seems to have the very same impact on the test loss.

Generally speaking, scaling depth-wise lead to smaller, but less efficient models. Indeed, modern hardware architecture can handle more efficiently large matrix products than many smaller matrix products. As shown in Table 5, width-scaled models are faster than depth-scaled models because the GPU can do more FLOP per second.

## 5 Conclusion

This work describes the behavior of decoder-only models on the multilingual multidomain machine translation task. We trained six models whose number of parameters range from 70M to 6.9B on sentence pairs in eight European languages. We found that scaling laws for machine translation cannot describe the general behavior of translation models, but they can still provide good estimation in a given domain, language-pair, and range of model sizes'. Indeed, We show that decoder-only models for translation tend to scale similarly as language models, as the Chinchilla law can also be applied to our models. As such, we recommend to train machine translation models using the same training recipes as large language models. While we think it is true for most, if not all, NLP tasks, more work need to be carried out to validate this hypothesis. However, we also highlight a critical limitation of scaling laws: they cannot generalize well beyond an unknown model and/or training dataset size. As models tend to be larger through time, it will be

extremely important to find ways to detect early unreasonable deviations of the "reference" scaling laws on which larger models are build.

We also show that models scaled width-wise appear to be more FLOP efficient than models scaled depth-wise, while reaching almost the same loss. Our experiments need to be continued in order to see when increasing the depth of the model starts to be more valuable than increasing its width. But, generally speaking, increasing the linearly both the depth and the width seems to be a good trade-off between efficiency and parameters count.

Efficient training requires packing as many sentence pairs as possible in a training batch. We discovered that unexpected biases can be introduced if proper masking is not applied, that is to say, if sequences can attend to previous ones. Since it is not possible with current state-of-the-art optimization methods, one must carefully format the training input data. We suggest dropping the *end-of-sentence* token, commonly used to signal the end of text generation, in favor of a *start-of-translation* token signaling the start of a new source sentence and, therefore, the end of the generated target sentence.

This study has been conducted on sentence-level pairs only. While this setup is a bit outdated, it is still the first time a comprehensive study has been made on multilingual machine translation using decoder-only architectures. Nevertheless, we expect decoder models to be easy to adapt to the document-level translation task, as one can simply finetune a sentence-level decoder with non-shuffled sentence pairs from a corpus of parallel documents.

## 6 Acknowledgment

## References

Duarte M. Alves, José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks.

Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric Hallahan, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler, Jason Phang, Shivanshu Purohit, Hailey Schoelkopf, Dashiell Stander, Tri Songz, Curt Tigges, Benjamin Thérien, Phil Wang, and Samuel Weinbach. 2023. GPT-NeoX: Large Scale Autoregressive Language Modeling in PyTorch.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359.

Patrick Fernandes, Behrooz Ghorbani, Xavier Garcia, Markus Freitag, and Orhan Firat. 2023. Scaling laws for multilingual neural machine translation. In *International Conference on Machine Learning*, pages 10053–10071. PMLR.

Zihao Fu, Wai Lam, Qian Yu, Anthony Man-Cho So, Shengding Hu, Zhiyuan Liu, and Nigel Collier. 2023.

Decoder-only or encoder-decoder? interpreting language model as a regularized encoder-decoder. *arXiv preprint arXiv:2304.04052*.

Yingbo Gao, Christian Herold, Zijian Yang, and Hermann Ney. 2022. Is encoder-decoder redundant for neural machine translation? In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 562–574.

Behrooz Ghorbani, Orhan Firat, Markus Freitag, Ankur Bapna, Maxim Krikun, Xavier Garcia, Ciprian Chelba, and Colin Cherry. 2021. Scaling laws for neural machine translation. *arXiv preprint arXiv:2109.07740*.

Mitchell A Gordon, Kevin Duh, and Jared Kaplan. 2021. Data and parameter scaling laws for neural machine translation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5915–5922.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, pages 30016–30030.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2023. Findings of the 2023 conference on machine translation (wmt23): Llms are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42.

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of machine translation summit x: papers*, pages 79–86.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model.

Thuat Nguyen, Chien Van Nguyen, Viet Dac Lai, Hieu Man, Nghia Trung Ngo, Franck Dernoncourt, Ryan A. Rossi, and Thien Huu Nguyen. 2024. CulturaX: A cleaned, enormous, and multilingual dataset for large language models in 167 languages. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4226–4237, Torino, Italia. ELRA and ICCL.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.

Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.

Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022a. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.

Ricardo Rei, Marcos Treviso, Nuno M Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José GC De Souza, Taisiya Glushkova, Duarte M Alves, Alon Lavie, et al. 2022b. Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task. *arXiv preprint arXiv:2209.06243*.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2024. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36.

Holger Schwenk, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021a. Wikimatrix: Mining 135m parallel sentences in 1620 language pairs from wikipedia. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Édouard Grave, Armand Joulin, and Angela Fan. 2021b. Ccmatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500.

Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. 2022. Beyond neural scaling laws: beating power law scaling via data pruning. *Advances in Neural Information Processing Systems*, 35:19523–19536.

Yi Tay, Mostafa Dehghani, Vinh Q Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Siamak Shakeri, Dara Bahri, Tal Schuster, et al. 2022. Ul2: Unifying language learning paradigms. *arXiv preprint arXiv:2205.05131*.

Kushal Tirumala, Daniel Simig, Armen Aghajanyan, and Ari Morcos. 2024. D4: Improving llm pretraining via document de-duplication and diversification. *Advances in Neural Information Processing Systems*, 36.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. 2020. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.

Junqiu Wei, Qun Liu, Yinpeng Guo, and Xin Jiang. 2021. Training multilingual pre-trained language model with byte-level subwords. *arXiv preprint arXiv:2101.09469*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*.

Michał Ziemski, Marcin Junczys-Dowmunt, and Bruno Pouliquen. 2016. The united nations parallel corpus v1. 0. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3530–3534.

## A    Full data distribution

Our models were trained on 11 language directions and 9 domains (8 are financial subdomains + general domain). The list 8 financial subdomains are given below:

**am** Asset Management

**ar** Annual Report

**corporateAction** Corporate Action Document

**equi** Equity Research

**ffs** Fund Fact Sheet

**kiid** Key Investor Information Document

**lifeInsurance** Life Insurance Document

**regulatory** Regulatory Document

## B    Prompt templates

We used the following system prompt to generate translation with Llama, Tower and Mistral:

> You are an expert translator. The user will ask you to produce translations, generate only the asked translation, do no justify nor explain anything.

We used the following instruction template to query the models:

> Translate from {SOURCE_LANGE} to {TARGET_LANG} the text below.
> {SOURCE_TEXT}

## C    Models' performances per direction

Performances of all models increase as parameters counts increase, regardless of the scoring method, as shown in Figures 10 and 11.



Figure 10: From top to bottom, BLEU, COMET and CometKiwi scores computed on the test dataset for all models and directions, on the general domain.

Figure 11: From top to bottom, BLEU, COMET and CometKiwi scores computed on the test dataset for all models and directions, averaged over all financial sub-domains.

# Shortcomings of LLMs for Low-Resource Translation: Retrieval and Understanding are Both the Problem

**Sara Court**
The Ohio State University
court.22@osu.edu

**Micha Elsner**
The Ohio State University
elsner.14@osu.edu

## Abstract

This work investigates the in-context learning abilities of pretrained large language models (LLMs) when instructed to translate text from a low-resource language into a high-resource language as part of an automated machine translation pipeline. We conduct a set of experiments translating Southern Quechua to Spanish and examine the informativity of various types of context retrieved from a constrained database of digitized pedagogical materials (dictionaries and grammar lessons) and parallel corpora. Using both automatic and human evaluation of model output, we conduct ablation studies that manipulate (1) context type (morpheme translations, grammar descriptions, and corpus examples), (2) retrieval methods (automated vs. manual), and (3) model type. Our results suggest that even relatively small LLMs are capable of utilizing prompt context for zero-shot low-resource translation when provided a minimally sufficient amount of relevant linguistic information. However, the variable effects of context type, retrieval method, model type, and language-specific factors highlight the limitations of using even the best LLMs as translation systems for the majority of the world's 7,000+ languages and their speakers.

## 1 Introduction

Despite great progress in the quality of today's state of the art machine translation (MT) systems, constraints on the amount and kinds of data available in the majority of the world's 7,000+ languages have led to yet another disparity in access and support for speakers of these languages: low-resource MT continues to be a major challenge (Hendy et al., 2023; Nicholas and Bhatia, 2023; Robinson et al., 2023; Stap and Araabi, 2023). Although many languages lack the kinds of large, standardized corpora necessary for traditional MT methods, recent work suggests it may be possible to leverage a smaller amount of existing resources, for example

pedagogical materials used for language instruction, to develop MT systems with Large Language Models (LLMs), albeit with varying results (Elsner and Needle, 2023; Tanzer et al., 2024; Zhang et al., 2024). These materials are often the result of community-driven or government-led initiatives to support language revitalization, reclamation, and mother-tongue education (Schreiner et al., 2020; Liu et al., 2022; Riestenberg et al., 2024). Such discrepancies in the needs and priorities of academic, commercial, and community-led efforts to develop digital resources and language technologies is what Gessler (2022) terms the "NLP Gap".

In this study, we investigate one way to lessen the NLP Gap, comparing LLMs' in-context learning abilities when translating from a low-resource language (a Peruvian variety of Southern Quechua) to a high-resource language (Spanish) using information retrieved from a database of pedagogical materials. We replicate results of earlier studies on a new language pair by comparing the effects of morpheme translations, sentences from a parallel corpus, and passages from a grammar instruction document on translation quality. We then conduct a more focused analysis by annotating translation outputs by hand using a modified MQM error typology (Burchardt, 2013). Finally, we conduct an ablation study on the effects of automated retrieval by manually constructing prompts using the same set of materials.

Our results suggest that while, unsurprisingly, translation quality improves with model size, such improvements seem to primarily be the result of previous exposure to the low-resource language during model pretraining, rather than an improved ability for the model to utilize prompt context, as evidenced by high scores in response to baseline (zero-shot) translation prompts. However, we also find evidence that in-context learning abilities may be inconsistent across different models of similar size. As found in previous studies, prompts contain-

ing morpheme and word-level translations reliably improve model outputs, but information from the grammar and corpus have a null or even negative effect on results. Human evaluation on a selection of outputs from two models – GPT-3.5 Turbo and GPT-4o – align with the quantitative measures we obtain using BLEURT (Sellam et al., 2020) as an automatic metric. Quantitative results also show an effect of automated retrieval on translation quality that is most evident in prompts containing morpheme translations and for models with lower baseline scores. Finally, we highlight a number of ethical concerns and limitations that arise from the proposed methods that are supported by our findings, and discuss the potential risks and challenges LLM-based methods for low-resource MT face moving forward.

## 2   LLMs for Machine Translation

Modern LLMs are now capable of translating many high-resource languages, but lack sufficient coverage of even modestly resourced languages to achieve comparable results without additional support (Kocmi et al., 2023). Retrieval-augmented generation (Rubin et al., 2022) may provide such support in the form of parallel sentences (Agrawal et al., 2022), dictionary definitions (Ghazvininejad et al., 2023; Lu et al., 2023) or other linguistic meta-knowledge such as a grammatical description. Retrieval-augmented methods offer exciting possibilities for low-resource translation, since the LLM might (in principle) be able to "teach itself" the language from learner-oriented resources produced by community members or language specialists.

Studies to date (Elsner and Needle, 2023; Reid et al., 2024; Zhang et al., 2024) experiment with four dimensions of variability: source language, LLM, type(s) of information retrieved, and retrieval method. Since the source languages in these studies have relatively little presence in public corpora or on the web, differing results across LLMs can tentatively be attributed to differences in their in-context learning and instruction-following abilities.

All studies find that word-level translations are helpful additions to prompts. Zhang et al. (2024) and Tanzer et al. (2024) also add sentence pairs from a parallel corpus, while Elsner and Needle (2023) add usage examples from a dictionary. Each improves results, although to a lesser degree. Elsner and Needle (2023) and Zhang et al. (2024) experiment with small fixed "grammar lesson" passages to provide explicit syntactic instruction, but find these ineffective. Tanzer et al. (2024) uses passages retrieved from a grammar book, also with relatively disappointing results. Reid et al. (2024) use the entire grammar book and a very long-context model to obtain better translations, but without exploring the role explicit grammar instruction actually plays in doing so.

Zhang et al. (2024) find that sentences from the corpus retrieved using BM25 embeddings (Robertson et al., 2009) work better than random ones. Tanzer et al. (2024), however, report that retrieval with longest common substring (LCS) matching outperforms embedding-based retrieval. Overall, the question of how to best retrieve relevant passages containing grammar material or sentences in a low-resource language is still open. This also complicates the interpretation of the mostly-negative results found for grammar passages. It is not clear whether these stem from poor retrieval, from the LLM's inability to process the retrieved content, or both. Moreover, although Reid et al. (2024) conducts human evaluation of the results for quality, to the best of our knowledge no study to date systematically investigates specific grammatical errors in the output.

Finally, each of these studies finds a significant decrease in LLMs' abilities to translate from a high-resource language into a low-resource language relative to experiments in the opposite direction. This is in line with McCoy et al. (2023), who find that while the accuracy of an LLM's output highly depends on the probability of both the input and the output text, output probability has a greater influence on model performance. We therefore focus this study on a single translation direction, instructing LLMs to output translations from a low-resource language into the language with which they are likely to have had more exposure during training, i.e, from Southern Quechua into Spanish, and leave the reverse direction for future work.

## 3   Quechuan Languages

Quechua is a family of languages Indigenous to the Andes in South America. This study focuses on varieties of Southern Quechua (S. Quechua, also known as *urin quechua* or *quechua sureño*) spoken in parts of Peru.[1] While previous studies investigated language-LLM pairs for which the baseline

---

[1] Unless noted otherwise, we use *Quechua* in this study to refer Southern Quechua and related varieties.

LLM lacked any pretrained knowledge, we find that newer LLMs can translate some S. Quechua sentences in a zero-shot setting. We expect this to be typical of many low-resource languages which, while often endangered, still may have some presence on the web.

Quechuan languages have by far the largest representation of all Indigenous Latin American languages in NLP research (Tonja et al., 2024) and are often included in ACL-affiliated workshops, datasets, and shared tasks (Cotterell et al., 2020; Ebrahimi et al., 2022, 2023). S. Quechua has a robust language toolkit (Rios, 2015), including the morphological parser we use in our pipeline. It has also been the subject of numerous studies on MT for both text and speech, developed in conjunction with monolingual and parallel corpora (Rios, 2015; Cardenas et al., 2018; Ortega et al., 2020; Zevallos et al., 2022). Nonetheless, such tools continue to face challenges, and Quechuan languages continue to lack the resources necessary to develop most of today's state of the art models.

Since Quechua is primarily spoken in South America, the majority of available digital resources, including all materials used in this study, use Spanish as the language of translation, explanation, and/or instruction. We therefore also use Spanish, rather than English or any other high-resource language, as the language of translation and prompting when testing our system.

## 3.1 Language-Specific Factors

While the proposed methods are general enough to be applied to any language pair, model outputs may reflect certain language-specific characteristics of the source and target languages, respectively. In this section, we provide a brief description of selected language-specific factors in S. Quechua as they relate to their translated Spanish counterparts. For a discussion of their potential effects on our results, please see Section 6.1.

### 3.1.1 Morphological Segmentation

S. Quechua is primarily agglutinating, i.e., much of the morphology may be described in terms of isomorphic form-meaning relationships, morphemes generally maintain a consistent form regardless of their phonological environment, and morpheme boundaries tend to be transparent. In contrast, morpheme segmentation in Spanish may be rendered opaque due to its fusional morphology and widespread use of conditioned allomorphy.

While LLMs are trained to process text via token-based rather than morpheme-based segmentation, it is possible that a lack of direct correspondence between the expression of morphosyntactic categories in S. Quechua and Spanish may affect a model's ability to leverage the information we provide as prompt context in our experiments. Correspondences in form and meaning across parallel usage examples may be particularly obscured, limiting the use of corpora designed for traditional MT methods. It may be possible to mitigate such issues with more advanced retrieval or prompting techniques, for example by explicitly instructing an LLM to conduct morphological analysis as part of the translation process, but we leave this for future work.

### 3.1.2 Syncretism and Polysemy

Although the language is primarily agglutinating, a number of morphemes in S. Quechua are syncretic, such that a given form may be used to express more than one grammatical category. For example, the 1SG.POSS marker, -*y*, shares the same form as both the 2SG.IMP marker and the infinitival marker, as illustrated in the following examples:

(1)    ñaña-y
       sister-1SG.POSS
       'my sister'

(2)    Mikhu-y!
       eat-2SG.IMP
       'Eat!'

(3)    mikhu-y-ta   muna-ni
       eat-INF-ACC want-1SG
       'I want to eat'

Similarly, words in S. Quechua may be polysemous, with the potential to express more than one meaning depending on their use in context. For example, the S. Quechua word *miski* (*misk'i*) may be translated as either *dulce* 'sweet' or *rico/a* 'delicious', and both *dulce* and *rico/a* are themselves polysemous in Spanish. *Dulce* may be used as an adjective, i.e., 'sweet', or a noun, i.e.,'candy', and *rico/a* may describe either richness in flavor, i.e., 'delicious', or in monetary wealth, i.e., 'rich'.

The exact forms displaying syncretism or polysemy must be identified on a language-specific basis, but the ambiguity they present poses a clear problem for our proposed methods in general, with

potential effects on both retrieval and generation. We discuss this issue further in Section 6.1.

### 3.1.3 Variation

Both S. Quechua and Spanish are characterized by extensive regional and dialectal variation. In S. Quechua, this includes differences in orthographic and/or phonological conventions as well as the specific lexical items and expression of morphosyntactic content. For example, the S. Quechua word for 'dog' may be rendered orthographically as *alqo*, *allqo*, *allku*, *allqu*, or *ashko*, and the additive suffix may be expressed as either *-pas* or *-pis*, depending on the community. Variation in the attested usage of specific lexical items and morphemes across communities is also common in S. Quechua. For example, the evidential marker *-mi /-m* is frequently attested in the Peruvian variety of S. Quechua used in this study, but essentially absent in many Bolivian varieties.

Variation across Spanish-speaking communities may also affect models' abilities to produce translations that are both accurate and appropriate. The Andean Spanish reference translations used in this study do not appear to affect the results of our automatic evaluation. However, were the proposed methods to be applied in a realistic setting, it would be especially important to assess the degree of alignment between any prescriptive linguistic standards that have been implicitly acquired by the LLM and the usage conventions of the language community or communities of interest.

## 4 Methods

### 4.1 Data

We conduct experiments on a collection of 50 pairs of S. Quechua - Spanish sentences sourced from one of the author's personal notes. These were selected to highlight a range of specific grammatical phenomena at multiple levels of difficulty— they include simple clauses and tenses (Example (4)), as well as more advanced constructions such as those involving past participles (Example (5)) and simultaneous events (Example (6).

(4)  qam allin-ta   tusu-nki
     you  good-ACC dance-2SG
     tu bailas bien
     'you dance well'

---

[TAREA] Traduce la siguiente frase del quechua al español. Responde sólo con la traducción:
quechua: kay wasiqa turiypam
español:

Figure 1: Example BASELINE prompt. English: *[TASK] Translate the following sentence from Quechua to Spanish. Respond only with the translation: Quechua: kay wasiqa turiypam; Spanish:*

(5)  awa-sqa-y   wali-qa   sumaq-mi
     knit-PP-1SG skirt-TOP great-ASSERTIVE
     la falda que tejí es linda
     'the skirt that I knit is pretty'

(6)  qam-qa    taki-ta     uyari-spa
     you-TOP song-ACC listen-SUBR
     wasi-yki-ta          picha-chka-nki
     house-2.POSS-ACC clean-PROG-2SG
     tú estás limpiando tu casa escuchando música
     'you're cleaning your house listening to music'

The first author, a foreign-language student of S. Quechua, received permission from her instructor to use notes from their lessons for the study. All sentence pairs were inspected by the instructor, a native bilingual speaker of both S. Quechua and Peruvian Spanish, to eliminate any errors and confirm the accuracy of all reference translations.

### 4.2 Prompt Construction

As a baseline, each sentence is inserted into a prompt template that instructs the model in Spanish to translate the S. Quechua sentence into Spanish and respond only with the translation (Figure 1). We automate a process for building on this template and compare the effects of adding information from three different sources to the prompt context.

#### 4.2.1 Morpheme Translations (MORPH)

We use a morphological parser (Rios, 2015) to segment each word of the source segment into morphemes, each with gloss symbols and a Spanish translation.[2] Some morphemes have multiple candidate meanings, all of which are retrieved. As an example, the word *rantikuq* is segmented as *ranti-ku-q* and glossed as "comprar.DB.VRoot-DB.VDeriv.+RflxInt+Ag.NS."

---

[2]We set aside valid concerns regarding the theoretical status of the *morpheme* for this study and define a morph(eme) loosely as a recognizable form-meaning pair that recurs in a language.

While numerous orthographic standards have been developed and promoted across Quechuan-speaking communities in South America, considerable variation in orthographic conventions may be found even within a particular community or variety (Rios and Castro Mamani, 2014). We discuss the implications of this for our results in Section 4.2.5.

We supplement the output from the parser using a Quechua-Spanish bilingual dictionary (Qheswa Simi Hamut'ana Kurak Suntur, 2005). We retrieve any dictionary entry whose headword exactly matches a morpheme in our segmentation. By default, we include all senses and any usage examples or contextual information in the dictionary entry as part of the prompt. We then concatenate the output of the parser with the retrieved dictionary entries and include this MORPH information as prompt context preceding the source sentence and baseline translation prompt.

### 4.2.2 Grammar Descriptions (GRAMMAR)

We also experiment with the inclusion of grammar lessons found in student-facing pedagogical materials, retrieving grammatical explanations relevant to each source sentence from a PDF document developed for students and teachers of S. Quechua (Pinto Tapia et al., 2005). The document is organized into short sections (1-3 sentences, plus paradigm tables or usage examples) that describe the particular grammatical concept associated with an affix in Quechua. For each source sentence, we retrieve sections associated with any affix listed in the document that is an exact match of a morpheme and include this in prompts using contextual information from the grammar. This improves on the methods described in Tanzer et al. (2024), who use LCS-based retrieval over an entire textbook, and Elsner and Needle 2023 and Zhang et al. 2024, whose grammatical description remains consistent across prompts regardless of the source text being translated.

### 4.2.3 Parallel Usage Examples (CORPUS)

Finally, we experiment with sentence-level examples from a S. Quechua-Spanish parallel corpus designed for traditional NLP tasks. We combine data made available via the AmericasNLP 2021 Shared Task on Open Machine Translation and the 2023 IWSLT shared task on low-resource SLT (Tiedemann, 2012; Agić and Vulić, 2019; Ortega et al., 2020; Mager et al., 2021; Agarwal et al.,

2023). For each source sentence, we retrieve the three best matches from the corpus using a LCS search against the full source sentence.

### 4.2.4 Combined Prompt Types

Combinations of information from all three sources yields 8 total conditions, including the baseline. An example prompt from each information source is given in Appendix E.

### 4.2.5 Manually Revised Prompts

To compute a soft upper bound on the improvements possible with better retrieval, we conduct an additional set of experiments using manually revised prompts. We first examine the content retrieved from the morphological parser, dictionary, and grammar document and remove all instances of ambiguity and irrelevant or misleading information from the prompt context.[3]

For example, many S. Quechua speakers use the term *runasimi* (lit: 'people mouth', 'the people's language'), as an endonym for the language. The parser, however, returns only the literal decomposition (*runa* 'ser humanos'/'people' and *simi* 'boca'/'mouth'), and the dictionary does not list *runasimi* as a headword but rather as one of eight different senses of *simi*. We thus remove all such irrelevant examples and translations from the prompt and retain only the content indicating a translation of *runasimi* in the linguistic sense.

We also manually retrieve content from the dictionary and grammar documents that were overlooked by the automated retriever. For example, the verb *yanuy* 'to cook' does not appear as a headword in the dictionary, but rather as a regional variant of *wayk'uy* 'to cook'. We also eliminate content from the grammar that was retrieved because of syncretism, or mistakes that cascaded from the morphological parser to result in irrelevant retrievals.

We manually parse each source sentence to only retrieve and include relevant information in the prompt context. All content in the revised prompts is sourced from the same material available to the automated retriever systems, and we do not add any additional information or use supplemental materials of any sort to create the revised prompts.

---

[3]We do not experiment with retrieval methods for corpus examples, which were retrieved using LCS match in both conditions. Improving on LCS-based retrieval remains an open question in low-resource LLM-MT, and we leave this for future work.

### 4.3 Models

We experiment with three proprietary models, GPT-3.5 Turbo (gpt-3.5-turbo-0125, Brown et al., 2020), GPT-4o (gpt-4o, Achiam et al., 2023), and Gemini 1.5 Pro (gemini-1.5-pro, Reid et al., 2024), and one open-source model, Llama 3 (llama-3-8b-instruct, AI@Meta, 2024). We use the pretrained models with their default settings, and do not adjust hyperparameters or conduct any finetuning as part of our experiments.

### 4.4 Evaluation

We conduct both automatic and human evaluations to identify trends in model errors and outputs in the various experimental conditions. We calculate BLEURT and BLEU scores as automatic metrics, and report mean BLEURT scores across items as the primary quantitative measure of translation quality for each of the conditions and models. We also use an adapted MQM schema to conduct qualitative human evaluation of the outputs of GPT-3.5 and GPT-4o for all prompts using automatic retrieval.

Each item selected for human evaluation is annotated by at least one of the authors by comparing the model's output to the source text and reference translation. We refer to the complete MQM typology to design our own four-dimensional framework of commonly attested errors in LLM-MT, each with a defined set of specific subtypes. Precise definitions and examples for all error categories and subtypes may be found in Appendix D.

Many of the categories in our schema are defined as in the core MQM framework. However, to capture some of the key behaviors reported in previous studies on LLM-MT and to evaluate the effects of prompt type on model outputs, we make the following adjustments. First, we utilize the Addition and Omission errors defined as Accuracy subtypes in the original MQM typology, but distinguish these from three additional subtypes: Substitution - Incorrect Subject, Substitution - Incorrect Tense/Aspect/Modality (TAM), and Substitution - Other. This is intended to capture LLM translations that differ from the source in terms of discrete lexical material or case, person, number, and/or TAM markings while otherwise maintaining the lexical and structural content needed to appropriately translate the source text. Although they are not the only grammatical phenomena that may be similarly misrendered, we select subject and TAM

markers for analysis as they are straightforward to identify and give a good indication of how well the LLMs cope with more abstract information about the meanings of functional morphemes.

Rather than including Mistranslation and MT Hallucination as Accuracy Error subtypes as in the original MQM typology, we define a separate Non-Translation category with three possible subtypes: Complete Mistranslation, Mistranslation with Lexical Correspondences, and Refusal. The third dimension of our typology, Model Error, was ultimately not used to classify any output in this study, but characterizes more generic model "misbehavior" such as failing to follow instructions, producing garbled text, or inappropriately generating content in the source language. Finally, Target Errors identify outputs that are ungrammatical, stylistically inappropriate, or semantically incoherent in the target language, regardless of their accuracy.

Detailed annotation guidelines were drafted and agreed upon to encourage consistency across annotators and experimental items. Annotators are instructed to identify and tag up to three specific errors for each translation output, with the exception of Target Errors, which do not count towards the three-error maximum. Each model output is also tagged for quality along a four-point scale as defined in Table 5.

Before proceeding with annotation over the larger dataset, both annotators also completed a test evaluation of the same 12 experimental items (96 sentences total) to assess inter-annotator agreement. Statistical measures ($\kappa = 0.72$ for quality judgments, $\alpha = 0.55$ for error categories) indicated some discrepancies in annotator judgments, especially for categories, since determining the three most important errors is especially subjective. These were identified and discussed, and agreement was ultimately deemed sufficient to proceed.

## 5 Results

### 5.1 Quality Metrics

We present BLEURT scores for prompts generated using automated retrieval in Table 1 and summarize human quality judgments for GPT-3.5 and GPT-4o with automated retrieval in Table 2. The complete distribution of BLEURT, BLEU, and human-annotated quality ratings for all of our experiments is provided in Appendix F. We find clear effects of LLM, prompt type, and retrieval method, as well as interactions among all three factors.

|        | GPT3.5 | GPT4o | Gem. | Lla3 |
|--------|--------|-------|------|------|
| BASE   | 0.19   | 0.66  | 0.56 | 0.15 |
| CORPUS | 0.27   | 0.59  | 0.49 | 0.19 |
| GRAM   | 0.23   | 0.56  | 0.55 | 0.17 |
| MORPH  | 0.44   | 0.54  | 0.61 | 0.39 |
| C+G    | 0.26   | 0.59  | 0.54 | 0.21 |
| C+M    | 0.44   | 0.59  | 0.59 | 0.36 |
| G+M    | 0.41   | 0.53  | 0.61 | 0.39 |
| C+G+M  | 0.43   | 0.57  | 0.61 | 0.15 |

Table 1: Mean BLEURT scores by LLM and prompt type. Shaded rows include morpheme contexts.

| LLM      | GPT-3.5 | GPT-4o |
|----------|---------|--------|
| BASELINE | 21      | 108    |
| CORPUS   | 43      | 101    |
| GRAMMAR  | 33      | 99     |
| MORPH    | 79      | 102    |
| C+G      | 41      | 101    |
| C+M      | 75      | 110    |
| G+M      | 68      | 100    |
| C+G+M    | 77      | 109    |

Table 2: Human-annotated quality ratings summarized as $3 \times high + 2 \times med + low$. Shaded rows include morpheme contexts.

|          | GPT3.5 | GPT4 | Gem. | Lla3 |
|----------|--------|------|------|------|
| G-AUTO   | 0.23   | 0.56 | 0.55 | 0.17 |
| G-MAN    | 0.24   | 0.58 | 0.54 | 0.15 |
| M-AUTO   | 0.44   | 0.54 | 0.61 | 0.39 |
| M-MAN    | 0.56   | 0.63 | 0.66 | 0.49 |
| CGM-AUTO | 0.43   | 0.57 | 0.61 | 0.15 |
| CGM-MAN  | 0.54   | 0.63 | 0.63 | 0.26 |

Table 3: Comparison of mean BLEURT scores for automatic versus manual retrieval of material in GRAMMAR, MORPH, and CORPUS-GRAMMAR-MORPH prompts.

Comparing across models, we find that Gemini and GPT-4o outperform Llama 3 and GPT-3.5 for every prompt type. This gap is highest for the least informative prompts, indicating that the Llama 3 and GPT-3.5 base models have relatively poor coverage of S. Quechua, while GPT-4o and Gemini have much better coverage. The effect is evident in both automatic and human quality evaluations.

Effects of prompt type are mediated by the quality of the pretrained model. Llama 3 and GPT-3.5 show a clear improvement in quality when MORPH information is included in the prompt. Gemini also improves when this information is added, but to a lesser extent. GPT-4o, on the other hand, performs best in response to the BASELINE (zero-shot) prompts, which attain the highest BLEURT scores across all models, prompt types, and retrieval methods evaluated in this study. In other words, providing additional information in the prompt's context actually *degrades* GPT-4o's ability to translate from S. Quechua to Spanish in all experimental conditions.

## 5.2 Effects of Automated Retrieval

To highlight the effects of automated retrieval on model output, we present BLEURT scores for a selection of prompt types and all four models in Table 3 (full scores may be found in Appendix F). The effect of manual retrieval for MORPH information is positive for all models, although this gap is smallest for Gemini (probably because its performance for these prompts is already highest). The effect for GRAMMAR prompts is either minor or negative.

## 5.3 Human Analysis of Translation Errors

The most common error type identified by the annotators is Substitution - Other, which includes a diverse assortment of lexical and phrasal incongruencies of varying degrees of severity. These are largely item-specific and therefore hard to characterize as a group. Using the error categories described in Section 4.4, we instead identify three more clearly interpretable phenomena and provide a detailed discussion of each in the following sections. We present counts for selected prompt types in Table 4, with examples in Appendix A and counts for all errors in Appendix G.

## 5.4 Mistranslations

Outright mistranslations are most common for GPT-3.5, making up 30 of the 50 responses in the BASE-

|  |  | BASE | MORPH | C+G+M |
|---|---|---|---|---|
| Mistranslation: complete + | GPT-3.5 | 45 | 11 | 12 |
| lexical correspondence | GPT-4o | 4 | 6 | 4 |
| Target Fluency: grammar + | GPT-3.5 | 0 | 14 | 10 |
| coherence + style | GPT-4o | 3 | 13 | 9 |
| Grammatical Divergence: | GPT-3.5 | 0 | 24 | 31 |
| subject + TAM | GPT-4o | 17 | 13 | 11 |

Table 4: Counts of human-annotated error types (per 50 sentences) by LLM and prompt type.

LINE condition. We also consider outputs that retain only minimal traces of the source content, which we label as Mistranslations with Lexical Correspondence. Approximately 1/3 of the 637 total errors tagged across all prompt types for GPT-3.5 are mistranslations of either type, roughly split between complete mistranslations and those with lexical correspondence (15.07% and 18.37%, respectively, of all errors tagged for GPT-3.5).

As reported in previous work, adding morpheme- and word-level translations to the prompt greatly reduces the rate of this kind of response. GPT-4o also produces drastically fewer mistranslations compared to its predecessor. However, it is notable that both models produce at least one mistranslation for each prompt type. In general, complete mistranslations are in fluent Spanish and contain no overt indications that something has been misrepresented. We return to the ethical implications of these errors in the Discussion.

We also note that many of the items tagged as Mistranslation with Lexical Correspondence show correspondence only for words that were already in Spanish in the source text. For example, some sentences contain Spanish loan names for the days of the week. While some of these errors are produced in deceptively fluent Spanish, we find many to be accompanied by semantic incoherence or ungrammaticality in the output. We discuss such target language fluency errors in the following section.

### 5.5 Target Fluency

Target Fluency errors occur when the output is not grammatical, coherent, or stylistically appropriate – for instance, if an output contains a nonsensical repetition or a verb with missing arguments. Outputs of this type bear a strong similarity to human "translationese" in that structural features of the source language may surface in the translation at the expense of naturalness (Koppel and Ordan, 2011; Freitag et al., 2019). Both GPT-3.5 and GPT-4o tend

to produce more such outputs when the prompt is more informative – 10 to 20% of the time (5-10 instances per 50) in prompts with morpheme translations.

### 5.6 Grammatical Divergence

We group misrendered verbal subjects and tense/aspect/morphology (TAM) markers together as Grammatical Divergence errors. Such errors are distinct from the Target Fluency errors described in the previous section— the Spanish output is grammatical, but fails to accurately reflect the content from the source. TAM divergences are much more prevalent than divergences in subject; for instance, only one of GPT-4o's 13 Grammatical Divergence errors in the MORPH condition misrender the subject marker.

Grammatical Divergence errors are annotated only for sentences that are not mistranslated outright, so GPT-3.5 produces none of these in the BASELINE condition. For more informative prompts, it is clear that GPT-4o is better than GPT-3.5 at translating both functional and lexical meanings. However, a relatively large number of sentences (over 20%) still contain such an error even with the highest performing model and prompt type. The relatively small drop in error between different prompt types for GPT-4o suggests that neither the corpus-based usage examples nor example paradigms and descriptions from the grammar document can fully prevent this type of error.

### 6 Discussion

We observe large differences between LLMs, both in terms of the overall quality of their generated translations as well as the effects of prompt type on their outputs. GPT-4o and Gemini, which have the highest baseline scores, benefit least from additional information— BLEURT scores actually decrease when CORPUS and GRAMMAR information is included. This occurs even with manually

curated prompts, suggesting it is not an effect of including irrelevant material. Nonetheless, the baseline results do not represent a ceiling on quality, since both models still produce errors in the BASELINE condition (GPT-4o produces 10 LOW-quality translations in our set of 50). These results suggest that even relevant grammar explanations, when written in prose with examples, do little to help the newest generation of LLMs to translate a low-resource language such as Southern Quechua.

Although GPT-4o and Gemini results are similar in many ways, we do find evidence for differences in their in-context learning abilities. Baseline prompts and the GPT-4o model produce the highest BLEURT scores across the dataset, but these outputs still show a number of errors characteristic of LLMs, particularly lexical substitution errors that are not necessarily corrected with the inclusion of more context. In contrast, Gemini, which has near-comparable performance across prompt types, shows an increase in scores when prompts include MORPH information, regardless of retrieval type, suggesting a greater ability to identify and utilize relevant word- and morph-level translations in the prompt's context. Previous work suggests that newer builds of GPT-4 are less capable of following instructions (Chen et al., 2023); such differences may be masked by the effects of pretraining when automatically evaluating translations. This suggests that researchers should continue to carefully select and compare among different LLMs when experimenting with retrieval-based translation.

## 6.1 Language-Specific Effects

We identify a number of translation errors of varying types that appear to be due to language-specific factors such as those discussed in Section 3.1. For example, we find an effect of polysemous lexical items for all prompt types in the outputs of both models on which we conduct human evaluation.

In the most straightforward cases, the model incorrectly generates an alternate sense of the word that is inappropriate given the content of the source sentence. We also find a number of instances in which the presence of such ambiguity has a cascading effect on lexical selection in the rest of the model's output. For example, when co-occuring with *miski* in the source text, the word *lawa* 'soup' is translated at times as *mazamorra*, a sweet porridge or pudding, *crema* 'cream', *miel* 'honey', *golosina* 'candy', or *dulces*, the nominal form of *dulce* meaning 'candy' or 'sweets.'

It may be possible to moderate such effects with additional refinement of the database structure and retrieval methods, which we leave for future work.

## 6.2 Ethical Concerns

Both our work and much of the previous work in this paradigm is motivated by the desire to close the "NLP Gap" among researchers, community members, and software developers interested in low-resource language technologies. Machine translation is listed as a welcome topic of research by some (though not all) members of American Indigenous communities (Mager et al., 2023), and is potentially an important tool for language learners (Jolley and Maimone, 2022). Even an imperfect translation system might be a useful tool for users with a clear understanding of its limitations. However, the systems evaluated in this work have two problematic tendencies that limit their potential for deployment in real community settings.

First, unfaithful translations often tend to be highly fluent (Section 5.4). While fluency ratings for older MT systems correlate well with accuracy scores, and have even been used as a proxy for overall translation quality (Gamon et al., 2005; Estrella et al., 2007), this correlation is reversed for our systems. LLMs are well-known for making false statements that seem plausible and authoritative (Bickmore et al., 2018; Dinan et al., 2021); this could be particularly problematic when they project illusions of expertise at the expense of an already marginalized group.

Second, some mistranslations identified in our study appear to draw on stereotypes of Indigenous groups (Appendix B). These are most apparent for the BASELINE system and GPT-3.5, but also (less frequently) occur with more informative prompts and better LLMs. Stereotypical sentences can involve flowery language with an emphasis on tradition or connectedness to nature (Erhart and Hall, 2019), as well as the unprompted addition of Indigenous Andean cultural customs and products (e.g., traditional medicine and chicha) to translations that are otherwise faithful to the source text. The overall effect is to exoticize Southern Quechua speakers and writers in ways that the original sentences do not. Similar stereotypes have also been noted in LLM-generated responses to open-ended prompts (Cheng et al., 2023; Delgado Solorzano and Toxtli, 2023; Shieh et al., 2024).

While we prompt models to output only the translation for evaluation purposes, models may have

some capacity to explain or qualify their translations and give reminders for responsible use of the technology. Should a retrieval-based translation system ever be deployed in a real-world setting for language learning, its developers should maximize transparency by presenting the content of any retrieved information and its source to the user along with the translation, reminding users directly of potential inaccuracies, and offering vetted resources for additional fact-checking when available.

# 7 Conclusion

Our results suggest a number of key limitations and concerns regarding the use of LLMs in a low-resource MT context, and have greater implications for our understanding of the seemingly "humanlike" conceptual, analytical, and in-context learning abilities of LLMs.

For the majority of the world's languages and their speakers, powering and supplying LLMs with enough pretraining data to overcome their limitations is not feasible. We therefore offer the following suggestions to those looking to develop low-resource LLM-MT: (1) improve data structures and methods for interacting with a language-specific database for retrieval-aided generation, (2) continue analysis of the mechanisms driving in-context learning in LLMs, for example by comparing ICL to the effects of finetuning (Dai et al., 2023), and (3) experiment with prompt structures and techniques, for example by altering the order of information (Liu et al., 2024) or by iteratively prompting the model to guide its reasoning towards a suitable translation (Wang et al., 2022).

Finally, we wish to emphasize the continued risks of prematurely deploying this or similar methods in any low-resource language community, particularly given the vulnerability and disproportionate lack of resources many such communities face in domains where these technologies would likely be used. As AI research continues to rapidly develop, we urge those conducting it to increase community engagement, amplify the voices of those traditionally at a disadvantage, and collaboratively develop research infrastructures that may lessen the NLP Gap (Brinklow, 2021). While there's still much to be done before low-resource LLM-MT may be safely implemented, we believe such a tool has the potential to empower speakers of any variety, including nonstandard varieties of high-resource languages such as English, to develop technologies that reflect their preferences and serve their unique needs.

# 8 Limitations

Limitations on the scope and replicability of this work may be attributed to one or more characteristics of the data and models used in this study, in addition to limitations inherent to the respective identities of its authors. First, the automatic metrics (i.e., BLEURT and BLEU scores) that we report are limited in their statistical validity. We have conducted some constrained tests to explore potential variance in scores, but expenses associated with text generation using proprietary models such as those developed by OpenAI and Google on a larger dataset may be prohibitive. This is compounded by the widely-acknowledged "black box" nature of the models powering both LLMs and BLEURT, as well as an increasing opacity with respect to the exact content and methods used to pretrain modern state of the art LLMs. For this reason, we focus our discussion on those results that show clear trends in both the quantitative and human evaluations we conduct.

There are also some constraints on our study and its methodology that are largely tied to linguistic factors, such as variation in orthography (and the need for digitized text-based resources as a prerequisite) as well as the lexical and grammatical variation that may be found in all languages, particularly the low-resource varieties we wish to support. We discuss some of these factors in Sections 3.1 and 6.1. Our results suggest it may be possible to guide the outputs of LLMs towards the specific usage conventions of a given community, but this is itself limited by the content of the materials used to develop the database from which prompt contexts are retrieved.

Neither of the authors is a native speaker of any Quechua or Spanish varieties, and only one is a student of these languages and has relationships to Quechua speakers and communities. While we have strived to be consistent in the Quechua and Spanish varieties used in our study (both the dictionary and grammar materials were provided by the same instructor who shared and proofread the 50 sentence pairs we use, and we select a morphological parser and corpora intended for use with Southern Quechua), variation is widespread among and within Quechua-speaking communities, and we do not have access to a dictionary, grammar,

morphological parser, and corpus developed by a unified and consistent set of authors. Future work should continue to explore ways to faithfully represent the diversity of linguistic conventions employed by communities interested in developing such technologies.

We acknowledge, as well, limitations that arise from the size of our dataset and database and the methods used to curate them. The 50 sentence pairs we use were selected to highlight a range of specific grammatical phenomena, not all of which were well represented in our database, and differ in their structural complexity. We are grateful for the guidance provided by the Quechua instructor whose lessons were a source for such examples and proofread the sentences before their inclusion in our experiments, but are limited by our status as non-native speakers. Human evaluation of model outputs was partially conducted using machine-translated English texts as references, but all annotations were inspected by the Spanish- and Quechua-speaking author who removed a small number of evaluations that reflected linguistic discrepancies between Quechua, Spanish, and English or inaccuracies in the machine-translated English.

## 9 Ethics Statement

We consulted the first author's Quechua instructor, Prof. Carmen Cazorla Zen, who gave us permission to use the sentences from the notes in this project and verified their accuracy. We cite the Quechua dictionary and grammar materials used to provide prompt information, and believe that our use of these materials is consonant with their original purpose. However, we do not distribute machine-readable versions of them as a contribution of this project, since this would violate the rights of the publisher. These materials were developed for use as pedagogical resources by institutions affiliated with the governments of Cuzco, Peru and Apurímac, Peru, respectively. Their authors were not contacted or consulted as part of the project.

We wish to acknowledge the delicate issue of academic *extractiveness* and its harmful impact on Indigenous and minority language communities and speakers. We are also aware of some of the controversial ideologies and policies associated with Qheswa Simi Hamut'ana Kuraq Suntur, the government-afilliated institution who published the dictionary we use in this study, and the potentially

negative effects of government-sponsored linguistic standardization more broadly (see, e.g., Coronel Molina (2008) for an analysis of the effects of the institution's ideologies on revitalization efforts in Peru). We do not endorse such policies, and have sought to avoid representing the diversity of Southern Quechua-speaking communities as a monolith. Instead, we hope our continued efforts to improve methods for low-resource translation will empower speakers of Southern Quechua and other Indigenous and minority languages to develop language technologies capable of representing their own community's unique language variety to serve the unique needs of its speakers.

There are numerous ethical issues related to the training and use of LLMs, such as labor issues and energy costs. While these issues are inextricable from the methods used in this project, we believe the potential impact of making low-resource translation viable and accessible to minority language communities who want them (our primary goal in this line of research) outweighs the problems inherent in using LLMs at all. We discuss the potential risks of deploying systems like the ones described here further in Section 6.2 of the main text.

## Acknowledgments

We thank Prof. Carmen Cazorla Zen, Professor of Quechua, for her help curating the data used in this study and for deepening our understanding of Southern Quechua and its speakers. We also thank Prof. Elvia Andía Grágeda, Professor of Quechua, for her instruction and advice, and the OSU Linguistics department for their feedback on a preliminary presentation of the work.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 technical report. arXiv preprint arXiv:2303.08774. See also: https://openai.com/index/hello-gpt-4o/.

Milind Agarwal, Sweta Agrawal, Antonios Anastasopoulos, Luisa Bentivogli, Ondřej Bojar, Claudia Borg, Marine Carpuat, Roldano Cattoni, Mauro Cettolo, Mingda Chen, William Chen, et al. 2023. FINDINGS OF THE IWSLT 2023 EVALUATION CAMPAIGN. In Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023), pages 1–61, Toronto, Canada (in-person and online). Association for Computational Linguistics.

Željko Agić and Ivan Vulić. 2019. JW300: A wide-coverage parallel corpus for low-resource languages. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 3204–3210, Florence, Italy. Association for Computational Linguistics.

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2022. In-context examples selection for machine translation. Preprint, arXiv:2212.02437.

AI@Meta. 2024. Llama 3 model card.

Timothy W Bickmore, Ha Trinh, Stefan Olafsson, Teresa K O'Leary, Reza Asadi, Nathaniel M Rickles, and Ricardo Cruz. 2018. Patient and consumer safety risks when using conversational assistants for medical information: An observational study of Siri, Alexa, and Google Assistant. J Med Internet Res, 20(9):e11510.

Nathan Thanyehténhas Brinklow. 2021. Indigenous language technologies: Anti-colonial oases in a colonizing (digital) world. WINHEC: International Journal of Indigenous Education Scholarship, (1):239–266.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Preprint, arXiv:2005.14165. See also: https://openai.com/index/new-embedding-models-and-api-updates/.

Aljoscha Burchardt. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In Proceedings of Translating and the Computer 35, London, UK. Aslib.

Ronald Cardenas, Rodolfo Zevallos, Reynaldo Baquerizo, and Luis Camacho. 2018. Siminchik: A speech corpus for preservation of Southern Quechua. ISI-NLP, 2:21.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023. Analyzing chatgpt's behavior shifts over time. In R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models.

Myra Cheng, Esin Durmus, and Dan Jurafsky. 2023. Marked personas: Using natural language prompts to measure stereotypes in language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1504–1532, Toronto, Canada. Association for Computational Linguistics.

Serafín M. Coronel Molina. 2008. Language Ideologies of the High Academy of the Quechua Language in Cuzco, Peru. Latin American and Caribbean Ethnic Studies, 3(3):319–340.

Ryan Cotterell, Christo Kirov, John Sylak-Glassman, Géraldine Walther, Ekaterina Vylomova, Arya D. McCarthy, Katharina Kann, Sabrina J. Mielke, Garrett

Nicolai, Miikka Silfverberg, David Yarowsky, Jason Eisner, and Mans Hulden. 2020. The CoNLL–SIGMORPHON 2018 shared task: Universal morphological reinflection. Preprint, arXiv:1810.07125.

Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. 2023. Why can GPT learn in-context? Language models implicitly perform gradient descent as meta-optimizers. Preprint, arXiv:2212.10559.

Cecilia Delgado Solorzano and Carlos Toxtli. 2023. Evaluating machine perception of Indigeneity: An analysis of ChatGPT's perceptions of Indigenous roles in diverse scenarios.

Emily Dinan, Gavin Abercrombie, A. Stevie Bergman, Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and Verena Rieser. 2021. Anticipating safety issues in E2E conversational AI: Framework and tooling. Preprint, arXiv:2107.03451.

Abteen Ebrahimi, Manuel Mager, Arturo Oncevay, Vishrav Chaudhary, Luis Chiruzzo, Angela Fan, John Ortega, Ricardo Ramos, Annette Rios, Ivan Vladimir Meza Ruiz, et al. 2022. AmericasNLI: Evaluating zero-shot natural language understanding of pretrained multilingual models in truly low-resource languages. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6279–6299, Dublin, Ireland. Association for Computational Linguistics.

Abteen Ebrahimi, Manuel Mager, Shruti Rijhwani, Enora Rice, Arturo Oncevay, Claudia Baltazar, María Cortés, Cynthia Montaño, John E. Ortega, Rolando Coto-solano, Hilaria Cruz, Alexis Palmer, and Katharina Kann. 2023. Findings of the AmericasNLP 2023 shared task on machine translation into Indigenous languages. In Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP), pages 206–219, Toronto, Canada. Association for Computational Linguistics.

Micha Elsner and Jordan Needle. 2023. Translating a low-resource language using GPT-3 and a human-readable dictionary. In Proceedings of the 20th SIGMORPHON workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 1–13, Toronto, Canada. Association for Computational Linguistics.

Ryan S Erhart and Deborah L Hall. 2019. A descriptive and comparative analysis of the content of stereotypes about Native Americans. Race and Social Problems, 11:225–242.

Paula Estrella, Andrei Popescu-Belis, and Maghi King. 2007. A new method for the study of correlations between MT evaluation metrics. In Proceedings of the 11th Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages: Papers, Skövde, Sweden.

Markus Freitag, Isaac Caswell, and Scott Roy. 2019. APE at scale and its implications on MT evaluation biases. In Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers), pages 34–44, Florence, Italy. Association for Computational Linguistics.

Michael Gamon, Anthony Aue, and Martine Smets. 2005. Sentence-level MT evaluation without reference translations: Beyond language modeling. In Proceedings of the 10th EAMT Conference: Practical applications of machine translation.

Luke Gessler. 2022. Closing the NLP gap: Documentary linguistics and NLP need a shared software infrastructure. In Proceedings of the Fifth Workshop on the Use of Computational Methods in the Study of Endangered Languages, pages 119–126, Dublin, Ireland. Association for Computational Linguistics.

Marjan Ghazvininejad, Hila Gonen, and Luke Zettlemoyer. 2023. Dictionary-based phrase-level prompting of large language models for machine translation. Preprint, arXiv:2302.07856.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? A comprehensive evaluation. Preprint, arXiv:2302.09210.

Jason R Jolley and Luciane Maimone. 2022. Thirty years of machine translation in language teaching and learning: A review of the literature. L2 Journal, 14(1):26–44.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, et al. 2023. Findings of the 2023 Conference on Machine Translation (WMT23): LLMs are here but not quite there yet. In Proceedings of the Eighth Conference on Machine Translation, pages 1–42, Singapore. Association for Computational Linguistics.

Moshe Koppel and Noam Ordan. 2011. Translationese and its dialects. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 1318–1326, Portland, Oregon, USA. Association for Computational Linguistics.

Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the Middle: How Language Models Use Long Contexts. Transactions of the Association for Computational Linguistics, 12:157–173.

Zoey Liu, Crystal Richardson, Richard Hatcher, and Emily Prud'hommeaux. 2022. Not always about you: Prioritizing community needs when developing endangered language technology. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3933–3944, Dublin, Ireland. Association for Computational Linguistics.

Hongyuan Lu, Haoyang Huang, Dongdong Zhang, Haoran Yang, Wai Lam, and Furu Wei. 2023. Chain-of-dictionary prompting elicits translation in large language models. Preprint, arXiv:2305.06575.

Manuel Mager, Elisabeth Mager, Katharina Kann, and Ngoc Thang Vu. 2023. Ethical considerations for machine translation of Indigenous languages: Giving a voice to the speakers. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4871–4897, Toronto, Canada. Association for Computational Linguistics.

Manuel Mager, Arturo Oncevay, Abteen Ebrahimi, John Ortega, Annette Rios, Angela Fan, Ximena Gutierrez-Vasques, Luis Chiruzzo, Gustavo Giménez-Lugo, Ricardo Ramos, et al. 2021. Findings of the AmericasNLP 2021 shared task on open machine translation for Indigenous languages of the Americas. In Proceedings of the First Workshop on Natural Language Processing for Indigenous Languages of the Americas, pages 202–217, Online. Association for Computational Linguistics.

R Thomas McCoy, Shunyu Yao, Dan Friedman, Matthew Hardy, and Thomas L Griffiths. 2023. Embers of autoregression: Understanding large language models through the problem they are trained to solve. arXiv preprint arXiv:2309.13638.

Gabriel Nicholas and Aliya Bhatia. 2023. Lost in translation: Large language models in non-english content analysis. Preprint, arXiv:2306.07377.

John E Ortega, Richard Castro Mamani, and Kyunghyun Cho. 2020. Neural machine translation with a polysynthetic low resource language. Machine Translation, 34(4):325–346.

Miguel Ángel Pinto Tapia, Luis Quispe Zúñiga, et al. 2005. Didáctica quechua i. Documento de trabajo, Dirección Regional de Educación Apurímac | Dirección Gestión Pedagógica.

Qheswa Simi Hamut'ana Kurak Suntur. 2005. Diccionario Quechua - Español - Quechua Qheswa - Español - Qheswa Simi Taqe, 2 edition. Multiservicios e Imprenta Edmundo Pantigozo EIRL, Cusco, Peru.

Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530.

Katherine J. Riestenberg, Ally Freemond, Brook Danielle Lillehaugen, and Jonathan N. Washington. 2024. Prioritizing Community Partners' Goals in Projects to Support Indigenous Language

Revitalization. In Decolonizing Linguistics. Oxford University Press.

Annette Rios. 2015. A basic language technology toolkit for Quechua. Ph.D. thesis, University of Zurich.

Annette Rios and Richard Castro Mamani. 2014. Morphological disambiguation and text normalization for southern Quechua varieties.

Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. Foundations and Trends in Information Retrieval, 3(4):333–389.

Nathaniel Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. ChatGPT MT: Competitive for high- (but not low-) resource languages. In Proceedings of the Eighth Conference on Machine Translation, pages 392–418, Singapore. Association for Computational Linguistics.

Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2655–2671, Seattle, United States. Association for Computational Linguistics.

Sylvia L.R. Schreiner, Lane Schwartz, Benjamin Hunt, and Emily Chen. 2020. Multidirectional leveraging for computational morphology and language documentation and revitalization. Language documentation and conservation, 14.

Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. BLEURT: Learning robust metrics for text generation. Preprint, arXiv:2004.04696.

Evan Shieh, Faye-Marie Vassel, Cassidy Sugimoto, and Thema Monroe-White. 2024. Laissez-faire harms: Algorithmic biases in generative language models. Preprint, arXiv:2404.07475.

David Stap and Ali Araabi. 2023. ChatGPT is not a good Indigenous translator. In Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP), pages 163–167, Toronto, Canada. Association for Computational Linguistics.

Garrett Tanzer, Mirac Suzgun, Eline Visser, Dan Jurafsky, and Luke Melas-Kyriazi. 2024. A benchmark for learning to translate a new language from one grammar book. Preprint, arXiv:2309.16575.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12), pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Atnafu Lambebo Tonja, Fazlourrahman Balouchzahi, Sabur Butt, Olga Kolesnikova, Hector Ceballos, Alexander Gelbukh, and Thamar Solorio. 2024. NLP progress in Indigenous Latin American languages. Preprint, arXiv:2404.05365.

Boshi Wang, Xiang Deng, and Huan Sun. 2022. Iteratively prompt pre-trained language models for chain of thought. Preprint, arXiv:2203.08383.

Rodolfo Zevallos, John Ortega, William Chen, Richard Castro, Núria Bel, Cesar Toshio, Renzo Venturas, Hilario Aradiel, and Nelsi Melgarejo. 2022. Introducing QuBERT: A large monolingual corpus and BERT model for Southern Quechua. In Proceedings of the Third Workshop on Deep Learning for Low-Resource Natural Language Processing, pages 1–13, Hybrid. Association for Computational Linguistics.

Chen Zhang, Xiao Liu, Jiuheng Lin, and Yansong Feng. 2024. Teaching large language models an unseen language on the fly. Preprint, arXiv:2402.19167.

## A  Example Errors

*The following section provides examples of errors analyzed in Section 5.3, one error per type.*

---

**Mistranslation: Complete Mistranslation**

---

**Model:** GPT-3.5 - BASELINE - AUTO

**Source:** qamqa taytaykipa munasqan lawata yanurqanki

**Gloss:**

qam-qa tayta-yki-pa          muna-sqa-n
you-TOP father-2.POSS-GEN like-SUBR-3SG
lawa-ta    yanu-rqa-nki
soup-ACC cook-PST-2SG

**Reference:** tú cocinaste la sopita que le gusta a tu papá
*you cooked the soup that your dad likes*

**Output:** yo te quiero mucho y mi amor por ti nunca se acabará
*I love you very much and my love for you will never end*

---

**Mistranslation: Lexical Correspondence**

---

**Model:** GPT-4o - C+G+M - AUTO

**Source:** ñuqaqa illani runakunawan rimayta munaptiy

**Gloss:**

ñuqa-qa illa-ni     runa-kuna-wan rima-y-ta
I-TOP    travel-1SG person-PL-with speak-ACC
muna-pti-y
like-SUBR-1SG

**Reference:** yo viajo porque me gustar hablar con las personas
*I travel because I like to talk to people*

**Output:** yo quiero hablar con las personas cuando amanece
*I want to talk to people when the sun rises*

---

**Target Fluency: Coherence**
*Unresolvable pronoun*

---

**Model:** GPT-3.5 - C+G+M -AUTO

**Source:** qamqa yachayta tukuptiyki, puñunki

**Gloss:**

qam-qa yacha-y-ta    tuku-pti-yki,    puñu-nki
you-TOP study-INF-ACC finish-SUBR-2SG, sleep-2SG

**Reference:** cuando termines de estudiar, tu duermes
*when you finish studying, you sleep*

**Output:** acabaré de aprender lo que tú duermes
*I will finish learning what you sleep*

---

**Target Fluency: Grammar**
*Missing causative and possessive*

---

**Model:** GPT-4o - CORPUS - AUTO

**Source:** ñuqaqa ñañayman wasita pichachini

**Gloss:**

ñuqa-qa ñaña-y-man           wasi-ta
I-TOP    sister-1SG.POSS-ILL house-ACC
picha-chi-ni
clean-CAUS-1SG

**Reference:** yo le hago limpiar la casa a mi hermana
*I make my sister clean the house*

**Output:** yo hermana casa limpio
*I sister house clean*

---

**Grammatical Divergence: TAM**
*Tense*

---

**Model:** GPT-3.5 - C+G+M - AUTO

**Source:** qam allinta tusunki

**Gloss:**

qam allin-ta     tusu-nki
you good-ACC dance-2SG

**Reference:** tu bailas bien
*you dance well*

**Output:** tú bailarás bien
*you will dance well*

## Grammatical Divergence: Subject
*'you' to 'I'*

**Model:** GPT-3.5 - MORPH - AUTO

**Source:** qamqa mamaykipa wasinta yanuq rinki

**Gloss:**

qam-qa mama-yki-pa wasi-n-ta
you-TOP mother-2.POSS-GEN house-3SG.POSS-ACC
yanu-q ri-nki
cook-AGT go-2SG

**Reference:** tú irás a cocinar a la casa de tu mamá
*you will go to cook at your mother's house*

**Output:** voy a cocinar en la casa de tu madre
*I'm going to cook at your mother's house*

## Grammatical Divergence: Subject
*'you' to 'the man'*

**Model:** GPT-4o - C+G+M - AUTO

**Source:** qamqa misita qawaspa qari warata awachkarqanki

**Gloss:**

qam-qa misi-ta qawa-spa qari wara-ta
you-TOP cat-ACC watch-SUBR male pants-ACC
awa-chka-rqa-nki
knit-PROG-PST-2SG

**Reference:** tú estabas tejiendo un pantalón de varón mirando al gato
*you were knitting a pair of men's pants looking at the cat*

**Output:** tú mirabas al gato mientras el hombre tejía el pantalón
*you looked at the cat while the man knitted the pants*

## Grammatical Divergence: TAM
*Missing conditional*

**Model:** GPT-4o - MORPH - AUTO

**Source:** qam runasimita yachapayaptiyki, allintam rimanki

**Gloss:**

qam runa-simi-ta yacha-paya-pti-yki,
you people-mouth-ACC learn-REP-COND-2SG,
allin-ta-m rima-nki
good-ACC-ASRT speak-2SG

**Reference:** si practicas el quechua lo hablarás bien
*if you practice Quechua you will speak it well*

**Output:** tú ya sabías el runasimi, hablas muy bien
*you already knew runasimi, you speak very well*

# B Unprompted Stereotypes

*The following section provides examples of errors involving cultural stereotypes as discussed in Section 6.2.*

---

**Model:** GPT-3.5 - BASELINE - AUTO

**Source:** payqa qatupi ñañanpa rantisqan apita rurarqan

**Gloss:**

pay-qa qatu-pi       ñaña-n-pa
3-TOP market-LOC sister-3SG.POSS-GEN
ranti-sqa-n       api-ta              rura-rqa-n
buy-SUBR-3SG mazamorra-ACC make-PST-3SG

**Reference:** ella preparó una mazamorra que compró su hermana en el mercado
*she prepared a mazamorra [pudding] that her sister bought at the market*

**Output:** el sol brilla en nuestra tierra, calienta y nos da vida
*the sun shines on our land, warms us and gives us life*

---

**Model:** GPT-3.5 - GRAMMAR

**Source:** paykunaqa yachachiqninkupa rakta medicina liwru apasqanta likanku

**Gloss:**

pay-kuna-qa yacha-chi-q-ni-nku-pa
3-PL-TOP       learn-CAUS-AGT-NI-3PL.POSS-GEN
rakta medicina liwru apa-sqa-n-ta              lika-nku
thick medicine book carry-SUBR-3SG-ACC read-3PL

**Reference:** ellos leerán el libro grueso de medicina que trajo su profesor
*they will read the thick medical book that their teacher brought*

**Output:** los médicos de la comunidad nos enseñaron a utilizar plantas medicinales
*the doctors from the community taught us how to use medicinal plants*

---

**Model:** GPT-4o - BASELINE

**Source:** qamkunapa taytamamaykichikta watuku-naykichik, paykunata purikuq apanaykichik

**Gloss:**

qam-kuna-pa tayta-mama-ykichik-ta
you-PL-GEN  father-mother-2PL.POSS-ACC
watuku-na-ykichik, pay-kuna-ta puri-ku-q
visit-must-2PL,       3-PL-ACC   walk-enjoy-AGT
apa-na-ykichik
bring-must-2PL

**Reference:** ustedes tienen que visitar a sus padres y llevarlos a pasear
*you have to visit your parents and take them for a walk*

**Output:** ustedes deben honrar a sus padres, guiarlos en su camino
*you must honor your parents, guide them on their path*

---

**Model:** Gemini 1.5 - BASELINE

**Source:** ñuqaqa upyayta munani

**Gloss:**

ñuqa-qa upya-y-ta         muna-ni
I-TOP     drink-INF-ACC want-1SG

**Reference:** yo quiero beber
*I want to drink*

**Output:** yo quiero beber chicha
*I want to drink chicha*

## C  Quality Descriptions

| Quality | Description |
|---------|-------------|
| High | Output is an accurate and/or acceptable translation of the source content. |
| Med | Output contains errors that prevent it from being an acceptable translation, but is generally high in quality otherwise. |
| Low | Output contains errors that prevent it from being an acceptable translation, with minor correspondences that vaguely identify it as relevant to the source. |
| None | Output does not appear to be relevant to the source. |

Table 5: Quality Descriptions

## D  Annotation Error Typology

| Dimension | Error | Description |
|-----------|-------|-------------|
| Accuracy | Addition | Translation includes information not present in the source, but does not result in the displacement of source content. |
| Accuracy | Omission | Translation is missing content from the source. |
| Accuracy | Substitution - Subject | The translated segment contains content identified as relevant to the source in other spans, but substitutes novel subject markers for those present in the source in the highlighted span; Classify an error as a "substitution" when the error appears to result in both Addition and Omission errors that cannot be distinguished into two distinct spans. |
| Accuracy | Substitution - TAM | The translated segment contains content identified as relevant to the source in other spans, but substitutes novel TAM for those present in the source in the highlighted span; Classify an error as a "substitution" when the error appears to result in both Addition and Omission errors that cannot be distinguished into two distinct spans. |
| Accuracy | Substitution - Other | Substitution errors that do not involve mistranslated subject markers or TAM. See above. |
| Accuracy | Overtranslation | Error occurring in the target content that is inappropriately more specific than the source content. |
| Accuracy | Undertranslation | Error occurring in the target content that is inappropriately less specific than the source content. |
| Target Error | Grammar | Other spans in the translated segment may be identified as relevant to the source, but the highlighted span is not grammatical in the target language. |
| Target Error | Coherence | Other spans in the translated segment may be identified as relevant to the source, but the highlighted span is unnatural or incoherent in the target language. |
| Target Error | Style/Register | Other spans in the translated segment may be identified as relevant to the source, but the highlighted span is produced in a style or register that is inappropriate given the content. |
| Non-Translation | Complete Mistranslation | The entire segment is coherent in the target language but the core predicate shows no immediate connection to the reference translation. |
| Non-Translation | Mistranslation - Lexical Correspondence | The entire segment is coherent in the target language but only minor correspondences to the reference translation may be identified. |
| Non-Translation | Refusal | Model does not attempt to translate into the target language, e.g., because it "does not understand". |
| Model error | Garbled | Output does not contain coherent text in the target language. |
| Model error | ChattyGPT | Output contains translated content, but is wordy, over-explanatory, and/or abruptly truncated. |

Table 6: Adapted MQM typology for human error annotation

## E   Example Prompts

The following are examples of prompts generated used automated retrieval from the database. English is included in italics for the reader, but was not provided to the models as part of the prompt.

---

**BASELINE**

---

[TAREA] Traduce la siguiente frase del quechua al español. Responde sólo con la traducción:
quechua: qam allinta tusunki
español:

– – – – – – – – – – – – – – – – – – – – – – – – – – – –

*[TASK] Translate the following sentence from Quechua to Spanish.  Respond only with the translation:*
*Quechua: You dance well*
*Spanish:*

---

**MORPHS-ONLY**

---

[CONTEXTO]
qam: [PrnPers+2sg]
allin: bueno [D̂B][NRoot]
ta: [+Acc][Cas]
tusu: bailar [VRoot][D̂B]
nki: [+2sg.Subj][VPers]
allin. adj. Bueno (término de aprobación). SINÓN: kusa. EJEM: allin p'unchay, buenos días: allin tuta, buenas noches; allin tutamanta, buena mañana, buenos días; allin inti chinkay, buenas tardes; allin iñiyniyoq, de buena fe, fiel, justo, íntegro: allin nunayoq, de espíritu bueno; allin puriq, de comportamiento bueno; allin puriy, comportamiento bueno; allin rikuy, tratamiento bueno; allin rikuq, el que trata bien; allin ruway, obrar bien, beneficiar; lo que se hace bien, beneficioso; allin ruwaq, el que hace bien; allin yuyay, pensar bien; pensamiento bueno; allin qolqeyoq, poseedor de plata fina; adinerado.
ta. s. Gram. Sufijo que desempeña los papeles de artículo y preposición. EJEM: llamata qatiy, arrea la llama; Urkusmanta hamuni, vengo de Urcos.

[TAREA] Traduce la siguiente frase . . .

– – – – – – – – – – – – – – – – – – – – – – – – – – – –

*[CONTEXT]*
*qam: [PrnPers+2sg]*
*allin: bueno [D̂B][NRoot]*
*ta: [+Acc][Cas]*

---

*tusu: bailar [VRoot][D̂B]*
*nki: [+2sg.Subj][VPers]*
*allin. adj. Good (term of approval).  SYN: kusa. EX: allin p'unchay, good day:  allin tuta, good evening; allin tutamanta, good morning, good day; allin inti chinkay, good afternoon; allin iñiyniyoq, good faith, faithful, just, upright: allin nunayoq, in good spirits; allin puriq, well behaved; allin puriy, good behavior; allin rikuy, good treatment; allin rikuq, one who treats others well; allin ruway, to do good, to benefit; one who does good, beneficial; allin ruwaq, one who does good; allin yuyay, think well; good thought; allin qolqeyoq, possessor of fine silver; wealthy.*
*ta.  s.  Gram.  Suffix that plays the roles of article and preposition. EX: llamata qatiy, herd the llama; Urkusmanta hamuni, I come from Urcos.*

*[TASK] Translate the following sentence . . .*

---

**GRAMMAR-ONLY**

---

[CONTEXTO]
ta: CASO ACUSATIVO. Su marca es –ta, esta es una marca de objeto directo con los verbos que no son de movimiento (quietud). Ejemplo:
Quyllur–ta qhawani Veo una estrella
T'anta–ta apay Lleva pan
Ñuqa quylluyta qhawani
Pedrucha t'antata rantin
En cambio con los verbos de movimiento –ta indica (hacia) que es igual a meta. Ejemplos:
Punu–ta rini Voy a Puno
Llaqta-ta risaq Iré al pueblo
Hamawt'anchis Punuta rinqa
Llanta umalliq llaqtata richkan
nki: FLEXIÓN DE TIEMPO. TIEMPO FUTURO. TIEMPO FUTURO. Los sufijos para cada una de las personas gramaticales son: saq, nki, nqa, sun, saqku, nkichis, nqaku; en singular y plural respectivamente.
Ejemplos:
Puklla-saq jugaré
Puklla-nki jugarás
Puklla-nqa jugará
Puklla-sun jugaremos
Puklla-saqku jugaremos
Puklla-nkichis Uds. jugarán
Puklla-nqaku ellos jugarán

[TAREA] Traduce la siguiente frase . . .

1350

[TAREA] Traduce la siguiente frase . . .

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

*[CONTEXT]*

*quechua: rimanakunapaq wawakunapa rimasqan simi aswan allinta takyachinaraq piwanpas maywanpas mana manchakuspa rimananpaq chaymi qillqanapaqpas ñawichanapaqpas aswan allin kanqa*

*Spanish: For this dialogue, knowing the language that the children speak would be important for them to express themselves without fear, and that is why writing and reading will be optimal.*

*quechua: kay tiqsipi sumaq rimanakunapaqa kawsayninchikmi allinta kallpachawanchik runaku-nahina allinta tiyanapaq chaymi ñuqanchikkqa allinta ñawichayta qillqayta yachananchik ñawpa ayllunchikkuna rurasqankuta maytukunapi tukuy puyñukunapi tiqsi muyu qhawarisqankuta*

*Spanish:To live in harmony we have to know our way of life well and then write and read to also value what our ancestors left us in each vision of the world.*

*quechua: winsislawcha chayarqamuptinsi tu-parquspanku allinta qatunakusqanku suwakuypi purinankupaq*

*español: cuando había llegado wenseslau y a su encuentro se habían reforzarón para andar a robar*

*[TASK] Translate the following sentence . . .*

---

*[CONTEXT]*

*ta: ACCUSATIVE CASE. Marked by –ta, this is a direct object marker with verbs that don't indicate movement. Example:*

*Quyllur–ta qhawani I see a star*

*T'anta–ta apay Bring bread*

*Ñuqa quylluyta qhawani*

*Pedrucha t'antata rantin*

*On the other hand, with verbs of motion -ta indicates (towards) the same goal Examples:*

*Punu–ta rini I go to Puno*

*Llaqta-ta risaq I will go to town*

*Hamawt'anchis Punuta rinqa*

*Llanta umalliq llaqtata richkan*

*nki: TENSE INFLECTION. FUTURE TENSE. FUTURE TENSE. The suffixes for each of the grammatical persons are: saq, nki, nqa, sun, saqku, nkichis, nqaku; in singular and plural respectively.*

*Ejemplos:*

*Puklla-saq jugaré*

*Puklla-nki jugarás*

*Puklla-nqa jugará*

*Puklla-sun jugaremos*

*Puklla-saqku jugaremos*

*Puklla-nkichis Uds. jugarán*

*Puklla-nqaku ellos jugarán*

*[TASK] Translate the following sentence . . .*

========================================

**CORPUS-ONLY**

========================================

[CONTEXTO]

quechua: rimanakunapaq wawakunapa rimasqan simi aswan allinta takyachinaraq piwanpas maywanpas mana manchakuspa rimananpaq chaymi qillqanapaqpas ñawichanapaqpas aswan allin kanqa

español: para este diálogo saber la lengua que dominan los niños sería importante para que ellos se expresen sin miedo de ahí será que la escritura y la lectura salga de manera óptima

quechua: kay tiqsipi sumaq rimanakunapaqa kawsayninchikmi allinta kallpachawanchik runaku-nahina allinta tiyanapaq chaymi ñuqanchikkqa allinta ñawichayta qillqayta yachananchik ñawpa ayllunchikkuna rurasqankuta maytukunapi tukuy puyñukunapi tiqsi muyu qhawarisqankuta

español: para vivir en armonía tenemos que conocer bien nuestra forma de vivir y luego escribir leer tambien a valorar lo que nos dejaron

## F Full Quality Scores

This section contains tables showing all automatic and human-annotated quality scores for each of our experiments. Table 7 contains the full set of BLEURT scores summarized in Tables 1 and 3 of the main text. Table 8 shows the corresponding BLEU scores for the same experiments. Table 9 and Table 10 contain the full set of the human-annotated scores summarized in Table 3.

| | GPT-3.5 | | GPT-4o | | Gemini-1.5 | | Llama 3 | |
|---|---|---|---|---|---|---|---|---|
| | auto | manual | auto | manual | auto | manual | auto | manual |
| BASELINE | 0.19 | 0.22 | 0.66 | 0.66 | 0.56 | 0.57 | 0.15 | 0.16 |
| CORPUS-ONLY | 0.27 | 0.29 | 0.59 | 0.61 | 0.49 | 0.47 | 0.19 | 0.18 |
| GRAMMAR-ONLY | 0.23 | 0.24 | 0.56 | 0.58 | 0.55 | 0.54 | 0.17 | 0.15 |
| MORPH-ONLY | 0.44 | 0.56 | 0.54 | 0.63 | 0.61 | 0.66 | 0.39 | 0.49 |
| CORPUS-GRAMMAR | 0.26 | 0.28 | 0.59 | 0.59 | 0.54 | 0.53 | 0.21 | 0.21 |
| CORPUS-MORPH | 0.44 | 0.52 | 0.59 | 0.64 | 0.59 | 0.64 | 0.36 | 0.38 |
| GRAMMAR-MORPH | 0.41 | 0.54 | 0.53 | 0.61 | 0.61 | 0.64 | 0.39 | 0.37 |
| CORPUS-GRAMMAR-MORPH | 0.43 | 0.54 | 0.57 | 0.63 | 0.61 | 0.63 | 0.15 | 0.26 |

Table 7: BLEURT scores for all LLMs and prompt types.

| | GPT-3.5 | | GPT-4o | | Gemini-1.5 Pro | | Llama 3 8B | |
|---|---|---|---|---|---|---|---|---|
| | auto | manual | auto | manual | auto | manual | auto | manual |
| BASELINE | 0.01 | 0.02 | 0.19 | 0.18 | 0.12 | 0.14 | 0.00 | 0.00 |
| CORPUS-ONLY | 0.02 | 0.02 | 0.16 | 0.22 | 0.14 | 0.13 | 0.02 | 0.01 |
| GRAMMAR-ONLY | 0.01 | 0.03 | 0.14 | 0.12 | 0.18 | 0.17 | 0.01 | 0.01 |
| MORPHS-ONLY | 0.06 | 0.08 | 0.12 | 0.13 | 0.15 | 0.18 | 0.03 | 0.05 |
| CORPUS-GRAMMAR | 0.01 | 0.01 | 0.14 | 0.17 | 0.12 | 0.08 | 0.01 | 0.01 |
| CORPUS-MORPHS | 0.05 | 0.08 | 0.19 | 0.18 | 0.17 | 0.17 | 0.02 | 0.04 |
| GRAMMAR-MORPHS | 0.03 | 0.04 | 0.11 | 0.10 | 0.15 | 0.16 | 0.02 | 0.01 |
| CORPUS-GRAMMAR-MORPHS | 0.04 | 0.04 | 0.16 | 0.16 | 0.17 | 0.20 | 0.00 | 0.01 |

Table 8: BLEU scores for all LLMs and prompt types.

### GPT-3.5 Turbo

| | None | Low | Med | High |
|---|---|---|---|---|
| BASELINE | 31 | 17 | 2 | 0 |
| CORPUS-ONLY | 18 | 23 | 8 | 1 |
| GRAMMAR-ONLY | 20 | 27 | 2 | 1 |
| MORPHS-ONLY | 3 | 22 | 16 | 9 |
| CORPUS-GRAMMAR | 18 | 23 | 9 | 0 |
| CORPUS-MORPH | 2 | 28 | 12 | 8 |
| GRAMMAR-MORPH | 3 | 29 | 13 | 5 |
| CORPUS-GRAMMAR-MORPH | 2 | 27 | 12 | 9 |

Table 9: Human quality annotation of GPT-3.5 outputs with automated retrieval (raw counts out of 50) by prompt type.

**GPT-4o**

|  | None | Low | Med | High |
|---|---|---|---|---|
| BASELINE | 0 | 10 | 20 | 20 |
| CORPUS-ONLY | 1 | 16 | 13 | 20 |
| GRAMMAR-ONLY | 0 | 17 | 16 | 17 |
| MORPHS-ONLY | 0 | 13 | 18 | 19 |
| CORPUS-GRAMMAR | 0 | 14 | 17 | 19 |
| CORPUS-MORPH | 0 | 10 | 17 | 23 |
| GRAMMAR-MORPH | 0 | 19 | 14 | 17 |
| CORPUS-GRAMMAR-MORPH | 0 | 9 | 20 | 21 |

Table 10: Human quality annotation of GPT-4o outputs with automated retrieval (raw counts out of 50) by prompt type.

## G   Full Error Counts

This section contains the full counts of annotated errors by category and prompt type.

**GPT-3.5 Turbo**

|  | BASE | C | G | M | C+G | C+M | G+M | C+G+M | TOTAL |
|---|---|---|---|---|---|---|---|---|---|
| None | 0 | 1 | 1 | 6 | 0 | 8 | 3 | 5 | 24 |
| Addition | 0 | 5 | 3 | 14 | 1 | 9 | 10 | 11 | 53 |
| Omission | 3 | 9 | 2 | 13 | 2 | 5 | 9 | 9 | 52 |
| Substitution - Subject | 0 | 3 | 0 | 7 | 0 | 9 | 9 | 12 | 40 |
| Substitution - TAM | 0 | 11 | 3 | 17 | 6 | 19 | 19 | 19 | 94 |
| Substitution - Other | 4 | 9 | 4 | 13 | 6 | 16 | 14 | 13 | 79 |
| Overtranslation | 1 | 1 | 1 | 4 | 0 | 2 | 3 | 2 | 14 |
| Undertranslation | 0 | 0 | 0 | 2 | 1 | 2 | 2 | 2 | 9 |
| Target Error - Grammar | 0 | 1 | 1 | 4 | 2 | 3 | 3 | 1 | 15 |
| Target Error - Coherence | 0 | 0 | 3 | 5 | 2 | 3 | 7 | 7 | 27 |
| Target Error - Style/Register | 0 | 3 | 0 | 5 | 2 | 3 | 1 | 2 | 16 |
| Complete Mistranslation | 30 | 19 | 21 | 2 | 18 | 2 | 2 | 2 | 96 |
| Mistranslation - Lexical Correspondence | 15 | 13 | 23 | 9 | 21 | 11 | 15 | 10 | 117 |
| Refusal | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| Total | 54 | 75 | 62 | 101 | 61 | 92 | 97 | 95 | 637 |

Table 11: Human error type annotation of GPT-3.5 outputs with automated retrieval (raw counts, up to 3 errors per sentence) by prompt type.

**GPT-4o**

| | BASE | C | G | M | C+G | C+M | G+M | C+G+M | TOTAL |
|---|---|---|---|---|---|---|---|---|---|
| None | 15 | 16 | 10 | 16 | 13 | 19 | 14 | 18 | 121 |
| Addition | 2 | 5 | 7 | 5 | 4 | 1 | 6 | 4 | 34 |
| Omission | 8 | 7 | 6 | 7 | 6 | 3 | 5 | 5 | 47 |
| Substitution - Subject | 1 | 2 | 0 | 1 | 2 | 1 | 2 | 2 | 11 |
| Substitution - Other | 22 | 24 | 22 | 18 | 19 | 18 | 17 | 20 | 160 |
| Substitution - TAM | 16 | 17 | 19 | 12 | 13 | 10 | 11 | 9 | 107 |
| Overtranslation | 2 | 1 | 0 | 2 | 2 | 2 | 1 | 2 | 12 |
| Undertranslation | 6 | 1 | 3 | 1 | 3 | 0 | 1 | 2 | 17 |
| Target Error - Grammar | 1 | 3 | 4 | 4 | 1 | 2 | 6 | 1 | 22 |
| Target Error - Coherence | 1 | 3 | 4 | 5 | 4 | 5 | 9 | 5 | 36 |
| Target Error - Style/Register | 1 | 2 | 3 | 4 | 4 | 2 | 4 | 3 | 23 |
| Complete Mistranslation | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| Mistranslation - Lexical Correspondence | 4 | 3 | 5 | 6 | 6 | 6 | 9 | 4 | 43 |
| Total | 79 | 85 | 83 | 81 | 77 | 69 | 85 | 75 | 634 |

Table 12: Human error type annotation of GPT-4o outputs with automated retrieval (raw counts, up to 3 errors per sentence) by prompt type.

# Introducing the NewsPaLM MBR and QE Dataset: LLM-Generated High-Quality Parallel Data Outperforms Traditional Web-Crawled Data

**Mara Finkelstein     David Vilar     Markus Freitag**
Google
{marafin,vilar,freitag}@google.com

## Abstract

Recent research in neural machine translation (NMT) has shown that training on high-quality machine-generated data can outperform training on human-generated data. This work accompanies the first-ever release of a LLM-generated, MBR-decoded and QE-reranked dataset with both sentence-level and multi-sentence examples.[1] We perform extensive experiments to demonstrate the quality of our dataset in terms of its downstream impact on NMT model performance. We find that training from scratch on our (machine-generated) dataset outperforms training on the (web-crawled) WMT'23 training dataset (which is 300 times larger), and also outperforms training on the top-quality subset of the WMT'23 training dataset. We also find that performing self-distillation by finetuning the LLM which generated this dataset outperforms the LLM's strong few-shot baseline. These findings corroborate the quality of our dataset, and demonstrate the value of high-quality machine-generated data in improving performance of NMT models.

## 1   Introduction

With the advent of large language models (LLMs), machine translation (MT) quality has improved dramatically (Kocmi et al., 2023a, 2024a), and performance tends to scale with model size (Gemini Team, 2024). While LLMs are now state-of-the-art translators, they are often impractical to use or serve, especially in high-traffic and/or resource-constrained settings. Thus, development of smaller, but still highly performant, MT models remains an active area of research. Recent work has shown that distillation of LLM translation quality, while requiring an expensive data generation process, is an effective approach (Li et al., 2024). In this work,

we introduce a new LLM-generated dataset called NewsPaLM, which we make freely available.

In addition to the size of the teacher model, another key determinant of the quality of machine-generated translation data is the decoding method used. While beam search and greedy decoding are the most common decoding methods used for NMT, Eikema and Aziz (2020a) showed that maximum *a posteriori* (MAP) decoding methods are suboptimal, and instead proposed Minimum Bayes Risk (MBR) decoding. Unlike MAP decoding, MBR decoding does not aim to produce the translation with the highest estimated model probability. Instead, it chooses the translation that is estimated to have the highest quality with respect to a utility metric. A follow-up study by Freitag et al. (2022) showed that MBR decoding with neural utility metrics significantly outperforms beam search decoding, according to expert-based human evaluation.

The main drawback of MBR decoding is its high computational cost. In particular, the algorithm requires that, for every input query, a large number $n$ of candidates be generated from the model, and then an (expensive) scoring function be computed on every pair of distinct candidates $(n_i, n_j)$, for a total of $O(n^2)$ computations. QE reranking (Fernandes et al., 2022) is a more efficient alternative to MBR decoding. This decoding method instead reranks the candidate model predictions using a neural quality estimation (QE) metric, and requires only $O(n)$ computations.

Finkelstein and Freitag (2023) showed that finetuning NMT models on MBR-decoded and QE-reranked datasets is an effective technique for distillation (while finetuning on beam search-decoded datasets is not) and that, given a LLM teacher, MBR and QE distillation can outperform finetuning on human-generated references.

In this work, we generate sentence-level parallel data using MBR decoding and multi-sentence parallel data using QE reranking. In addition to detailing

---

[1]The dataset can be found at https://github.com/google-research/google-research/tree/master/newspalm_mbr_qe.

our dataset creation method, we also perform extensive experiments to demonstrate the quality of our dataset in terms of its downstream impact on NMT model performance.

Our contributions can be summarized as follows:

- We release our LLM-generated, sentence-level and multi-sentence, MBR and QE translation dataset.
- We demonstrate that our dataset is high-quality, by using it to train NMT models from scratch and comparing performance against baselines using human-generated parallel data. This is the first work to pretrain NMT models on MBR and QE data.
- We show that training on our dataset outperforms training on the web-crawled WMT'23 training dataset (which is 300 times larger than ours). Moreover, our dataset also outperforms (by an even larger margin) when compared against quality-based filtering of the WMT'23 dataset to match the size of our dataset.
- We also demonstrate our dataset's quality by performing self-distillation (using the PaLM-2 LLM from which this data was generated), and show that this outperforms the LLM's strong few-shot baseline. To our knowledge, this is the first work to investigate MBR finetuning a LLM.
- We investigate the effect of sentence-level versus multi-sentence MBR and QE training data on NMT model performance as a function of sequence length, and investigate the tradeoff between dataset size and model quality, during both pretraining and finetuning.

## 2  NewsPaLM Dataset

This paper accompanies a dataset release of sentence-level and multi-sentence English-German and German-English parallel data, generated from the (monolingual) Newscrawl corpus as made available for the WMT evaluation campaigns[2] using the *PaLM-2 Bison* LLM (Anil et al., 2023). We detail below the steps to create this dataset, which we call NewsPaLM.

The dataset construction process consisted of four steps, as described in the following sections.

### 2.1  Source-side Data Collection: Newscrawl

To construct the English and German source-side datasets, we first collected all Newscrawl data from 2007 to 2022, released as part of the WMT'23 Machine Translation Shared Task (Kocmi et al., 2023a). This is a large corpus of crawled news, with about 398 million and 507 million lines for English and German, respectively. For both of these languages, document-split versions of the dataset (with document boundaries intact) are available.

We collected both the sentence-level and document-level versions of the datasets. Basic preprocessing had already been applied to the sentence-level version, including removing lines with no ASCII letters and deduplication. This preprocessing was not applied to the document-level version. We performed minimal additional cleaning to fix incorrectly encoded characters.

### 2.2  Construction of "Blobs"

We used the document-split versions of the datasets to construct multi-sentence (i.e. "blob-level") examples. We refer to these examples as blobs, rather than paragraphs, since they do not respect paragraph boundaries but, rather, simply represent the concatenation of contiguous sentences up to a maximum length. In particular, we joined headlines using the separator "\n\n", and otherwise joined sentences with spaces, up to a maximum length of 512 tokens (using the PaLM-2 tokenizer; Anil et al. (2023)). The blobs respect document boundaries, each blob contains only complete sentences (no sentence fragments), and each blob may or may not contain a headline (depending on where in the document the blob comes from).

### 2.3  Cluster-Based Text Selection

The size of the Newscrawl full dataset and the high computational cost of the decoding techniques (§2.4) makes it impractical to process all the available data. In order to reduce the size of the data, while at the same time ensuring diversity in the samples, we follow a clustering-based sample selection approach. As a first step, we embed the source side of the data using XLM-RoBERTa (Conneau et al., 2020). We then apply Recursive Agglomerative Clustering (RAC) (Sumengen et al., 2021), an unsupervised clustering algorithm which is an efficient extension of Hierarchical Agglomerative Clustering. These algorithms are initialized by defining a set of clusters, each containing a single point from

|              | EN → DE | DE → EN |
|--------------|---------|---------|
| Sentence-level | 3,287 | 3,264 |
| Blob-level     | 3,826 | 4,017 |

Table 1: Number of defined clusters per dataset.

the original data points. The guiding principle is to iteratively merge the two clusters which are closest to each other, until some stopping criterion is met, e.g. a maximum distance between the clusters to be merged. Note that this algorithm requires the number of clusters to be chosen as a hyperparameter, unlike other clustering algorithms like $k$-means which have the advantage that the number of clusters is defined by the algorithm itself. We selected the number of clusters shown in Table 1 for each of the data sets.

Once the clusters have been defined, we sample uniformly from them. In this way, we ensure that the diversity of the original dataset is maintained in the reduced sample.

## 2.4 MBR Decoding and QE Reranking

The preceding steps handle preparation of source-side data. To generate the target-side data from these sources, we used the *PaLM-2 Bison* LLM (Anil et al., 2023), 5-shot prompted with ICL examples from the newstest2021 test set (Akhbardeh et al., 2021). Note that unlike previous work which also used PaLM-2 to generate translation data for distillation (Finkelstein and Freitag, 2023), here we do not finetune on the translation task prior to data generation.

A key component of our data generation process is the decoding method. We generated the sentence-level data using MBR decoding and the blob-level data using QE reranking. Both MBR decoding and QE reranking can be decomposed into two steps: candidate list generation (Section § 2.4.1) and scoring (Section § 2.4.2).

### 2.4.1 Candidate List Generation

The first step in the decoding process is to generate a list of candidate model outputs, given a source segment. In this work, we used a candidate size of 512 and generated candidate translations using epsilon sampling (Hewitt et al., 2022) with $\varepsilon = 0.02$, which was shown to be the best sampling method for MBR decoding in Freitag et al. (2023).

### 2.4.2 MBR and QE scoring

Next, the best output is chosen based on a utility function. This step is where MBR decoding and QE reranking diverge. For MBR decoding, we use a reference-based utility metric $u_{mbr}(h, r)$, which estimates the quality of a candidate translation $h$ conditioned on a reference translation $r$. Formally, given a set of hypotheses $\mathcal{H}$, the Minimum Bayes Risk (MBR) translation $h^{mbr}$ is selected from the candidates in $\mathcal{H}$ according to

$$h^{mbr} = \arg\max_{h \in \mathcal{H}} \frac{1}{|\mathcal{H}|} \sum_{y \in \mathcal{H}} u_{mbr}(h, y).$$

For QE reranking, on the other hand, we use a reference-free (QE) utility metric $u_{qe}(h, s)$, which estimates the quality of a candidate translation $h$ conditioned on the source $s$, rather than on the reference. We select the best QE translation $h^{qe}$ of the source $s$ from the candidates in $\mathcal{H}$ as

$$h^{qe} = \arg\max_{h \in \mathcal{H}} u_{qe}(h, s)$$

In this work, we used the *BLEURT* (Sellam et al., 2020) utility metric for MBR decoding and the *MetricX-QE* (Juraska et al., 2023) utility metric for QE reranking. Note that the maximum context length for *BLEURT* (candidate and reference combined) is 512, while for *MetricX-QE* (candidate and source combined), it is 1024. Given that the blob-level source-side data alone can contain up to 512 tokens, we could not use *BLEURT* as the utility function for this data. MBR decoding with *MetricX* is prohibitively expensive, hence our decision to perform QE reranking instead.

As a baseline against which to compare data generated using these state-of-the-art decoding methods, we also created accompanying sentence-level and blob-level datasets from the same source-side data using greedy decoding.

## 2.5 Dataset Statistics

Here we briefly present basic statistics about the four datasets we created (sentence-level and blob-level versions, for the en→de and de→en language pairs). See Appendix A for additional dataset statistics.

Table 2 shows the size (in number of examples) of each dataset. Note that each dataset has about 800 thousand to one million examples.

Table 3 shows the average length of source and target examples (in number of tokens, as defined by the Moses tokenizer) per dataset. For

|  | **EN → DE** | **DE → EN** |
|---|---|---|
| MBR SENT-LEVEL | 998,435 | 1,022,344 |
| QE BLOB-LEVEL | 925,829 | 769,028 |

Table 2: Number of examples per dataset.

|  |  | **Source** | **Target** |
|---|---|---|---|
| EN → DE | SENTENCES | 37.5 | 39.8 |
|  | BLOBS | 364.5 | 339.8 |
| DE → EN | SENTENCES | 77.3 | 88.3 |
|  | BLOBS | 288.4 | 323.4 |

Table 3: Average source and target lengths per dataset, computed using the Moses tokenizer.

English-German, the blob-level examples are about ten times longer than the sentence-level examples, while for German-English, they are about four times longer. Figure 1 shows the distribution of target example lengths for English-German. Note that the blob-level data distribution is shifted to the right of the sentence-level data distribution, as expected.

## 3 Experimental Setup

We perform a series of pretraining and finetuning experiments to validate the quality of our NewsPaLM dataset, and to contextualize its performance with respect to a much larger dataset of human-generated data. All of our experiments are performed on both English-German (en→de) and German-English (de→en).



Figure 1: Distribution of English-German MBR sentence-level versus QE blob-level target lengths (computed using the Moses tokenizer).

### 3.1 Datasets

#### 3.1.1 Training Data

As a baseline against which to compare our NewsPaLM dataset, we use the parallel WMT'23 training data (Kocmi et al., 2023b), which consists of 296 million sentence-level examples. A subset of this data (consisting of about 3 million sentences, from Europarl, News Commentary, and Rapid documents) contains document boundaries, which we use to construct blob-level examples using a procedure similar to the blob-level dataset creation process described in Section §2.2. That is, we partition the sentences into contiguous blocks, each of which has a total number of tokens up to a token limit of 512 (for each of source and target). In our experiments, this WMT'23 data is only used for pretraining.

The remainder of our pretraining and finetuning data comes from our (machine-generated) NewsPaLM dataset, described in Section §2. As an additional baseline, we compare the MBR-decoded and QE-reranked versions of this dataset against the greedy-decoded version. Note that for both language pairs (en→de and de→en), the sentence-level and blob-level NewsPaLM data combined contains less than 2 million examples (Table 2).

#### 3.1.2 Development and Test Sets

For both language pairs, we use the sentence-level and paragraph-level versions of the newstest2021 test set (Farhad et al., 2021), as well as the (sentence-level) generalMT2022 test set (Kocmi et al., 2022), as our development sets for checkpoint picking. We report all results on the WMT'23 (Kocmi et al., 2023b) and WMT'24 (Kocmi et al., 2024b) test sets. Note that the WMT'23 and WMT'24 en→de test sets are paragraph-level.

### 3.2 Models

For both language pairs (en→de and de→en), we use a 602 million parameter Transformer encoder-decoder architecture, implemented in *Pax*[3]. The model has 8 encoder and 8 decoder layers (rather than 6), but otherwise is similar to the *transformer-big* setting in Vaswani et al. (2017), with model dimension of 1024, hidden dimension of 8192, and 16 multi-attention heads. We train without label smoothing. For each language pair, we use a bilingual vocabulary of 32k subword units trained on the WMT'23 training dataset (Kocmi et al., 2023b).

---

[3] https://github.com/google/paxml

The best (base and incremental) checkpoints were chosen to maximize *BLEURT* (Sellam et al., 2020) on the development set.

We also experiment with self-distillation of the *PaLM-2 Bison* (Anil et al., 2023) LLM, which is the model used to generate our datasets (see Section §2.4). We compare self-distillation (finetuning) against 5-shot prompting of this model (using the same ICL examples as during NewsPaLM dataset generation).

### 3.3 Evaluation

We evaluate our models on four automatic metrics: *MetricX* (Juraska et al., 2023), *Comet20* (Rei et al., 2020), *Comet22* (Rei et al., 2022), and *BLEURT* (Sellam et al., 2020). Note that for *MetricX*, lower scores are better, while for the remaining metrics, higher scores are better. Since the MBR data is generated using *BLEURT* as the utility function, and the QE data is generated using *MetricX*, the MBR-finetuned models may overfit to the *BLEURT* metric, while the QE-finetuned models may overfit to the *MetricX* metric. Thus, we primarily depend on the *Comet\** metrics to measure model quality.

## 4 Results

### 4.1 Pretraining

We first experiment with training bilingual (en→de and de→en) encoder-decoder translation models (as described in §3.2) from scratch, to compare our NewsPaLM dataset (as described in §3.1.1) against the WMT'23 training dataset. As shown in Table 4, **pretraining on the NewsPaLM QE blob-level dataset** (which contains less than one million examples; Table 2) **outperforms pretraining on the entire WMT'23 training dataset**, which is more than 300 times larger. The NewsPaLM QE dataset achieves a *Comet22* score of 80.62 on the English-German WMT'23 test set (row 2c), while the WMT'23 training dataset achieves a score of 78.79 (row 1a).

Note that **training on the MBR sentence-level data** (row 2b) **underperforms training on the QE blob-level data** (row 2c). As shown in Figure 2, this is mostly due to a large drop in performance on longer sequence lengths. Thus, **exposure to multi-sentence data during training is essential to perform well on paragraph-level test sets**. Also note that during pretraining, we see no additional gains from mixing in the MBR sentence-level data relative to using the QE blob-level data

only (rows 2c versus 2d).

We also experiment with pretraining on the greedy-decoded version of our NewsPaLM dataset, to compare against pretraining on the MBR-decoded and QE-reranked versions. Interestingly, the former (pretraining on the greedy-decoded data) outperforms the latter (pretraining on the MBR-decoded and QE-reranked versions), as shown in rows 2a versus 2d in Table 4. Based on manual inspection of examples, we hypothesize that the MBR-decoded and QE-reranked data is more free-style and harder for the model to learn than the greedy-decoded data. This is illustrated in Table 10 in the Appendix. If this were the case, the model would perform better by first learning the "easier" data (during pretraining), then adapting to the more free-style data during finetuning. We test this hypothesis by comparing two model training curricula: For the first, we pretrain on the greedy-decoded data and finetune on the MBR-decoded and QE-reranked data. For the second, we do the opposite: pretraining on the MBR-decoded and QE-reranked data and finetuning on the greedy decoded data. As we hypothesized, the **former model training curriculum (MBR and QE finetuning from the greedy-pretrained checkpoint) performed better** (Table 5).

We have seen that pretraining on a small and clean, synthetically-produced dataset (NewsPaLM) can outperform finetuning on a large and noisy, human-generated one (WMT'23 training dataset). However, previous work such as Peter et al. (2023) has shown that MT model performance can be boosted by selecting a high-quality subset of a large and noisy training corpus, using data selection techniques such as QE filtering. Thus, we perform QE filtering (using the *BLEURT-QE* metric, as in Peter et al. (2023)) to select the highest-quality examples in the WMT'23 (sentence-level) training dataset, while reducing its size to exactly match that of our (sentence-level) NewsPaLM dataset (of about one million examples). As shown in row 1b versus row 2b in Table 4, **training on the QE-filtered WMT'23 dataset substantially underperforms training on our MBR-decoded NewsPaLM dataset (of the same size)**, and also underperforms training on the full WMT'23 dataset. Note that this result does not contradict previous work showing the benefit of data filtering, since previous work did not reduce the dataset to such a small fraction (0.3%) of the original size. Thus, our

Figure 2: Comparison of pretraining performance on `NewsPaLM` MBR sentence-level dataset versus `NewsPaLM` QE blob-level dataset, bucketed by source length (WMT'23 en→de test set). Note that performance of the model trained on the blob-level data is stable across segment lengths, while performance of the model trained on the sentence-level data declines as segment length increases (according to both *MetricX* and *Comet22* metrics).

`NewsPaLM` dataset is highly efficient (which is one indicator of its quality), and its efficiency cannot be matched be selecting a high-quality subset of a large, noisy corpus.

## 4.2 Finetuning

Next, we experiment with how the different variants of our `NewsPaLM` dataset (and mixtures thereof) behave during finetuning (and whether this behavior differs from that observed during pretraining). Unless otherwise indicated, we initialize from the checkpoint pretrained on the WMT'23 training data (row 1a in Table 4). We report en→de and de→en results on the WMT'23 test set in Table 6, and refer the reader to Table 12 in Appendix B for pretraining and finetuning results on the WMT'24 en→de test set.

As shown in Table 6, **MBR and QE finetuning** (row 1d) **outperforms greedy finetuning** (row 1a), using the same mixture proportions (9:1) for the sentence-level and blob-level data. As shown in Table 5 and discussed in §4.1, MBR and QE finetuning from the greedy-pretrained checkpoint outperforms greedy finetuning from the MBR and QE-pretrained checkpoint as well. Also, note that for en→de, MBR and QE finetuning from the checkpoint pretrained on the WMT'23 training data (row 1d in Table 6) slightly underperforms initializing from the checkpoint pretrained on the greedy `NewsPaLM` dataset (row 1a in Table 5) according to the WMT'23 test set, but the opposite is the case according to the WMT'24 test set (see Tables 12 and 13 in Appendix B) and based on the de→en results on the WMT'23 test set.

Unlike during pretraining, **finetuning on the MBR sentence-level data outperforms finetuning on the QE blob-level data** (rows 1b versus 1c in Table 6), and we see no additional gains from mixing in the QE blob-level data relative to using the MBR sentence-level data only (rows 1b versus 1d). We hypothesize that the model learns to use long context (from the blob-level data) during pretraining, and it doesn't forget during finetuning, so blob-level data is less important during this stage.

We also experiment with finetuning the *PaLM-2 Bison* (Anil et al., 2023) LLM (as described in §3.2), which is the teacher model used to generate our `NewsPaLM` dataset. As shown in Table 6, **self-distillation via MBR (and QE) finetuning does indeed improve performance over the LLM's strong few-shot baseline** (rows 2a vs 2b). As with the encoder-decoder model, finetuning *PaLM-2 Bison* on the MBR data outperforms finetuning on the QE data (and outperforms finetuning on a mixture of the MBR and QE data). The improvement in performance of *PaLM-2 Bison* due to MBR finetuning is observed across all source length buckets (Figure 4 in Appendix B) and all domains in the WMT'23 and WMT'24 test sets (Tables 14 and 15 in Appendix B), despite the MBR data being sentence-level only and coming primarily from the news domain. MBR finetuning the *PaLM-2 Bison* model also outperforms MBR finetuning the much smaller encoder-decoder student (rows 1b vs 2b in Table 6), as expected.

## 4.3 Ablations

### 4.3.1 Effect of Dataset Size

Given the expense of creating LLM-generated, MBR-decoded datasets such as the ones presented in this work, we investigate how model performance scales with dataset size during both pretraining and finetuning. We randomly sample 25% of the MBR-decoded `NewsPaLM` dataset (for both en→de and de→en), then train on the subsampled dataset. As shown in Figure 3 (and Table 16 in Appendix B), *finetuning* **on the subsampled dataset only took a small performance hit relative to finetuning on the full dataset, but** *pretraining* **took a large performance hit**. Thus, it is likely that pretraining performance would continue to improve had we generated a larger `NewsPaLM` dataset, while we would be unlikely to observe substantial incremental improvements in finetuning performance by increasing the dataset size. Also note

| | Model | MetricX ↓ | COMET22 ↑ |
|---|---|---|---|
| | 1a) WMT'23 (all) | 4.20 | 78.79 |
| | 1b) WMT'23 (sentence-level, *BLEURT-QE* filtered) | 16.69 | 43.18 |
| *en→de* | 2a) Greedy sentence-level + blob-level (9:1) | **2.60** | **81.67** |
| | 2b) MBR sentence-level | 6.39 | 72.05 |
| | 2c) QE blob-level | 2.82 | 80.62 |
| | 2d) MBR sentence-level + QE blob-level (9:1) | 2.99 | 79.68 |
| | 1a) WMT'23 (all) | 5.55 | 82.41 |
| | 1b) WMT'23 (sentence-level, *BLEURT-QE* filtered) | 14.80 | 57.27 |
| *de→en* | 2a) Greedy sentence-level + blob-level (9:1) | **3.47** | **83.30** |
| | 2b) MBR sentence-level | 4.97 | 80.55 |
| | 2c) QE blob-level | 4.01 | 82.33 |
| | 2d) MBR sentence-level + QE blob-level (9:1) | 3.95 | 82.02 |

Table 4: Pretraining performance (WMT'23 test set).

| | Model | MetricX ↓ | COMET22 ↑ |
|---|---|---|---|
| *en→de* | 1a) MBR + QE finetuning (from greedy-pretrained ckpt) | **2.11** | **82.78** |
| | 1b) Greedy finetuning (from MBR + QE-pretrained ckpt) | 2.63 | 81.48 |
| *de→en* | 1a) MBR + QE finetuning (from greedy-pretrained ckpt) | **3.10** | **84.05** |
| | 1b) Greedy finetuning (from MBR + QE-pretrained ckpt) | 3.60 | 83.08 |

Table 5: Comparison of pretraining on NewsPaLM greedy data, then finetuning on NewsPaLM MBR and QE data, versus vice-versa (WMT'23 test set).

that the stability in finetuning performance under subsampling held up despite using the most efficient subset selection method (random, as opposed to e.g., QE filtering), another indicator supporting the high quality of our NewsPaLM dataset.

### 4.3.2 Effect of Cluster-based Data Selection

As described in §2.3, we used a clustering-based approach (sampling uniformly over the computed clusters) to select the subset of Newscrawl data which we used to generate the NewsPaLM dataset. To isolate the effect of our sample selection technique, we compare its performance against sampling uniformly from the original Newcrawl dataset distribution (i.e., without taking cluster information into account). Since our NewsPaLM dataset contains the subset of Newscrawl examples selected by sampling uniformly over the clusters, we approximate the above comparison as follows, selecting 25% of the NewsPaLM dataset in both cases:

- We sample uniformly from NewsPaLM to approximate cluster-guided sampling from the



Figure 3: Comparison of model performance when pretraining and finetuning on the full versus subsampled NewsPaLM MBR dataset (WMT'23 test set). The subsampled dataset is 25% of the size of the full dataset, and was sampled randomly. Note that pretraining performance drops substantially when training on the subsampled dataset (for both en→de and de→en), while finetuning performance is minimally affected.

| Model | MetricX ↓ | COMET22 ↑ |
|---|---|---|
| **en→de** | | |
| 1a) Greedy sentence-level + blob-level (9:1) | 2.59 | 81.49 |
| 1b) MBR sentence-level | 2.30 | **82.69** |
| 1c) QE blob-level | 2.45 | 81.83 |
| 1d) MBR sentence-level + QE blob-level (9:1) | **2.26** | 82.52 |
| 2a) PaLM-2 five-shot (no finetuning) | 1.62 | 84.54 |
| 2b) PaLM-2 MBR sentence-level | **1.14** | **85.64** |
| 2c) PaLM-2 QE blob-level | 1.47 | 84.77 |
| 2d) PaLM-2 MBR sentence-level + QE blob-level (9:1) | 1.17 | 85.54 |
| **de→en** | | |
| 1a) Greedy sentence-level + blob-level (9:1) | 3.12 | 84.14 |
| 1b) MBR sentence-level | 2.91 | **84.57** |
| 1c) QE blob-level | 2.99 | 84.27 |
| 1d) MBR sentence-level + QE blob-level (9:1) | **2.82** | 84.53 |
| 2a) PaLM-2 five-shot (no finetuning) | 2.25 | 85.36 |
| 2b) PaLM-2 MBR sentence-level | **1.91** | **86.26** |
| 2c) PaLM-2 QE blob-level | 2.03 | 85.81 |
| 2d) PaLM-2 MBR sentence-level + QE blob-level (9:1) | 1.92 | 86.24 |

Table 6: Finetuning performance (WMT'23 test set). Unless otherwise indicated, performance is reported for the encoder-decoder model. For finetuning, this model was initialized from the checkpoint pretrained on the full WMT'23 training dataset (row 1a in Table 4). Results for *PaLM-2 Bison* few-shot prompting versus self-distillation using NewsPaLM MBR and QE data are reported in rows 2a-d.

full Newscrawl dataset.

- We use the original cluster sizes of the full Newscrawl dataset (computed prior to selecting the NewsPaLM subset), and sample from NewsPaLM according to this distribution. This approximates sampling from the full Newscrawl dataset *without* taking cluster information into account. Note that because the original cluster distribution was highly skewed, with most of the examples belonging to the top few clusters, we could not exactly match the original distribution while sampling 25% of the NewsPaLM dataset, but we chose the distribution to be the one which was closest to the original.

As shown in Table 7, using the cluster information in the subsampling procedure marginally improves results for en→de pretraining and finetuning, and for de→ en finetuning. (There is no clear signal for de→en pretraining; according to *MetricX*, using the cluster information helps, while according to *Comet22*, it hurts.)

## 5 Discussion

Training on LLM-generated, MBR-decoded and QE-reranked datasets is an established technique for leveraging monolingual data to improve NMT model quality (Finkelstein and Freitag (2023), Wang et al. (2024)). While this technique is highly effective, generating such datasets remains a substantial bottleneck, and is often prohibitively expensive. This work accompanies the first-ever open-source release of a LLM-generated, sentence-level and blob-level MBR and QE dataset. We measure the quality of our dataset in terms of its downstream impact on NMT model performance, both when training a NMT model from scratch and when fine-tuning.

We find that training from scratch on our MBR-decoded and QE-reranked NewsPaLM dataset outperforms training on the entire WMT'23 training dataset (which is 300 times larger), and also outperforms training on the top-quality subset of the WMT'23 training data (selected via QE filtering, and matching the size of our dataset). Moreover, we find that NMT models are unable to generalize well to multi-sentence queries without exposure to such data at training time, motivating the inclusion of blob-level data in our dataset.

We also find that MBR and QE finetuning outperform finetuning on the greedy-decoded version of our dataset. Unlike Finkelstein and Freitag (2023),

| | Model | MetricX ↓ | COMET22 ↑ |
|---|---|---|---|
| | 1a) PT: Sampling uniformly over clusters | **11.02** | **59.75** |
| | 1b) PT: Sampling uniformly from original Newscrawl distribution | 11.60 | 59.09 |
| *en→de* | | | |
| | 2a) FT: Sampling uniformly over clusters | **2.55** | **82.11** |
| | 2b) FT: Sampling uniformly from original Newscrawl distribution | **2.55** | 81.88 |
| | 1a) PT: Sampling uniformly over clusters | **7.41** | 75.91 |
| | 1b) PT: Sampling uniformly from original Newscrawl distribution | 7.58 | **76.55** |
| *de→en* | | | |
| | 2a) FT: Sampling uniformly over clusters | **2.96** | **84.33** |
| | 2b) FT: Sampling uniformly from original Newscrawl distribution | 3.09 | 83.92 |

Table 7: Comparison of model performance when trained on subsampled NewsPaLM data with and without cluster-based data selection (WMT'23 test set). PT stands for Pretraining and FT, for Finetuning. Random subsampling of NewsPaLM approximates sampling uniformly across the Newscrawl clusters, while subsampling according to the Newscrawl cluster distribution approximates discarding cluster information and sampling randomly according to the original data distribution.

which only performs MBR and QE self-distillation using a small encoder-decoder NMT model, here we show that self-MBR and self-QE finetuning are effective for the much stronger *PaLM-2 Bison* LLM as well.

Finally, we show via subsampling experiments on our NewsPaLM dataset that pretraining versus finetuning performance scale very differently with dataset size: While finetuning performance only took a small hit when reducing our dataset to 25% of its original size, pretraining performance took a large hit. However, note that the full NewsPaLM dataset is already orders of magnitude smaller than most datasets used for NMT model training, including the WMT'23 training dataset.

## 6 Related Work

While MT research has traditionally relied on MAP decoding or generating $k$-best lists through beam search for MBR decoding, Eikema and Aziz (2020b) proposed an approximation of MBR decoding via unbiased sampling. Their method aims to address the limitations of MAP decoding (Eikema and Aziz, 2020b; Müller and Sennrich, 2021; Eikema and Aziz, 2022) by demonstrating that samples drawn from the NMT model align more faithfully with training data statistics when compared to beam search. Freitag et al. (2022) showed that using neural metrics results in significant improvements in translation quality. To the best of our knowledge, this is the first work that applies MBR decoding beyond the sentence level for the task of machine translation.

While the improvements in translation quality afforded by MBR are widely acknowledged, its high computational cost limits its application in practice. Different approaches have been proposed to speed up MBR computation, e.g. (Eikema and Aziz, 2022; Cheng and Vlachos, 2023; Jinnai and Ariu, 2024; Vamvas and Sennrich, 2024; Tomani et al., 2024). Similar in spirit to MBR decoding, QE-rescoring approaches (Fernandes et al., 2022) also directly optimize a utility function, with linear-time cost.

We approach the efficiency problem from a different perspective, carrying out a one-off expensive MBR decoding run, which can then be re-used for training and finetuning other models via knowledge distillation (Buciluǎ et al., 2006; Hinton et al., 2015). This technique has been a successful way to improve smaller systems by leveraging the capacities of bigger models, while retaining higher computational efficiency. The technique has been applied to numerous NLP tasks, including neural machine translation (Kim and Rush, 2016; Tan et al., 2018; Zhang et al., 2019; Jooste et al., 2022; Wan et al., 2024, inter alia). In the current era of LLMs, these models provide prime candidates for leveraging their impressive capabilities to improve other models. Yoo et al. (2021) use GPT3 for data augmentation in different classification tasks, in addition to using soft-labels predicted by the language model. Hsieh et al. (2023) propose to use "rationales" generated by PaLM to train a much smaller T5 model, achieving comparable or even superior performance. Li et al. (2024) use "selective" distillation to generate synthetic data using a variant of LLaMA-7B, expanding the coverage of training data for a translation model. Closer to our work, Finkelstein and Freitag (2023) propose to use MBR on a LLM to generate high-quality translations with which to train a dedicated translation

model. As reference, (Xu et al., 2024) provides a much more comprehensive survey of knowledge distillation approaches using LLMs.

Another dimension related to our work is the area of data selection for NMT training. While a big amount of work has been focused on filtering noise from web-crawled data (e.g. Zaragoza-Bernabeu et al., 2022), there are also approaches aimed at improving the translation quality by limiting the training data to high-quality samples. Carpuat et al. (2017) use semantic divergence to select the most relevant portion of the training data, while (Peter et al., 2023) use QE metrics on the training data to select only high-quality sentence pairs. Xu et al. (2023) indeed show that only a small amount of high-quality multilingual and parallel data is needed for obtaining state-of-the art translation results finetuning a LLM. A similar approach was used by (Alves et al., 2024) to finetune LLaMA for translation and translation-related tasks.

One can also find different examples of clustering for data selection for NLP tasks. Aharoni and Goldberg (2020) showed that automatic clustering techniques can adequately recover semantic information from text corpora. Yu et al. (2023) use clustering for data selection to finetuning a LLM. Related to these approaches, nearest-neighbor machine translation (Khandelwal et al., 2021) uses distance measures between examples to select examples closer to the sentence to translate in an additional module of a translation system. (Agrawal et al., 2023) and (Vilar et al., 2023) use similar approaches to construct prompts for LLMs.

## 7 Conclusion

In this work, we have described the dataset creation process for the first-ever release of a LLM-generated MBR and QE dataset. We have shown that this dataset can be used to build a small and efficient, but high-quality, NMT model from scratch. In fact, training on this dataset outperforms training on the much larger, human-generated WMT'23 dataset. We have also shown that this dataset can improve NMT performance during finetuning, both for an encoder-decoder system and via self-distillation for an already highly performant LLM. We hope that this dataset will enable further distillation research by the wider community, even by those without resources to generate datasets from large teacher models using expensive decoding techniques.

There are many avenues for future work. This work presented the first investigation of multi-sentence (i.e., blob-level) QE finetuning, and a natural next step would be to move to the document level. The dataset creation process described here could also be continued iteratively, by generating a new MBR and QE dataset from the same LLM teacher, but after finetuning on the original version of the dataset (or the version from the previous iteration). While this would be expensive, it would likely yield further incremental improvements in dataset quality. Finally, there remain many open questions regarding how to optimally perform distillation of a stronger teacher model into a weaker student using MBR and QE data. For instance, rather than finetuning on a uniform mixture of all examples in the dataset, the student model's perplexity on these examples could be taken into account to select a subset of examples and/or to determine the optimal progression of examples to expose the student to during finetuning.

## Limitations

The (target-side) data generation process was expensive, due to both using a LLM and a costly decoding method. For MBR dataset creation, computation of each dataset example required generation of $n$ outputs from the LLM teacher model, and then $O(n^2)$ forward passes through the utility function, where $n$ is the candidate size. For QE dataset generation, $O(n)$ forward passes through the utility function were required per example. Thus, the dataset construction method proposed here is not easily scalable to other language pairs and/or source-side data in the absence of substantial computing resources.

## References

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. In-context examples selection for machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8857–8873, Toronto, Canada. Association for Computational Linguistics.

Roee Aharoni and Yoav Goldberg. 2020. Unsupervised domain clusters in pretrained language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7747–7763, Online. Association for Computational Linguistics.

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88, Online. Association for Computational Linguistics.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. Palm 2 technical report.

Cristian Buciluǎ, Rich Caruana, and Alexandru Niculescu-Mizil. 2006. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541.

Marine Carpuat, Yogarshi Vyas, and Xing Niu. 2017. Detecting cross-lingual semantic divergence for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 69–79, Vancouver. Association for Computational Linguistics.

Julius Cheng and Andreas Vlachos. 2023. Faster minimum bayes risk decoding with confidence-based pruning.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Bryan Eikema and Wilker Aziz. 2020a. Is map decoding all you need? the inadequacy of the mode in neural machine translation. *arXiv preprint arXiv:2005.10283*.

Bryan Eikema and Wilker Aziz. 2020b. Is MAP decoding all you need? the inadequacy of the mode in neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Bryan Eikema and Wilker Aziz. 2022. Sampling-based approximations to minimum Bayes risk decoding for neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10978–10993, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Akhbardeh Farhad, Arkhangorodsky Arkady, Biesialska Magdalena, Bojar Ondřej, Chatterjee Rajen, Chaudhary Vishrav, Marta R Costa-jussa, España-Bonet Cristina, Fan Angela, Federmann Christian, et al. 2021. Findings of the 2021 conference on machine translation (wmt21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88. Association for Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Mara Finkelstein and Markus Freitag. 2023. Mbr and qe finetuning: Training-time distillation of the best and most expensive decoding methods. *arXiv preprint arXiv:2309.10966*.

Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023. Epsilon sampling rocks: Investigating sampling strategies for minimum bayes risk decoding for machine translation. *arXiv preprint arXiv:2305.09860*.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Gemini Team. 2024. Gemini: A family of highly capable multimodal models.

John Hewitt, Christopher D Manning, and Percy Liang. 2022. Truncation sampling as language model desmoothing. *arXiv preprint arXiv:2210.15191*.

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.

Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8003–8017, Toronto, Canada. Association for Computational Linguistics.

Yuu Jinnai and Kaito Ariu. 2024. Hyperparameter-free approach for faster minimum bayes risk decoding.

Wandri Jooste, Rejwanul Haque, and Andy Way. 2022. Knowledge distillation: A method for making neural machine translation more efficient. *Information*, 13(2):88.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google Submission to the WMT 2023 Metrics Shared Task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2021. Nearest neighbor machine translation. In *International Conference on Learning Representations*.

Yoon Kim and Alexander M. Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327, Austin, Texas. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024a. Preliminary wmt24 ranking of general mt systems and llms.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023a. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2023b. Findings of the 2023 conference on machine translation

(wmt23): Llms are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2024b. Preliminary wmt24 ranking of general mt systems and llms. *arXiv preprint arXiv:2407.19884*.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Jiahuan Li, Shanbo Cheng, Shujian Huang, and Jiajun Chen. 2024. MT-PATCHER: Selective and extendable knowledge distillation from large language models for machine translation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6445–6459, Mexico City, Mexico. Association for Computational Linguistics.

Mathias Müller and Rico Sennrich. 2021. Understanding the properties of minimum Bayes risk decoding in neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 259–272, Online. Association for Computational Linguistics.

Jan-Thorsten Peter, David Vilar, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, and Markus Freitag. 2023. There's no data like better data: Using qe metrics for mt data filtering. *arXiv preprint arXiv:2311.05350*.

Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. *arXiv preprint arXiv:2009.09025*.

Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*.

Baris Sumengen, Anand Rajagopalan, Gui Citovsky, David Simcha, Olivier Bachem, Pradipta Mitra, Sam Blasiak, Mason Liang, and Sanjiv Kumar. 2021. Scaling hierarchical agglomerative clustering to billion-sized datasets. *arXiv preprint arXiv:2105.11653*.

Xu Tan, Yi Ren, Di He, Tao Qin, Zhou Zhao, and Tie-Yan Liu. 2018. Multilingual neural machine translation with knowledge distillation. In *International Conference on Learning Representations*.

Christian Tomani, David Vilar, Markus Freitag, Colin Cherry, Subhajit Naskar, Mara Finkelstein, Xavier Garcia, and Daniel Cremers. 2024. Quality-aware translation models: Efficient generation and quality estimation in a single model.

Jannis Vamvas and Rico Sennrich. 2024. Linear-time minimum bayes risk decoding with reference aggregation.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for translation: Assessing strategies and performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15406–15427, Toronto, Canada. Association for Computational Linguistics.

Yuxian Wan, Wenlin Zhang, Zhen Li, Hao Zhang, and Yanxia Li. 2024. Dual knowledge distillation for neural machine translation. *Computer Speech & Language*, 84:101583.

Jun Wang, Eleftheria Briakou, Hamid Dadkhahi, Rishabh Agarwal, Colin Cherry, and Trevor Cohn. 2024. Don't throw away data: Better sequence knowledge distillation. *arXiv preprint arXiv:2407.10456*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*.

Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. 2024. A survey on knowledge distillation of large language models. *arXiv preprint arXiv:2402.13116*.

Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-Woo Lee, and Woomyoung Park. 2021. GPT3Mix: Leveraging large-scale language models for text augmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2225–2239, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Yue Yu, Rongzhi Zhang, Ran Xu, Jieyu Zhang, Jiaming Shen, and Chao Zhang. 2023. Cold-start data selection for better few-shot language model fine-tuning: A prompt-based uncertainty propagation approach. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2499–2521, Toronto, Canada. Association for Computational Linguistics.

Jaume Zaragoza-Bernabeu, Gema Ramírez-Sánchez, Marta Bañón, and Sergio Ortiz Rojas. 2022. Bicleaner AI: Bicleaner goes neural. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 824–831, Marseille, France. European Language Resources Association.

Biao Zhang, Deyi Xiong, Jinsong Su, and Jiebo Luo. 2019. Future-aware knowledge distillation for neural machine translation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(12):2278–2287.

# A    Additional Dataset Statistics

Table 8 shows the average source-to-target length ratios for each of the NewsPaLM MBR and QE datasets. (See Table 3 in Section 2 for the average source lengths and average targets lengths of the NewsPaLM datasets.)

|  | Sentence-level | Blob-level |
|---|---|---|
| EN → DE | 0.9604 | 1.0812 |
| DE → EN | 0.9009 | 0.9729 |

Table 8: Source-to-target length ratios per dataset, computed using the Moses tokenizer.

# B    Additional Results

Table 9 is an extension of Table 4 in Section 4.1, and shows pretraining results across all metrics (including *BLEURT* and *Comet20*) on the WMT'23 en→de and de→en test sets. Table 10 illustrates the stylistic differences between greedy and MBR decoding. Table 11 is an extension of Table 6 in Section 4.2, and shows finetuning results on the WMT'23 test set across all metrics. Table 12 shows all en→de pretraining and finetuning results on the WMT'24 test set. (Note that there does not exist a WMT'24 de→en test set.) Table 13 is the companion to Table 5 in Section 4, but on the WMT'24 (rather than WMT'23) test set.

Tables 14 and 15 show the pretraining and finetuning results on the en→de WMT'23 and WMT'24 test sets, respectively, broken out by domain. For WMT'23, the domains are *Mastodon*, *News*, *Speech*, and *User Review*. For WMT'24, the domains are *Literary*, *News*, *Social*, and *Speech*. Note that the models pretrained and finetuned on our NewsPaLM dataset perform especially strongly on the *News* domain, but the gains aren't limited to this domain.

Figure 4 shows the *PaLM-2 Bison* few-shot versus self-MBR-finetuned results on the en→de WMT'23 test set, bucketed by source segment length. Note that the gains in performance from self-distillation are consistent across all segment length buckets.

Table 16 accompanies Figure 3 in Section 4.3.1, and shows the results of the NewsPaLM subsampling ablations across all metrics.



Figure 4: *PaLM-2 Bison* few-shot versus NewsPaLM MBR-finetuned performance bucketed by source length (en→de WMT'23 test set). Note that self-MBR finetuning (on sentence-level data only) improves performance across all source length buckets.

| | Model | BLEURT ↑ | MetricX ↓ | COMET20 ↑ | COMET22 ↑ |
|---|---|---|---|---|---|
| *en→de* | 1a) WMT'23 (all) | 64.11 | 4.20 | 42.52 | 78.79 |
| | 1b) WMT'23 (sentence-level, *BLEURT-QE* filtered) | 29.33 | 16.69 | -1.14 | 43.18 |
| | 2a) Greedy sentence-level + blob-level (9:1) | **67.75** | **2.60** | **51.20** | **81.67** |
| | 2b) MBR sentence-level | 56.93 | 6.39 | 21.86 | 72.05 |
| | 2c) QE blob-level | 67.18 | 2.82 | 45.72 | 80.62 |
| | 2d) MBR sentence-level + QE blob-level (9:1) | 66.34 | 2.99 | 43.20 | 79.68 |
| *de→en* | 1a) WMT'23 (all) | 69.96 | 5.55 | 51.52 | 82.41 |
| | 1b) WMT'23 (sentence-level, *BLEURT-QE* filtered) | 44.80 | 14.80 | -64.27 | 57.27 |
| | 2a) Greedy sentence-level + blob-level (9:1) | **70.94** | **3.47** | **54.86** | **83.30** |
| | 2b) MBR sentence-level | 67.21 | 4.97 | 44.12 | 80.55 |
| | 2c) QE blob-level | 69.68 | 4.01 | 51.00 | 82.33 |
| | 2d) MBR sentence-level + QE blob-level (9:1) | 69.07 | 3.95 | 49.71 | 82.02 |

Table 9: Pretraining performance (WMT'23 test set).

| | |
|---|---|
| Source | "While President Trump has not commented on Lewis' death, White House press secretary Kayleigh McEnany said in a tweet Saturday, [. . . ] |
| Greedy | Während Präsident Trump sich nicht zu Lewis' Tod geäußert hat, sagte die Pressesprecherin des Weißen Hauses, Kayleigh McEnany, am Samstag in einem Tweet: [. . . ] |
| MBR | Präsident Trump hat bislang keinen Kommentar zu Lewis' Tod abgegeben. Weißes Hauses-Pressesprecherin Kayleigh McEnany äußerte sich hingegen auf Twitter: [. . . ] |
| Explanation | Greedy decoding closely follows the English sentence structure. MBR decoding produces a more natural German sentence, and also decides to introduce a sentence break. |

| | |
|---|---|
| Source | "If we follow the Hungarian philosopher Gáspár Miklós Tamás, with his very broad definition of fascism as ""a break with the enlightenment tradition of citizenship as a universal entitlement,"" the similarities sharpen." |
| Greedy | Wenn wir dem ungarischen Philosophen Gáspár Miklós Tamás mit seiner sehr weiten Definition von Faschismus als „Bruch mit der aufklärerischen Tradition der Staatsbürgerschaft als universellem Recht" folgen, werden die Ähnlichkeiten deutlicher. |
| MBR | Wenn man dem ungarischen Philosophen Gàspár Miklós Tamás folgt, der in seiner sehr breiten Definition des Faschismus eine „Verletzung der aufgeklärten Tradition des Bürgerstatus als universalem Recht" sieht, dann wird die Ähnlichkeit noch deutlicher. |
| Explanation | The MBR translation uses an impersonal form, which is preferred in German. The subordinate clause about the definition of fascism is also reformulated, resulting in a more natural structure. |

Table 10: Comparison between translations generated with greedy and MBR decoding.

| | Model | BLEURT ↑ | MetricX ↓ | COMET20 ↑ | COMET22 ↑ |
|---|---|---|---|---|---|
| | 1a) Greedy sentence-level + blob-level (9:1) | 68.31 | 2.59 | 51.91 | 81.49 |
| | 1b) MBR sentence-level | **70.65** | 2.30 | **55.38** | **82.69** |
| *en→de* | 1c) QE blob-level | 68.19 | 2.45 | 51.97 | 81.83 |
| | 1d) MBR sentence-level + QE blob-level (9:1) | 70.35 | **2.26** | 55.07 | 82.52 |
| | 2a) PaLM-2 five-shot (no finetuning) | 72.34 | 1.62 | 60.62 | 84.54 |
| | 2b) PaLM-2 MBR sentence-level | **74.38** | **1.14** | **64.86** | **85.64** |
| | 2c) PaLM-2 QE blob-level | 72.31 | 1.47 | 61.03 | 84.77 |
| | 2d) PaLM-2 MBR sentence-level + QE blob-level (9:1) | 74.21 | 1.17 | 64.43 | 85.54 |
| | 1a) Greedy sentence-level + blob-level (9:1) | 72.61 | 3.12 | 59.52 | 84.14 |
| | 1b) MBR sentence-level | **73.56** | 2.91 | **61.47** | **84.57** |
| *de→en* | 1c) QE blob-level | 73.02 | 2.99 | 59.73 | 84.27 |
| | 1d) MBR sentence-level + QE blob-level (9:1) | 73.47 | **2.82** | 61.04 | 84.53 |
| | 2a) PaLM-2 five-shot (no finetuning) | 74.73 | 2.25 | 64.72 | 85.36 |
| | 2b) PaLM-2 MBR sentence-level | **76.20** | **1.91** | **68.41** | **86.26** |
| | 2c) PaLM-2 QE blob-level | 75.56 | 2.03 | 66.55 | 85.81 |
| | 2d) PaLM-2 MBR sentence-level + QE blob-level (9:1) | 76.18 | 1.92 | 68.12 | 86.24 |

Table 11: Finetuning performance (WMT'23 test set). Unless otherwise indicated, performance is reported for the encoder-decoder model. For finetuning, this model was initialized from the checkpoint pretrained on the full WMT'23 training dataset (row 1a in Table 9). Results for *PaLM-2 Bison* few-shot prompting versus self-distillation using NewsPaLM MBR and QE data are reported in rows 2a-d.

| Model | BLEURT ↑ | MetricX ↓ | COMET20 ↑ | COMET22 ↑ |
|---|---|---|---|---|
| 1a) PT: WMT'23 (all) | 65.08 | 3.15 | 29.18 | 77.79 |
| 1b) PT: WMT'23 (sentence-level, Bleurt-QE filtered) | 34.62 | 13.06 | -90.56 | 49.38 |
| 1c) PT: Greedy sentence-level + blob-level (9:1) | 64.78 | 2.95 | 27.88 | 77.81 |
| 1d) PT: MBR sentence-level | 55.82 | 5.32 | -0.20 | 69.88 |
| 1e) PT: QE blob-level | 63.80 | 3.17 | 22.34 | 76.72 |
| 1f) PT: MBR sentence-level + QE blob-level (9:1) | 64.27 | 3.13 | 21.32 | 75.93 |
| 2a) FT: Greedy sentence-level + blob-level (9:1) | 67.94 | 2.26 | 39.22 | 80.20 |
| 2b) FT: MBR sentence-level | 70.33 | 1.96 | 41.96 | 80.96 |
| 2c) FT: QE blob-level | 68.14 | 2.27 | 37.07 | 79.88 |
| 2d) FT: MBR sentence-level + QE blob-level (9:1) | 70.04 | 2.00 | 41.83 | 80.81 |
| 2e) FT: PaLM-2 five-shot (no finetuning) | 72.37 | 1.28 | 49.39 | 83.51 |
| 2f) FT: PaLM-2 MBR sentence-level | 73.94 | 1.05 | 53.99 | 84.44 |
| 2g) FT: PaLM-2 QE blob-level | 71.34 | 1.40 | 46.16 | 82.84 |
| 2h) FT: PaLM-2 MBR sentence-level + QE blob-level (9:1) | 73.83 | 1.05 | 53.57 | 84.31 |

Table 12: Pretraining and finetuning performance (en→de WMT'24 test set). The PT prefix indicates pretrained models, and the FT prefix indicates finetuned models. Unless otherwise indicated, performance is reported for the encoder-decoder model. For finetuning, this model was initialized from the checkpoint pretrained on the full WMT'23 training dataset (row 1a). Results for *PaLM-2 Bison* few-shot prompting versus self-distillation using NewsPaLM MBR and QE data are reported in rows 2e-h.

| Model | BLEURT ↑ | MetricX ↓ | COMET20 ↑ | COMET22 ↑ |
|---|---|---|---|---|
| MBR + QE finetuning (from greedy-pretrained ckpt) | **68.12** | **2.41** | **33.75** | **79.37** |
| Greedy finetuning (from MBR + QE-pretrained ckpt) | 65.61 | 2.91 | 27.50 | 77.66 |

Table 13: Comparison of pretraining on NewsPaLM greedy data, then finetuning on NewsPaLM MBR and QE data, versus vice-versa (en→de WMT'24 test set).

| Model | Mastodon | | News | | Speech | | User Review | |
|---|---|---|---|---|---|---|---|---|
| | MetricX ↓ | COMET22 ↑ | MetricX ↓ | COMET22 ↑ | MetricX ↓ | COMET22 ↑ | MetricX ↓ | COMET22 ↑ |
| 1a) PT: WMT'23 (all) | 4.09 | 79.26 | 3.97 | 80.46 | 3.87 | 78.55 | 5.19 | 75.14 |
| 1b) PT: Greedy sentence-level + blob-level (9:1) | 2.33 | 81.98 | 1.59 | 84.72 | 3.44 | 79.30 | 3.74 | 79.10 |
| 1c) PT: MBR sentence-level | 4.83 | 72.22 | 4.41 | 77.98 | 8.09 | 70.78 | 10.88 | 64.17 |
| 1d) PT: QE blob-level | 2.70 | 80.27 | 1.75 | 84.96 | 3.63 | 77.37 | 3.71 | 78.68 |
| 1e) PT: MBR sentence-level + QE blob-level (9:1) | 2.73 | 79.60 | 1.79 | 84.28 | 3.82 | 77.90 | 4.39 | 75.02 |
| 2a) FT: Greedy sentence-level + blob-level (9:1) | 2.50 | 81.32 | 2.05 | 83.73 | 2.94 | 80.40 | 3.23 | 79.69 |
| 2b) FT: MBR sentence-level | 2.03 | 82.87 | 1.87 | 84.94 | 2.75 | 80.75 | 3.01 | 81.15 |
| 2c) FT: QE blob-level | 2.36 | 81.10 | 1.86 | 84.42 | 2.72 | 80.44 | 3.21 | 81.15 |
| 2d) FT: MBR sentence-level + QE blob-level (9:1) | 2.17 | 82.21 | 1.76 | 84.91 | 2.54 | 80.94 | 2.87 | 81.40 |
| 2e) FT: PaLM-2 five-shot (no finetuning) | 1.40 | 84.86 | 1.15 | 86.11 | 2.45 | 82.47 | 1.83 | 83.82 |
| 2f) FT: PaLM-2 MBR sentence-level | 0.97 | 86.00 | 0.88 | 86.48 | 1.70 | 83.33 | 1.22 | 86.21 |

Table 14: Per-domain results on en→de WMT'23 test set. The PT prefix indicates pretrained models, and the FT prefix indicates finetuned models. Unless otherwise indicated, performance is reported for the encoder-decoder model. For finetuning, this model was initialized from the checkpoint pretrained on the full WMT'23 training dataset (row 1a). Results for *PaLM-2 Bison* few-shot prompting versus self-distillation using NewsPaLM MBR data are reported in rows 2e-f.

| Model | Literary | | News | | Social | | Speech | |
|---|---|---|---|---|---|---|---|---|
| | MetricX ↓ | COMET22 ↑ | MetricX ↓ | COMET22 ↑ | MetricX ↓ | COMET22 ↑ | MetricX ↓ | COMET22 ↑ |
| 1a) PT: WMT'23 (all) | 3.79 | 75.95 | 2.86 | 81.20 | 2.78 | 76.44 | 3.67 | 80.04 |
| 1b) PT: Greedy sentence-level + blob-level (9:1) | 6.11 | 68.47 | 1.40 | 84.42 | 2.44 | 77.36 | 2.22 | 82.80 |
| 1c) PT: MBR sentence-level | 9.79 | 59.40 | 3.39 | 79.32 | 4.44 | 68.81 | 4.61 | 75.44 |
| 1d) PT: QE blob-level | 6.53 | 68.01 | 1.38 | 84.01 | 2.69 | 75.95 | 2.39 | 81.38 |
| 1e) PT: MBR sentence-level + QE blob-level (9:1) | 6.11 | 68.47 | 1.40 | 84.42 | 2.44 | 77.36 | 2.22 | 82.80 |
| 2a) FT: Greedy sentence-level + blob-level (9:1) | 3.12 | 77.89 | 1.64 | 84.00 | 2.13 | 78.85 | 2.18 | 82.69 |
| 2b) FT: MBR sentence-level | 2.83 | 78.56 | 1.40 | 84.78 | 1.79 | 79.70 | 1.96 | 83.36 |
| 2c) FT: QE blob-level | 3.13 | 77.38 | 1.75 | 84.13 | 2.10 | 78.42 | 2.26 | 82.51 |
| 2d) FT: MBR sentence-level + QE blob-level (9:1) | 2.91 | 78.26 | 1.36 | 84.88 | 1.85 | 79.56 | 1.96 | 83.11 |
| 2e) FT: PaLM-2 five-shot (no finetuning) | 1.42 | 82.24 | 1.08 | 84.89 | 1.23 | 82.79 | 1.38 | 85.25 |
| 2f) FT: PaLM-2 MBR sentence-level | 1.24 | 83.55 | 0.86 | 85.79 | 0.99 | 83.76 | 1.09 | 85.75 |

Table 15: Per-domain results on en→de WMT'24 test set. The PT prefix indicates pretrained models, and the FT prefix indicates finetuned models. Unless otherwise indicated, performance is reported for the encoder-decoder model. For finetuning, this model was initialized from the checkpoint pretrained on the full WMT'23 training dataset (row 1a). Results for *PaLM-2 Bison* few-shot prompting versus self-distillation using NewsPaLM MBR data are reported in rows 2e-f.

| | Model | BLEURT ↑ | MetricX ↓ | COMET20 ↑ | COMET22 ↑ |
|---|---|---|---|---|---|
| *en→de* | 1a) PT: Full dataset | 56.93 | 6.39 | 21.86 | 72.05 |
| | 1b) PT: Subsampled dataset | 43.73 | 11.02 | -25.72 | 59.75 |
| | 2a) FT: Full dataset | 70.65 | 2.30 | 55.38 | 82.69 |
| | 2b) FT: Subsampled dataset | 70.00 | 2.55 | 53.91 | 82.11 |
| *de→en* | 1a) PT: Full dataset | 67.21 | 4.97 | 44.12 | 80.55 |
| | 1b) PT: Subsampled dataset | 60.43 | 7.41 | 22.89 | 75.91 |
| | 2a) FT: Full dataset | 73.56 | 2.91 | 61.47 | 84.57 |
| | 2b) FT: Subsampled dataset | 73.15 | 2.96 | 59.35 | 84.33 |

Table 16: Comparison of model performance when pretraining and finetuning on the full versus subsampled NewsPaLM MBR dataset (WMT'23 test set). The subsampled dataset is 25% of the size of the full dataset, and was sampled randomly. The PT prefix indicates pretrained models, and the FT prefix indicates finetuned models.

# Is Preference Alignment Always the Best Option to Enhance LLM-Based Translation? An Empirical Analysis

**Hippolyte Gisserot-Boukhlef**[1,4]    **Ricardo Rei**[2]    **Emmanuel Malherbe**[1]
**Céline Hudelot**[4]    **Pierre Colombo**[3,4]    **Nuno M. Guerreiro**[2,4,5,6]
[1]Artefact Research Center    [2]Unbabel    [3]Equall
[4]MICS, CentraleSupélec, Université Paris-Saclay    [5]Instituto de Telecomunicações
[6]Instituto Superior Técnico & Universidade de Lisboa (Lisbon ELLIS Unit)
hippolyte.gisserot-boukhlef@centralesupelec.fr

## Abstract

Neural metrics for machine translation (MT) evaluation have become increasingly prominent due to their superior correlation with human judgments compared to traditional lexical metrics. Researchers have therefore utilized neural metrics through quality-informed decoding strategies, achieving better results than likelihood-based methods. With the rise of Large Language Models (LLMs), preference-based alignment techniques have gained attention for their potential to enhance translation quality by optimizing model weights directly on preferences induced by quality estimators. This study focuses on Contrastive Preference Optimization (CPO) and conducts extensive experiments to evaluate the impact of preference-based alignment on translation quality. Our findings indicate that while CPO consistently outperforms Supervised Fine-Tuning (SFT) on high-quality data with regard to the alignment metric, it may lead to instability across downstream evaluation metrics, particularly between neural and lexical ones. Additionally, we demonstrate that relying solely on the base model for generating candidate translations achieves performance comparable to using multiple external systems, while ensuring better consistency across downstream metrics.[1]

## 1 Introduction

Neural metrics for machine translation evaluation that are trained to mimic human preferences, such as BLEURT (Sellam et al., 2020), COMET (Rei et al., 2020, 2022a), or Metric-X (Juraska et al., 2023), have become increasingly prevalent. These metrics offer greater accuracy and better reflect human judgments compared to traditional lexical metrics (Mathur et al., 2020; Kocmi et al., 2021; Freitag et al., 2022b; Kocmi et al., 2024) like

BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005) or chrF (Popović, 2015), which mainly consider lexical overlap with a reference text. As such, researchers have attempted to leverage these improvements by integrating them directly into translation systems.

One appealing strategy to incorporate quality information to improve downstream translation performance involves using decoding strategies such as N-Best reranking and Minimum Bayes Risk (MBR) decoding (Kumar and Byrne, 2002, 2004; Eikema and Aziz, 2020; Fernandes et al., 2022; Freitag et al., 2022a). These techniques rely on generating multiple candidates to maximize a given quality metric at inference time, and research has shown that they consistently yield better results than likelihood-based decoding techniques (Eikema and Aziz, 2020; Koehn and Knowles, 2017; Ott et al., 2018).

With the rise of decoder-only LLMs in MT, quality-informed fine-tuning techniques have gained significant attention. Unlike decoding-based methods that inject quality information at inference time, fine-tuning modifies model weights using training sets induced with quality information. These approaches include filtering parallel training data based on a quality metric (Alves et al., 2024), distilling gains from more expensive quality-aware methods such as MBR (Finkelstein et al., 2024), or employing preference-based alignment techniques (Rafailov et al., 2024; Xu et al., 2024a), where the model learns preferences induced by quality metrics between candidate translations typically generated by multiple systems. In this work, we focus specifically on the latter.

Alignment techniques represent a paradigm shift from quality-aware inference time approaches, as they optimize the metric of interest *indirectly*. Understanding the impact of these approaches on translation quality is thus a relevant problem. While some studies have examined quality-

---

[1]All relevant preference datasets and aligned models, along with detailed evaluation metrics, are available at https://huggingface.co/collections/artefactory/translation-alignment-analysis.

informed decoding techniques and their influence on translation output (Amrhein and Sennrich, 2022), there is still a gap in our understanding of how preference-based fine-tuning affects translation quality.

In this work, we aim to bridge this gap by examining the properties of preference-based alignment techniques, with a particular focus on Contrastive Preference Optimization (CPO) (Xu et al., 2024a), which has been used successfully to achieve very competitive translation performance. Our analysis seeks to describe the effects of preference-based fine-tuning on downstream performance, specifically regarding alignment effectiveness, the interactions between optimized and non-optimized metrics, and the impact of using multiple candidate translation systems for generating preference data. Through extensive experimentation, we find that:

- Preference-based alignment globally outperforms Supervised Fine-Tuning (SFT) on high-quality data in terms of maximizing the alignment metric.

- However, preference-based alignment is highly sensitive to the choice of candidate systems used for generating preference data, affecting both the alignment metric and downstream metric consistency.

- Aligning a model using its own translations achieves performance comparable to employing multiple external systems, while ensuring better metric consistency and allowing for improved control over the alignment process.

## 2 Background

### 2.1 Quality-Informed Translation

Along with human evaluation, lexical metrics like BLEU (Papineni et al., 2002), chrF (Popović, 2015), METEOR (Banerjee and Lavie, 2005), and ROUGE (Lin, 2004) have long been used for translation evaluation. However, human evaluation is costly, and lexical metrics have been shown to correlate poorly with human judgements.

More recently, some neural metrics have emerged as a preferred method to mimic human preferences without relying on expensive human evaluation. The intuitive approach involves training an encoder model on human-annotated source-translation-reference triplets. Among the metrics most frequently mentioned in the literature are

BLEURT (Yan et al., 2023), COMET (Rei et al., 2020), CometKiwi (Rei et al., 2022b), xCOMET (Guerreiro et al., 2023), and Metric-X (Juraska et al., 2023). They can be divided into two families: *reference-based* metrics, that include a human-written gold reference as an input to the scoring model, and *reference-free* metrics, which only require access to the source sentence and the generated translation. These neural metrics have proven particularly effective at scoring translations and achieve much higher correlation with human judgments than their lexical counterparts (Mathur et al., 2020; Kocmi et al., 2021; Freitag et al., 2022b; Kocmi et al., 2024).

These neural metrics have also been leveraged to improve translation models through decoding strategies. The approach involves sampling various candidate translations, scoring them according to a given metric, and selecting the one with the highest score. This methodology is exemplified by MBR decoding in the reference-based setting and N-best reranking in the reference-free setting (Fernandes et al., 2022; Freitag et al., 2022a).

### 2.2 Quality-Based Fine-Tuning

With the recent rise of decoder-only LLMs applied to translation tasks (Zhu et al., 2023; Jiao et al., 2023; Hendy et al., 2023; Kocmi et al., 2023; Freitag et al., 2023; Xu et al., 2023; Alves et al., 2023; Xu et al., 2024a; Alves et al., 2024), and with automatic metrics increasingly reflecting human judgments (Sellam et al., 2020; Rei et al., 2020; Juraska et al., 2023), quality-based fine-tuning has gained considerable traction. This approach shifts the objective from selecting the best candidate translation according to a metric at inference time to directly updating model weights through fine-tuning to produce the desired translations. A straightforward approach is to perform SFT on high-quality translations, evaluated and then filtered with respect to a metric of interest (Alves et al., 2024).

Another attractive alternative is Preference Optimization (PO) (Simianer, 2018; Rafailov et al., 2024; Xu et al., 2024a; Yang et al., 2023; Xu et al., 2024b; Wu et al., 2024), which focuses on learning preferences between chosen and rejected translations rather than simply increasing the likelihood of high-quality sentences. A popular PO method is Direct Preference Optimization (DPO) (Rafailov et al., 2024), which aims to maximize a scaled likelihood gap between a chosen and a rejected option.

More recently, CPO (Xu et al., 2024a) has emerged as a promising alternative, incorporating an SFT term into the DPO loss, effectively combining the strengths of both methods. Moreover, by removing the reference policy from the learning objective, it improves training efficiency.

## 3 Experimental Setup

Here, we detail our experimental setup, explaining how we built the preference data, and train and evaluate the models.

### 3.1 Preference Data

**Preference datasets.** To build a preference dataset, one needs candidate translations, an evaluation metric $m$ to score these translations, and a method to select chosen and rejected hypotheses. We denote a candidate dataset by

$$\mathcal{D} = \{(x_i, \mathcal{Y}_i)\}_{i=1}^N,$$

where $x_i$ denotes the source sentence and $\mathcal{Y}_i$ is a set of candidate translations. One can then derive a preference dataset,

$$\mathcal{D}_{pref} = \{(x_i, y_i^r, y_i^c)\}_{i=1}^N,$$

where $y_i^c \in \mathcal{Y}_i$ (chosen hypothesis) is a translation preferred to $y_i^r \in \mathcal{Y}_i$ (rejected hypothesis) according to a metric $m$ and a given selection method.

**Multi-system approach.** In the multi-system scenario, we follow the setting outlined by Xu et al. (2024a). Candidate translations are generated using three different systems, namely ALMA-13B-LoRA (the base model we aim to align, referred to as Base) (Xu et al., 2023), GPT-4 (OpenAI, 2023), and the human-written gold reference (referred to as Ref). Formally, for all data samples,

$$\mathcal{Y}_i^{multi} = \left\{ y_i^{Ref}, y_i^{Base}, y_i^{GPT\text{-}4} \right\}.$$

Then, for each sample, the three translations are evaluated with regard to $m$. The one with highest (resp. lowest) score is selected as the chosen (resp. rejected) hypothesis. Formally,

$$y_i^c = \underset{y \in \mathcal{Y}_i^{multi}}{\arg\max}\, m\,(y) \,\wedge\, y_i^r = \underset{y \in \mathcal{Y}_i^{multi}}{\arg\min}\, m\,(y)$$

**Mono-system approach.** In the mono-system setting, we solely rely on the base model for candidate generation. For each source sentence, $K = 50$ candidates are top-$p$-sampled ($p = 0.6$) with a temperature $\tau = 0.9$,[2] and are then ranked based on

evaluation metric $m$. For all samples, this results in a set of candidates

$$\mathcal{Y}_i^{mono} = \{y_i^1, \cdots, y_i^K\},$$

where $y_i^1 \preceq \cdots \preceq y_i^K$ are sorted in increasing quality order, with no loss of generality. Preference pairs are then derived to ensure that $y_i^r \preceq y_i^{Base} \preceq y_i^c$ holds for all samples. Further details on the construction of mono-system preference datasets are given in Section 5 and Appendix B.1.

**Source dataset.** We rely on the FLORES-200-based (Team et al., 2022) dataset used in Xu et al. (2024a) as a primary data source. It includes over 20000 translation pairs spanning six languages (English (en), Czech (cs), German (de), Icelandic (is), Russian (ru), and Chinese (zh)) and covering ten language directions, either into-English (xx-en) or out-of-English (en-xx).

**Alignment metrics.** In line with Xu et al. (2024a), we rely on reference-free neural metrics, namely xCOMET-QE-XXL (Guerreiro et al., 2023) (referred to as xCOMET-QE), and the WMT'23 version of CometKiwi-XXL (Rei et al., 2023) (denoted by CometKiwi), as well as on a reference-based lexical metric, chrF (Popović, 2015).

### 3.2 Training

**Learning objective.** We focus our diagnosis on CPO (Xu et al., 2024a), which combines a preference term with a likelihood term and achieves state-of-the-art performance in preference-based metric alignment for translation tasks. The empirical loss function is formally expressed as:

$$\mathcal{L}_{CPO} = -\frac{1}{N} \sum_{i=1}^N \left[ \log \sigma \left( \beta \log \frac{\pi_\theta\left(y_i^c|x_i\right)}{\pi_\theta\left(y_i^r|x_i\right)} \right) \right] + \mathcal{L}_{SFT},$$

where $\mathcal{L}_{SFT} = -\frac{1}{N} \sum_{i=1}^N \left[ \log \pi_\theta\left(y_i^c|x_i\right) \right]$ is the negative-log-likelihood loss applied to chosen translations, $\pi_\theta$ is the model to fine-tune, $\sigma$ is the sigmoid function and $\beta$ is a hyperparameter. In our experiments, CPO alignment is consistently compared to vanilla SFT on chosen translations.[3]

**Training parameters.** We replicate the exact same parameters as the ones outlined by Xu et al. (2024a). ALMA-13B-LoRA is LoRA fine-tuned

---

[2]These are the default parameters used in the ALMA paper (Xu et al., 2023, 2024a).

[3]All our models are trained using the code implementation provided by Xu et al. (2024a).

| | xx-en | | | en-xx | | |
|---|---|---|---|---|---|---|
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
| **Base** | 87.80 | 80.86 | 58.53 | 91.91 | 81.17 | 49.49 |
| *Preferences induced with* **xCOMET-QE** | | | | | | |
| SFT | • 89.13 | 81.49 | 59.82 | • 92.38 | 81.67 | 50.28 |
| CPO | • 89.95 | 81.89 | 59.83 | • 92.75 | 83.60 | 47.69 |
| *Preferences induced with* **CometKiwi** | | | | | | |
| SFT | 89.26 | • 81.70 | 60.01 | 92.44 | • 81.93 | 50.49 |
| CPO | 89.82 | • 82.04 | 60.22 | 92.19 | • 83.64 | 48.11 |
| *Preferences induced with* **chrF** | | | | | | |
| SFT | 87.61 | 80.82 | • 56.97 | 92.20 | 81.70 | • 50.30 |
| CPO | 78.51 | 75.62 | • 45.32 | 88.89 | 80.99 | • 42.50 |

Table 1: Comparison between SFT on preferred translations and CPO in the multi-system setting, using xCOMET-QE, CometKiwi and chrF as alignment metrics. The same 3 metrics are reported for evaluation, separately for into-English (xx-en) and out-of-English (en-xx) translations on the WMT'22 dataset. Green shades indicate metric improvements over the base model, while red shades indicate metric decreases. We represent with (•) scenarios where the preference metric matches the evaluation metric. Values in *italic* font denote statistically significant differences between SFT- and CPO-based alignment at the 5% level, based on one-tailed paired Student's $t$-tests.

with rank 16 for one epoch, starting with a learning rate of $10^{-4}$, using inverse square root decay and a batch size of 128. The $\beta$ parameter of the CPO objective function is set equal to 0.1, in line with the original DPO paper by Rafailov et al. (2024).

### 3.3 Evaluation

**Inference setup.** Following other works on LLM-based translation (Alves et al., 2024; Briakou et al., 2024), all generations at inference time are produced using greedy decoding, as it provides maximum computational efficiency while preserving high output quality.[4]

**Evaluation datasets.** We evaluate our approaches on the WMT'22 test dataset, which consists of 17471 source-reference pairs and includes the same ten language pairs as the preference data. Evaluations on WMT'23 test data are provided in Appendix A.

**Evaluation metrics.** We use the same three metrics used to create the preference datasets: xCOMET-QE, CometKiwi, and chrF. Additional evaluation metrics are reported in Appendix A, specifically the reference-based version of Metric-X-Large (referred to as Metric-X) (Juraska et al., 2023), and BLEU (Papineni et al., 2002).

---

[4]Inference is performed using the vLLM library (Kwon et al., 2023).

## 4 Multi-System Preference Fine-Tuning

We begin our analysis by focusing on the multi-system setting (Xu et al., 2024a), in which the chosen and rejected options are derived from a pool of three candidate systems consisting of ALMA-13B-LoRA (base model), GPT-4, and the gold reference.

### 4.1 Top-Level Analysis

**Neural-based alignment improves downstream performance.** Table 1 shows that when aligning with neural metrics (xCOMET-QE or CometKiwi), both SFT on preferred translations and CPO consistently improve performance on the alignment metric across language pairs. We also observe that aligning on xCOMET-QE improves results on CometKiwi, and vice-versa. We hypothesize this may be the result of high correlation between different neural metrics, as they are typically trained on similar data. Overall, these results demonstrate that alignment-based techniques can achieve similar objectives to those of quality-aware decoding approaches like MBR, even though the target metric is only indirectly optimized.

**CPO induces adverse metric effects.** In Table 1, we observe that when aligning with neural metrics, CPO yields significantly greater improvements on the alignment metric compared to SFT. The inclusion of the reject option seems to offer additional benefits over the traditional SFT objective

| | xx-en | | | en-xx | | |
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
|---|---|---|---|---|---|---|
| **Base** | 87.80 | 80.86 | 58.53 | 91.91 | 81.17 | 49.49 |
| *Optimization via* **SFT** | | | | | | |
| *Preferences induced with* **xCOMET-QE** | | | | | | |
| All systems | 89.13 | 81.49 | 59.82 | 92.38 | 81.67 | 50.28 |
| No Base | *89.41* | 81.56 | *60.26* | 92.32 | 81.65 | *50.52* |
| No Ref | *89.32* | *81.58* | *60.08* | 92.22 | *81.33* | *50.05* |
| No GPT-4 | *88.44* | *81.15* | 58.86 | 92.33 | 81.74 | *50.06* |
| *Preferences induced with* **chrF** | | | | | | |
| All systems | 87.61 | 80.82 | 56.97 | 92.20 | 81.70 | 50.30 |
| No Ref | 89.21 | 81.49 | *60.17* | 91.99 | *80.96* | *50.57* |
| *Optimization via* **CPO** | | | | | | |
| *Preferences induced with* **xCOMET-QE** | | | | | | |
| All systems | 89.95 | 81.89 | 59.83 | 92.75 | 83.60 | 47.69 |
| No Base | *89.59* | *81.73* | 59.94 | 92.74 | *83.13* | 48.54 |
| No Ref | 89.91 | 81.86 | *60.59* | 92.44 | *81.97* | *50.67* |
| No GPT-4 | *88.81* | *81.35* | 57.91 | 92.22 | *83.16* | 46.82 |
| *Preferences induced with* **chrF** | | | | | | |
| All systems | 78.51 | 75.62 | 45.32 | 88.89 | 80.99 | 42.50 |
| No Ref | *89.26* | *81.52* | *60.63* | *90.83* | *79.37* | *51.11* |

Table 2: Impact of the systems used for candidate generation on WMT'22 performance in the multi-system setting after undergoing SFT and CPO optimization. Values in *italic* font denote statistically significant differences between all-systems-based alignment and alignment with one system removed, at the $5\%$ significance level, based on one-tailed paired Student's $t$-tests. Evaluation metrics and color codes are the same as in Table 1.

in this context. However, aligning with CPO also introduces adverse effects between neural and lexical metrics for out-of-English translations. More specifically, and consistent with the findings of Xu et al. (2024a), aligning on neural metrics negatively impacts lexical metrics. Importantly, this is further evidence to support recommendations provided in (Kocmi et al., 2024): even though, in most cases, neural and lexical MT evaluation metrics should be positively correlated, we should employ caution when using the same metric for evaluation that was used during training/inference. Nevertheless and perhaps more interestingly, it turns out SFT does not produce such effects, raising the question of whether these contradictory evaluation dynamics seen with CPO stem from the learning objective itself or the mix of candidate systems used.

**Lexical alignment fails to improve downstream performance.** Table 1 shows that preference-based lexical alignment[5] behaves differently com-

pared to neural alignment. Specifically, SFT results are roughly stagnant, showing a slight decrease in chrF for into-English translations and a slight increase for out-of-English translations. In contrast, CPO results in a steep drop across the metric board for both into- and out-of-English translations. Using the gold reference as the chosen system appears to impair downstream performance, especially when performing alignment using CPO.

### 4.2 Impact of the Candidate Systems

We now turn to investigating how much the success of alignment-based fine-tuning depends on the choice of the candidate systems. Unless otherwise specified, we use xCOMET-QE as the alignment metric and examine the performance impact of withdrawing systems from the candidate pool. We perform SFT and CPO on the newly created datasets. We report results in Table 2.

**The choice of the candidate systems impacts alignment performance.** Table 2 shows that for both SFT- and CPO-based methods, removing systems from the pool of candidates significantly af-

---

[5]When performing alignment using a lexical metric like chrF, the chosen translation is by definition the gold reference as long as it is present in the pool of candidates. The translation with the lowest chrF score among the remaining systems

is then rejected.

| | xx-en | | | en-xx | | |
|---|---|---|---|---|---|---|
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
| **Base** | 87.80 | 80.86 | 58.53 | 91.91 | 81.17 | 49.49 |
| *Chosen system set to* **Base** | | | | | | |
| SFT | • 88.17 | 81.08 | 58.91 | • 91.94 | 81.21 | 49.35 |
| CPO | • 87.94 | 81.02 | 58.62 | • 91.75 | 81.06 | 48.56 |
| *Chosen system set to* **Ref** | | | | | | |
| SFT | • 88.04 | 81.06 | 57.73 | • 92.35 | 81.94 | 50.12 |
| CPO | • 81.95 | 77.86 | 48.75 | • 86.97 | 80.01 | 39.81 |
| *Chosen system set to* **GPT-4** | | | | | | |
| SFT | • 89.81 | 81.67 | 60.53 | • 91.96 | 80.83 | 50.73 |
| CPO | • 89.69 | 80.99 | 60.42 | • 90.50 | 78.81 | 50.22 |

Table 3: Impact of imposing the chosen system on WMT'22 downstream performance in the multi-system setting. Values in *italic* font denote statistically significant differences between SFT- and CPO-based alignment at the 5% significance level, based on one-tailed paired Student's $t$-tests. Evaluation metrics and color codes are the same as in Table 1.

fects performance on the alignment metric. This is particularly the case for out-of-English translation with CPO optimization. Notably, removing GPT-4 has the strongest negative impact on downstream xCOMET-QE. This is expected as it is the highest-quality system among the system candidates (see Table 11 in Appendix B).

**Some candidate systems can be harmful to preference-based alignment.** In Section 4.1, we observed CPO negatively impacts en-xx chrF when aligning on neural metrics, unlike SFT on preferred translations. Table 2 suggests this may stem from including gold references in the candidate system pool: removing them eliminates this adverse effect. We also noted in Section 4.1 that lexical alignment fails to improve downstream chrF, with sharp decreases with CPO. This issue is resolved by removing gold references. Overall, candidate system choice affects alignment effectiveness and downstream metric consistency, with CPO showing higher sensitivity to preference settings than SFT.

### 4.3 Impact of the Chosen System

To complement findings from Section 4.2 and further characterize the sensitivity of preference-based alignment, we propose examining downstream performance when the chosen system is fixed to a single system. We create three preference datasets based on xCOMET-QE, in which we either impose the base model, reference or GPT-4 as the chosen system. When applicable, the rejected translation

is selected from the remaining systems (if one has a lower xCOMET-QE than the chosen system); otherwise, the sample is discarded.

**CPO is not robust to the preference setting.** In contrast to the observations made in Section 4.1, Table 3 shows that, under this setup, CPO fails to outperform SFT for both xx-en and en-xx translations. When systematically choosing base translations, CPO is unable to surpass the trivial SFT setting where the base model is fine-tuned on its own translations.[6] Moreover, downstream CPO performance significantly declines when gold references are chosen, underperforming the non-aligned model across all metrics, even including the alignment metric. These results reinforce the claims made in Section 4.2 and indicate a lack of robustness of CPO compared to SFT. In the following section (Section 5), we demonstrate that this instability observed with CPO can be mitigated by using a more normalized preference setting, relying only on the base model for candidate generation.

### 5 Mono-System Preference Fine-Tuning

So far, we have exclusively focused on multi-system alignment, which involves using external models for candidate generation and preference dataset building. Although this approach is common for metric alignment (Luong and Manning,

---

[6]As expected, performing SFT on a model's own greedy predictions has minimal impact on downstream performance.

| | xx-en | | | en-xx | | |
|---|---|---|---|---|---|---|
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
| **Base** | 87.80 | 80.86 | 58.53 | 91.91 | 81.17 | 49.49 |
| *Optimization via* **SFT** | | | | | | |
| Multi-system | ● 89.13 | 81.49 | 59.82 | ● 92.38 | 81.67 | 50.28 |
| Mono-system | ● *88.51* | *81.29* | *59.05* | ● *92.17* | 81.54 | *49.41* |
| *Optimization via* **CPO** | | | | | | |
| Multi-system | ● 89.95 | 81.89 | 59.83 | ● 92.75 | 83.60 | 47.69 |
| Mono-system | ● *89.35* | 81.80 | *59.52* | ● *92.69* | 82.91 | *49.02* |
| Mono-system (opt.) | ● *89.58* | 81.97 | *59.65* | ● *92.87* | 83.47 | *49.11* |

Table 4: Comparison between multi- and mono-system fine-tuning on WMT'22 test data. Alignment is performed on xCOMET-QE for both SFT and CPO. Mono-system (opt.) denotes the model fine-tuned on optimized mono-system preference data. Values in *italic* font denote statistically significant differences between multi-system- and mono-system-based alignment at the $5\%$ significance level. Evaluation metrics and color codes are the same as in Table 1, based on one-tailed paired Student's $t$-tests.

2015; Sennrich et al., 2016; Xu et al., 2024a), some works have shown that a model can be aligned effectively using only its own outputs (Yang et al., 2023; Yuan et al., 2024; Dubey et al., 2024). In this section, we propose to take a closer look at this strategy and identify its potential advantages and disadvantages compared to the multi-system approach. We use xCOMET-QE as the alignment metric. To ensure a fair comparison, we first generate the mono-system dataset to approximately replicate the properties of the multi-system dataset regarding the alignment metric.[7] Details on the construction of mono-system preference datasets are given in Section 3 and Appendix B.1.

### 5.1 Comparison With Multi-System Alignment

**Mono-system alignment improves downstream performance.** Table 4 shows that performing SFT and CPO on a mono-system dataset using xCOMET-QE for alignment results in improved downstream performance across all neural metrics compared to the base model, as observed in the multi-system scenario (Section 4.1). This finding highlights the effectiveness of alignment techniques even when using only the model's own translations for candidate generation, without needing access to high-quality external systems. This is particularly relevant in practical scenarios in which such access may be limited or unavailable.

**CPO consistently outperforms SFT on neural metrics.** Similar to when relying on multiple systems for candidate generation, we observe in Table 4 that CPO outperforms SFT regarding downstream performance on neural metrics. This finding reinforces the observation made in Section 4.1 and tends to confirm the superiority of the CPO objective over SFT on preferred translations in optimizing neural-based alignment performance.

**Mono-system alignment slightly underperforms multi-system alignment.** Table 4 shows that while mono-system alignment increases downstream performance on neural metrics, the improvement levels are not as high as in the multi-system setting. Despite the mono- and multi-system preference datasets being built with the same alignment metric properties, having translations from different distributions, particularly from GPT-4 (cf. Section 4.2 and Table 2), appears to add value for achieving optimized alignment effectiveness.

**Removing external systems almost eliminates the adverse metric effects observed with CPO.** In Section 4.1, we showed that multi-system neural alignment using CPO greatly impacts lexical performance for out-of-English translations. Table 4 demonstrates that mono-system alignment almost completely mitigates these negative effects. While there is still a slight decrease in en-xx chrF, it is much smaller compared to the multi-system scenario. This confirms the findings from Sections 4.2 and 4.3 that CPO is sensitive to the preference setting, but also shows that relying solely on candidate translations from the base model limits adverse ef-

---

[7]The created mono-system dataset has an average rejected/-chosen xCOMET-QE of 87.8/97.3, compared to 87.9/97.2 for the multi-system dataset (Table 11).

Figure 1: Impact of chosen and rejected option quality on downstream performance, using xCOMET-QE for alignment and evaluation. The chart is derived by linearly interpolating results from nine preference datasets (points A to I), each with different average rejected and chosen qualities. Test performance on WMT'22 (average across all language pairs) is reported in brackets. Example: point C (avg. rejected xCOMET-QE: 75.4, avg. chosen: 98.2) achieves 90.9 xCOMET-QE on WMT'22 test data.

fects on downstream metric consistency. A possible explanation is that candidate translations from the same system distribution tend to have similar properties, thereby reducing the likelihood of observing high lexical instability when performing alignment based on a neural metric like xCOMET-QE.

**The mono-system approach offers better control over the alignment process.** Specifically, mono-system alignment provides more fine-grained control over the respective qualities of the chosen and rejected options. This setting allows for tuning these qualities to maximize post-alignment performance, which is not possible when using a limited number of external systems. This aspect is further explored in the following section (Section 5.2).

## 5.2 Optimizing the Preference Data

In this final experiment, we examine how the quality of chosen and rejected options affects downstream performance. We build nine preference datasets, each with varying average xCOMET-QE scores for chosen and rejected options. The hypotheses' average qualities are categorized into three groups: High, Mid, and Low. As detailed in Section 3.1, the quality of the chosen (resp. rejected) option is always ensured to be above (resp. below) the quality of the base translation. The statistics of the created datasets are summarized in

Appendix B.1 (Table 11).

**The respective qualities of the rejected and chosen options have a significant impact on post-CPO performance.** Figure 1 highlights the need to closely monitor the qualities of chosen and rejected options to fully leverage the mono-system approach. Specifically, several properties of preference data were found to negatively impact post-CPO performance: (i) a chosen option of too low quality, (ii) an extremely low or high quality of the rejected option, and (iii) too wide a gap between the qualities of the rejected and chosen options.

**Optimizing preference data yields competitive performance to multi-system setting.** Figure 1 shows that for effective metric alignment with CPO, the rejected option's quality should be moderate (neither too high nor too low), while the chosen option's quality should be as high as possible. Specifically, optimal test performance was obtained with rejected options average around 90% ($\Delta = -10\%$) of the base model's quality, and chosen options averaging around 105% ($\Delta = +5\%$). Under this scenario, we show that performance levels can match those in the multi-system setting while maintaining consistency with lexical scores (Table 4). However, these results also highlight the complexity of achieving optimal preference-based alignment and get the most of the reject option.

## 6 Conclusion

Our experiments revealed several key findings. Firstly, we showed that preference-based alignment, specifically using CPO, globally outperforms SFT on high-quality data in terms of improving neural evaluation metrics. However, we identified significant drawbacks when relying on multiple systems for preference data generation, revealing adverse effects between neural and lexical metrics, and highlighting a lack of robustness in preference-based alignment compared to the SFT approach. Finally, we showed that using candidate translations all originating from the same system distribution, specifically the base model, can be an effective strategy for gaining more control over preference-based fine-tuning. This approach achieves performance comparable to using multiple external systems while ensuring better consistency across evaluation metrics. In a nutshell, while preference-based alignment techniques hold promise for improving MT quality, careful consideration must be given to the choice of candidate translations, the learning objective, and the potential trade-offs regarding downstream metric consistency.

## Limitations

In this work, we conducted extensive experiments to assess the impact of preference-based fine-tuning on downstream translation quality. For efficiency and practicality, we focused on the experimental setup detailed by Xu et al. (2024a), which utilizes three systems for candidate generation. Similarly, we used the same evaluation metrics and datasets. Future experiments could benefit from validating our findings using different model families, a broader range of alignment and evaluation metrics, and additional translation datasets, for instance including other languages.

Additionally, in the mono-system setting, we explored the impact of varying the qualities of chosen and rejected options and derived general insights on optimizing preference data. Further research could involve using different datasets, models, and alignment metrics to characterize more precisely the factors that influence downstream performance in this specific scenario. This approach could lead to a deeper mathematical understanding of the elements that affect performance in preference-based fine-tuning, resulting in more robust and scalable optimization techniques.

Finally, our evaluation relied on automatic metrics, both lexical and neural, with the latter closely approximating human judgments but still being unable to fully replace them. Given their imperfect correlation with human preferences, future work could benefit from additional human evaluation of outputs obtained via the approaches we studied to get an even deeper understanding of post-alignment downstream performance dynamics.

## Ethics Statement

Our work aims to investigate the mechanisms of model alignment to enhance transparency in the field of automatic translation. We believe this effort improves the interpretability of model outputs, which is beneficial for ethical considerations. Additionally, our analysis is distinctly multilingual, with an emphasis on low-resource languages, contributing to expanding the scope of MT. We have identified no potential negative societal impacts from our work.

## Acknowledgements

## References

Duarte M Alves, Nuno M Guerreiro, João Alves, José Pombal, Ricardo Rei, José GC de Souza, Pierre Colombo, and André FT Martins. 2023. Steering large language models for machine translation with finetuning and in-context learning. *arXiv preprint arXiv:2310.13448*.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Chantal Amrhein and Rico Sennrich. 2022. Identifying weaknesses in machine translation metrics through minimum bayes risk decoding: A case study for comet. *arXiv preprint arXiv:2202.05148*.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Eleftheria Briakou, Jiaming Luo, Colin Cherry, and Markus Freitag. 2024. Translating step-by-step: Decomposing the translation process for improved translation quality of long-form texts. *arXiv preprint arXiv:2409.06790*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou

U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models.

Bryan Eikema and Wilker Aziz. 2020. Is MAP decoding all you need? the inadequacy of the mode in neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.

Mara Finkelstein, Subhajit Naskar, Mehdi Mirzazadeh, Apurva Shah, and Markus Freitag. 2024. Mbr and qe finetuning: Training-time distillation of the best and most expensive decoding methods.

Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022a. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022b. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. *arXiv preprint arXiv:2301.08745*.

Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore. Association for Computational Linguistics.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 478–494, Online. Association for Computational Linguistics.

Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024. Navigating the metrics maze: Reconciling score magnitudes and accuracies. *arXiv preprint arXiv:2401.06760*.

Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.

Shankar Kumar and Bill Byrne. 2002. Minimum bayes-risk word alignments of bilingual texts. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing (EMNLP 2002)*, pages 140–147.

Shankar Kumar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 169–176, Boston, Massachusetts, USA. Association for Computational Linguistics.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Minh-Thang Luong and Christopher Manning. 2015. Stanford neural machine translation systems for spoken language domains. In *Proceedings of the 12th International Workshop on Spoken Language Translation: Evaluation Campaign*, pages 76–79, Da Nang, Vietnam.

Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4984–4997, Online. Association for Computational Linguistics.

R OpenAI. 2023. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2(5).

Myle Ott, Michael Auli, David Grangier, and Marc'Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3956–3965. PMLR.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Maja Popović. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.

Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.

Patrick Simianer. 2018. *Preference Learning for Machine Translation*. Ph.D. thesis.

NLLB Team, Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, et al. 2022. No language left behind: Scaling human-centered machine translation (2022). *URL https://arxiv. org/abs/2207.04672*.

Qiyu Wu, Masaaki Nagata, Zhongtao Miao, and Yoshimasa Tsuruoka. 2024. Word alignment as preference for machine translation. *arXiv preprint arXiv:2405.09223*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024a. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *arXiv preprint arXiv:2401.08417*.

Nuo Xu, Jun Zhao, Can Zu, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024b. Advancing translation preference modeling with rlhf: A step towards cost-effective solution. *arXiv preprint arXiv:2402.11525*.

Yiming Yan, Tao Wang, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Mingxuan Wang. 2023. BLEURT has universal translations: An analysis of automatic metrics by minimum risk training. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,

pages 5428–5443, Toronto, Canada. Association for Computational Linguistics.

Guangyu Yang, Jinghong Chen, Weizhe Lin, and Bill Byrne. 2023. Direct preference optimization for neural machine translation with minimum bayes risk decoding. *arXiv preprint arXiv:2311.08380*.

Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Multilingual machine translation with large language models: Empirical results and analysis. *arXiv preprint arXiv:2304.04675*.

# A Additional Results

In this section, we present results on WMT'23 test data. The findings in Tables 5, 6, 7 and 8 support the observations discussed in the main text for the WMT'22 dataset. In Tables 9 and 10, we also provide additional insights, split by language pairs, and include extra metrics, specifically Metric-X and BLEU.

| | xx-en | | | en-xx | | |
|---|---|---|---|---|---|---|
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
| **Base** | 88.00 | 77.74 | 52.30 | 86.19 | 73.08 | 47.31 |
| *Preferences induced with* **xCOMET-QE** | | | | | | |
| SFT | 88.96 | 78.46 | 53.30 | 87.07 | 73.99 | 48.38 |
| CPO | 89.77 | 78.95 | 53.47 | 88.09 | 76.75 | 44.29 |
| *Preferences induced with* **CometKiwi** | | | | | | |
| SFT | 89.03 | 78.57 | 53.53 | 87.11 | 74.21 | 48.45 |
| CPO | 89.58 | 79.16 | 53.97 | 87.25 | 76.71 | 44.48 |
| *Preferences induced with* **chrF** | | | | | | |
| SFT | 87.91 | 77.62 | 51.20 | 86.95 | 73.96 | 48.14 |
| CPO | 81.79 | 72.38 | 41.46 | 83.21 | 74.76 | 37.96 |

Table 5: Comparison between SFT on preferred translations and CPO in the multi-system setting on WMT'23 test data. Notations and formatting are the same as in Table 1.

| | xx-en | | | en-xx | | |
|---|---|---|---|---|---|---|
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
| **Base** | 88.00 | 77.74 | 52.30 | 86.19 | 73.08 | 47.31 |
| *Optimization via* **SFT** | | | | | | |
| *Preferences induced with* **xCOMET-QE** | | | | | | |
| All systems | 88.96 | 78.46 | 53.30 | 87.07 | 73.99 | 48.38 |
| No Base | 89.07 | 78.53 | 53.57 | 86.94 | 73.70 | 48.52 |
| No-Ref | 89.05 | 78.47 | 53.39 | 87.04 | 73.60 | 48.65 |
| No GPT-4 | 88.29 | 78.02 | 52.62 | 87.04 | 74.08 | 48.03 |
| *Preferences induced with* **chrF** | | | | | | |
| All systems | 87.91 | 77.62 | 51.20 | 86.95 | 73.96 | 48.14 |
| No Ref | 88.89 | 78.47 | 53.51 | 86.65 | 73.02 | 49.04 |
| *Optimization via* **CPO** | | | | | | |
| *Preferences induced with* **xCOMET-QE** | | | | | | |
| All systems | 89.77 | 78.95 | 53.47 | 88.09 | 76.75 | 44.29 |
| No Base | 89.52 | 78.54 | 53.44 | 87.66 | 75.84 | 45.27 |
| No Ref | 89.57 | 79.26 | 54.18 | 87.41 | 74.46 | 48.88 |
| No GPT-4 | 89.16 | 78.46 | 51.94 | 87.45 | 76.62 | 43.30 |
| *Preferences induced with* **chrF** | | | | | | |
| All systems | 81.79 | 72.38 | 41.46 | 83.21 | 74.76 | 37.96 |
| No Ref | 88.79 | 78.73 | 54.21 | 85.40 | 71.82 | 49.59 |

Table 6: Impact of candidate systems on WMT'23 downstream performance in the multi-system setting. Notations and formatting are the same as in Table 2.

| | xx-en | | | en-xx | | |
|---|---|---|---|---|---|---|
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
| **Base** | 88.00 | 77.74 | 52.30 | 86.19 | 73.08 | 47.31 |
| | *Chosen system set to* **Base** | | | | | |
| SFT | ● 88.07 | 77.93 | 52.52 | ● 86.52 | 73.27 | 47.52 |
| CPO | ● 88.05 | 77.95 | *52.24* | ● *86.68* | *73.75* | *46.54* |
| | *Chosen system set to* **Ref** | | | | | |
| SFT | ● 88.33 | 77.92 | 51.75 | ● 87.29 | 74.57 | 47.86 |
| CPO | ● *84.06* | 74.22 | *44.53* | ● *81.01* | *73.55* | *34.64* |
| | *Chosen system set to* **GPT-4** | | | | | |
| SFT | ● 89.57 | 79.06 | 54.08 | ● 86.70 | 73.18 | 49.23 |
| CPO | ● *88.99* | *78.64* | *53.95* | ● *85.14* | *71.40* | *48.68* |

Table 7: Impact of the chosen system on WMT'23 downstream performance in the multi-system setting. Notations and formatting are the same as in Table 3.

| | xx-en | | | en-xx | | |
|---|---|---|---|---|---|---|
| | Neural | | Lexical | Neural | | Lexical |
| | xCOMET-QE | CometKiwi | chrF | xCOMET-QE | CometKiwi | chrF |
| **Base** | 88.00 | 77.74 | 52.30 | 86.19 | 73.08 | 47.31 |
| | *Optimization via* **SFT** | | | | | |
| Multi-system | ● 88.96 | 78.46 | 53.30 | ● 87.07 | 73.99 | 48.38 |
| Mono-system | ● *88.55* | *78.17* | *52.74* | ● *86.75* | *73.87* | *47.43* |
| | *Optimization via* **CPO** | | | | | |
| Multi-system | ● 89.77 | 78.95 | 53.47 | ● 88.09 | 76.75 | 44.29 |
| Mono-system | ● *89.33* | *78.78* | *53.17* | ● *87.94* | *76.01* | *46.65* |
| Mono-system (opt.) | ● *89.36* | *78.92* | *53.28* | ● *88.50* | *76.87* | *46.48* |

Table 8: Comparison between multi- and mono-system fine-tuning on WMT'23 test data. Notations and formatting are the same as in Table 4.

Table 9 — top block (cs-en, en-cs, de-en, en-de)

| | cs-en | | | | | en-cs | | | | | de-en | | | | | en-de | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Neural | | Metric-X | Lexical | | Neural | | Metric-X | Lexical | | Neural | | Metric-X | Lexical | | Neural | | Metric-X | Lexical | |
| | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU |
| **Base** | 83.42 | 82.57 | 2.00 | 65.81 | 41.25 | 90.80 | 81.96 | 1.48 | 53.58 | 27.42 | 93.33 | 83.48 | 2.03 | 55.24 | 29.02 | 96.32 | 80.87 | 1.22 | 56.95 | 27.65 |
| **SFT** | | | | | | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 86.18 | 83.17 | 1.98 | 67.36 | 42.81 | 91.49 | 82.69 | 1.42 | 55.02 | 28.58 | 93.88 | 83.85 | 2.00 | 56.32 | 29.86 | 96.57 | 81.21 | 1.19 | 57.62 | 28.22 |
| No Base | 86.58 | 83.01 | 1.99 | 67.79 | 42.99 | 91.61 | 82.63 | 1.42 | 55.51 | 28.84 | 93.88 | 83.88 | 2.00 | 56.62 | 30.08 | 96.63 | 81.66 | 1.19 | 57.83 | 28.15 |
| No Ref | 86.43 | 83.20 | 1.97 | 67.50 | 42.93 | 91.28 | 82.10 | 1.45 | 54.94 | 28.75 | 93.88 | 83.96 | 1.99 | 56.63 | 30.06 | 96.41 | 81.10 | 1.20 | 57.49 | 28.08 |
| No GPT-4 | 85.23 | 83.03 | 2.02 | 66.55 | 42.21 | 91.53 | 82.78 | 1.42 | 54.52 | 28.07 | 93.41 | 83.59 | 2.02 | 55.41 | 29.10 | 96.49 | 81.22 | 1.20 | 57.38 | 28.14 |
| Chosen = Base | 84.64 | 82.87 | 2.03 | 66.39 | 41.99 | 90.77 | 81.69 | 1.50 | 53.52 | 27.67 | 93.55 | 83.61 | 2.03 | 55.82 | 29.48 | 96.43 | 81.14 | 1.22 | 56.97 | 27.80 |
| Chosen = Ref | 83.78 | 82.83 | 2.04 | 65.74 | 41.62 | 91.36 | 82.65 | 1.40 | 54.35 | 27.73 | 93.35 | 83.69 | 2.02 | 54.67 | 28.48 | 96.67 | 81.60 | 1.19 | 57.32 | 27.99 |
| Chosen = GPT-4 | 86.97 | 82.69 | 2.00 | 67.62 | 41.99 | 91.14 | 81.33 | 1.48 | 56.27 | 28.97 | 94.22 | 83.95 | 1.98 | 57.02 | 30.10 | 96.29 | 80.92 | 1.21 | 58.05 | 28.00 |
| → CometKiwi | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 86.17 | 83.31 | 1.97 | 67.50 | 43.08 | 91.57 | 83.01 | 1.40 | 55.29 | 28.69 | 94.00 | 84.07 | 1.98 | 56.40 | 29.79 | 96.58 | 81.47 | 1.19 | 57.77 | 28.13 |
| → chrF | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 83.63 | 82.85 | 2.06 | 65.22 | 40.77 | 91.42 | 82.45 | 1.41 | 54.76 | 28.11 | 92.86 | 83.53 | 2.03 | 54.26 | 28.29 | 96.60 | 81.63 | 1.18 | 57.50 | 28.06 |
| No Ref | 86.15 | 82.91 | 1.99 | 67.67 | 43.23 | 91.29 | 81.85 | 1.48 | 55.83 | 29.24 | 93.91 | 83.85 | 2.00 | 56.58 | 30.04 | 96.32 | 80.74 | 1.22 | 57.97 | 28.21 |
| *Mono-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 85.00 | 83.11 | 1.99 | 66.38 | 41.81 | 91.04 | 82.23 | 1.45 | 53.27 | 27.42 | 93.71 | 83.80 | 2.00 | 55.84 | 29.27 | 96.46 | 81.16 | 1.21 | 57.08 | 27.94 |
| **CPO** | | | | | | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 87.40 | 83.58 | 1.94 | 67.52 | 42.54 | 90.86 | 84.58 | 1.40 | 50.58 | 23.44 | 94.22 | 84.10 | 1.94 | 56.21 | 29.44 | 97.34 | 83.38 | 1.12 | 55.95 | 26.20 |
| No Base | 87.59 | 83.18 | 1.95 | 67.90 | 41.86 | 91.60 | 82.63 | 1.41 | 56.16 | 28.97 | 94.28 | 84.07 | 1.98 | 56.98 | 30.02 | 96.57 | 82.02 | 1.18 | 57.86 | 27.66 |
| No Ref | 86.68 | 83.22 | 1.96 | 67.46 | 42.84 | 90.94 | 84.14 | 1.38 | 51.48 | 24.91 | 93.97 | 84.06 | 1.96 | 56.32 | 29.80 | 97.16 | 82.67 | 1.12 | 56.53 | 27.20 |
| No GPT-4 | 84.95 | 83.22 | 1.99 | 65.99 | 41.31 | 90.45 | 83.91 | 1.39 | 49.32 | 22.04 | 93.65 | 83.88 | 1.96 | 54.95 | 28.36 | 97.31 | 83.27 | 1.12 | 55.30 | 25.17 |
| Chosen = Base | 84.10 | 82.69 | 2.05 | 65.94 | 41.77 | 90.42 | 81.80 | 1.51 | 54.29 | 26.71 | 93.41 | 83.48 | 2.04 | 55.53 | 29.34 | 96.57 | 79.91 | 1.12 | 47.32 | 21.40 |
| Chosen = Ref | 71.82 | 79.56 | 2.22 | 54.77 | 28.59 | 79.03 | 80.03 | 1.59 | 40.19 | 14.25 | 89.61 | 81.47 | 2.08 | 47.32 | 21.40 | 96.57 | 79.91 | 1.12 | 48.46 | 18.46 |
| Chosen = GPT-4 | 87.70 | 81.77 | 2.07 | 66.72 | 39.42 | 89.73 | 78.83 | 1.57 | 56.21 | 27.28 | 94.17 | 83.06 | 2.03 | 57.19 | 29.47 | 95.68 | 79.61 | 1.27 | 57.68 | 26.23 |
| → CometKiwi | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 86.74 | 83.46 | 1.95 | 67.67 | 42.42 | 90.21 | 84.81 | 1.43 | 51.23 | 23.53 | 94.20 | 84.22 | 1.93 | 56.45 | 29.61 | 97.19 | 83.71 | 1.13 | 56.24 | 25.41 |
| → chrF | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 66.31 | 76.41 | 2.53 | 51.03 | 25.35 | 83.47 | 80.25 | 1.72 | 43.35 | 15.80 | 87.27 | 79.30 | 2.29 | 44.28 | 19.34 | 96.87 | 81.49 | 1.14 | 51.01 | 19.73 |
| No Ref | 86.65 | 82.99 | 1.98 | 68.33 | 43.19 | 89.83 | 79.38 | 1.59 | 57.08 | 28.51 | 93.94 | 83.67 | 2.02 | 56.68 | 29.76 | 95.84 | 80.01 | 1.27 | 58.68 | 27.66 |
| *Mono-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 85.99 | 83.66 | 1.93 | 67.10 | 42.09 | 91.40 | 83.99 | 1.39 | 52.46 | 26.30 | 94.01 | 84.14 | 1.95 | 55.97 | 29.18 | 96.89 | 82.29 | 1.15 | 56.53 | 27.25 |
| Optimized | 86.51 | 83.84 | 1.90 | 67.08 | 41.73 | 91.47 | 84.28 | 1.37 | 52.40 | 25.92 | 94.26 | 84.32 | 1.94 | 56.15 | 29.26 | 97.13 | 82.78 | 1.12 | 56.72 | 27.09 |

Table 9 — middle block (is-en, en-is, ru-en, en-ru)

| | is-en | | | | | en-is | | | | | ru-en | | | | | en-ru | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU |
| **Base** | 76.22 | 85.36 | 1.89 | 59.72 | 35.34 | 89.15 | 80.68 | 2.40 | 53.31 | 23.49 | 89.68 | 80.92 | 1.82 | 62.72 | 35.29 | 92.77 | 82.62 | 2.04 | 51.93 | 25.86 |
| **SFT** | | | | | | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 79.24 | 85.88 | 1.82 | 62.05 | 37.49 | 89.18 | 80.63 | 2.38 | 53.14 | 22.90 | 90.59 | 81.35 | 1.77 | 64.18 | 36.97 | 93.18 | 83.18 | 1.96 | 53.02 | 26.79 |
| No Base | 80.12 | 86.09 | 1.84 | 62.77 | 38.11 | 88.65 | 80.19 | 2.45 | 53.23 | 22.95 | 90.85 | 81.40 | 1.77 | 64.60 | 37.20 | 93.02 | 83.10 | 1.97 | 53.26 | 26.91 |
| No Ref | 79.67 | 85.94 | 1.83 | 62.52 | 37.97 | 88.74 | 80.10 | 2.47 | 52.93 | 22.84 | 90.82 | 81.50 | 1.77 | 64.47 | 37.24 | 92.90 | 82.73 | 2.04 | 52.73 | 26.29 |
| No GPT-4 | 77.64 | 85.63 | 1.87 | 60.33 | 35.79 | 89.43 | 81.05 | 2.34 | 53.38 | 23.27 | 90.04 | 81.15 | 1.80 | 63.20 | 36.19 | 93.13 | 83.05 | 1.95 | 52.82 | 26.65 |
| Chosen = Base | 76.85 | 85.60 | 1.87 | 59.94 | 35.60 | 88.67 | 80.35 | 2.41 | 52.58 | 22.53 | 89.87 | 80.99 | 1.82 | 63.08 | 36.01 | 92.87 | 82.64 | 2.05 | 52.13 | 26.22 |
| Chosen = Ref | 76.86 | 85.52 | 1.87 | 59.88 | 35.41 | 89.44 | 81.24 | 2.30 | 53.45 | 23.60 | 89.39 | 80.92 | 1.80 | 61.40 | 34.43 | 93.29 | 83.43 | 1.92 | 52.74 | 26.54 |
| Chosen = GPT-4 | 81.03 | 86.01 | 1.84 | 63.09 | 37.61 | 88.26 | 79.01 | 2.56 | 53.17 | 22.44 | 91.11 | 81.45 | 1.77 | 64.91 | 37.39 | 92.47 | 82.44 | 2.27 | 53.53 | 26.96 |
| → CometKiwi | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 79.38 | 86.03 | 1.81 | 62.26 | 37.69 | 89.06 | 80.95 | 2.37 | 53.42 | 23.30 | 90.78 | 81.52 | 1.75 | 64.31 | 37.09 | 93.20 | 83.30 | 1.96 | 53.22 | 26.96 |
| → chrF | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 76.29 | 85.36 | 1.89 | 59.07 | 34.61 | 88.99 | 80.46 | 2.35 | 53.09 | 23.20 | 88.63 | 80.51 | 1.84 | 60.40 | 33.44 | 93.11 | 83.31 | 1.93 | 52.80 | 26.77 |
| No Ref | 79.56 | 85.96 | 1.83 | 62.19 | 37.70 | 88.34 | 79.65 | 2.54 | 53.22 | 22.88 | 90.79 | 81.40 | 1.77 | 64.56 | 37.30 | 92.67 | 82.50 | 2.06 | 53.29 | 26.90 |
| *Mono-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 77.44 | 85.77 | 1.85 | 60.19 | 35.69 | 89.11 | 80.46 | 2.38 | 52.95 | 23.16 | 90.25 | 81.27 | 1.78 | 63.52 | 36.22 | 93.14 | 83.07 | 2.00 | 52.13 | 26.40 |
| **CPO** | | | | | | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 80.73 | 86.12 | 1.81 | 63.01 | 38.10 | 89.28 | 82.70 | 2.08 | 51.70 | 21.29 | 91.14 | 81.57 | 1.72 | 63.88 | 35.78 | 94.52 | 85.51 | 1.71 | 50.60 | 23.97 |
| No Base | 81.03 | 85.91 | 1.84 | 63.29 | 37.90 | 89.15 | 80.18 | 2.47 | 53.47 | 22.46 | 91.11 | 81.61 | 1.77 | 64.66 | 36.55 | 92.89 | 83.40 | 1.95 | 53.27 | 26.27 |
| No Ref | 79.90 | 86.07 | 1.80 | 62.63 | 37.81 | 89.67 | 82.02 | 2.10 | 51.95 | 22.34 | 91.06 | 81.57 | 1.73 | 64.20 | 36.49 | 94.30 | 84.85 | 1.81 | 51.28 | 25.44 |
| No GPT-4 | 77.75 | 85.63 | 1.84 | 60.46 | 35.73 | 88.70 | 82.37 | 2.06 | 51.02 | 20.54 | 90.04 | 81.07 | 1.78 | 61.58 | 33.91 | 94.44 | 85.41 | 1.70 | 49.69 | 23.09 |
| Chosen = Base | 75.89 | 85.47 | 1.93 | 59.27 | 35.04 | 88.83 | 80.38 | 2.37 | 52.21 | 22.51 | 89.76 | 81.01 | 1.81 | 62.86 | 35.46 | 92.57 | 82.11 | 2.13 | 51.33 | 25.64 |
| Chosen = Ref | 62.78 | 82.24 | 2.09 | 51.44 | 25.64 | 76.02 | 77.65 | 2.46 | 43.95 | 13.38 | 83.04 | 77.41 | 2.01 | 51.27 | 23.21 | 93.39 | 82.63 | 1.73 | 42.24 | 15.76 |
| Chosen = GPT-4 | 81.02 | 85.15 | 1.92 | 62.85 | 36.17 | 86.51 | 76.73 | 2.88 | 52.54 | 20.88 | 90.95 | 80.91 | 1.83 | 64.68 | 36.06 | 91.13 | 80.75 | 2.21 | 53.43 | 25.58 |
| → CometKiwi | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 80.83 | 86.17 | 1.79 | 63.15 | 37.98 | 88.15 | 82.29 | 2.30 | 52.23 | 21.16 | 91.06 | 81.80 | 1.72 | 64.15 | 36.14 | 93.99 | 85.44 | 1.77 | 51.04 | 23.54 |
| → chrF | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 57.08 | 80.43 | 2.40 | 47.93 | 22.21 | 81.59 | 78.80 | 2.34 | 46.13 | 14.70 | 79.00 | 75.29 | 2.25 | 47.38 | 20.28 | 93.52 | 83.83 | 1.82 | 45.51 | 17.75 |
| No Ref | 80.13 | 85.77 | 1.85 | 62.87 | 37.63 | 86.51 | 77.57 | 2.87 | 53.35 | 22.12 | 90.89 | 81.43 | 1.79 | 64.57 | 36.80 | 91.54 | 81.04 | 2.18 | 53.86 | 26.18 |
| *Mono-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 79.33 | 86.05 | 1.80 | 61.78 | 37.42 | 89.75 | 82.17 | 2.14 | 52.71 | 22.69 | 90.78 | 81.44 | 1.75 | 63.50 | 35.72 | 93.79 | 83.99 | 1.87 | 52.03 | 26.09 |
| Optimized | 79.16 | 86.02 | 1.78 | 62.25 | 37.42 | 89.81 | 82.81 | 2.11 | 52.79 | 22.60 | 90.96 | 81.60 | 1.72 | 63.69 | 35.57 | 93.95 | 84.77 | 1.84 | 52.18 | 26.00 |

Table 9 — bottom block (zh-en, en-zh, xx-en, xx-xx)

| | zh-en | | | | | en-zh | | | | | xx-en | | | | | xx-xx | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU | xCOMET-QE | CometKiwi | Metric-X | chrF | BLEU |
| **Base** | 89.49 | 74.32 | 4.10 | 51.25 | 21.69 | 89.10 | 79.48 | 2.30 | 33.62 | 34.52 | 87.80 | 80.86 | 2.42 | 58.53 | 31.78 | 91.91 | 81.17 | 1.83 | 49.49 | 28.28 |
| **SFT** | | | | | | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 90.09 | 75.50 | 3.98 | 51.82 | 21.72 | 89.84 | 80.12 | 2.23 | 34.07 | 35.01 | 89.13 | 81.49 | 2.37 | 59.82 | 32.92 | 92.38 | 81.67 | 1.77 | 50.28 | 28.91 |
| No Base | 90.29 | 75.71 | 3.96 | 52.30 | 21.97 | 89.80 | 79.92 | 2.27 | 34.14 | 35.09 | 89.41 | 81.56 | 2.37 | 60.26 | 33.19 | 92.32 | 81.65 | 1.79 | 50.52 | 29.00 |
| No Ref | 90.25 | 75.55 | 3.98 | 51.98 | 21.80 | 89.99 | 79.98 | 2.24 | 33.79 | 34.86 | 89.32 | 81.58 | 2.35 | 58.86 | 32.17 | 92.22 | 81.33 | 1.77 | 50.06 | 28.76 |
| No GPT-4 | 89.68 | 74.71 | 4.11 | 51.11 | 21.42 | 89.57 | 80.24 | 2.23 | 33.89 | 34.86 | 88.44 | 81.15 | 2.42 | 58.56 | 32.17 | 92.33 | 81.74 | 1.77 | 50.05 | 28.77 |
| Chosen = Base | 89.41 | 74.70 | 4.14 | 51.38 | 21.65 | 89.30 | 79.79 | 2.26 | 33.20 | 34.16 | 88.17 | 81.08 | 2.43 | 58.91 | 32.21 | 91.94 | 81.21 | 1.83 | 49.35 | 28.26 |
| Chosen = Ref | 90.21 | 74.68 | 4.05 | 49.70 | 20.34 | 89.52 | 80.41 | 2.23 | 34.42 | 35.20 | 88.04 | 81.06 | 2.41 | 57.73 | 31.20 | 92.35 | 81.94 | 1.75 | 50.12 | 28.73 |
| Chosen = GPT-4 | 90.60 | 76.38 | 3.87 | 52.71 | 22.06 | 89.74 | 79.50 | 2.30 | 33.86 | 34.78 | 89.81 | 81.67 | 2.34 | 60.53 | 33.02 | 91.96 | 80.83 | 1.85 | 50.73 | 28.89 |
| → CometKiwi | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 90.25 | 75.82 | 3.94 | 52.21 | 21.98 | 90.09 | 80.45 | 2.22 | 34.23 | 35.16 | 89.26 | 81.70 | 2.34 | 60.01 | 33.06 | 92.44 | 81.93 | 1.77 | 50.49 | 29.03 |
| → chrF | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 90.06 | 74.28 | 4.10 | 48.66 | 19.82 | 89.26 | 80.01 | 2.26 | 34.75 | 35.52 | 87.61 | 80.82 | 2.44 | 56.97 | 30.56 | 92.20 | 81.70 | 1.77 | 50.30 | 28.91 |
| No Ref | 90.06 | 75.61 | 4.01 | 52.39 | 22.13 | 89.49 | 79.38 | 2.33 | 33.89 | 34.94 | 89.21 | 81.49 | 2.38 | 60.17 | 33.23 | 91.99 | 80.96 | 1.86 | 50.57 | 29.06 |
| *Mono-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 89.74 | 74.84 | 4.04 | 51.39 | 21.49 | 89.57 | 80.24 | 2.24 | 33.45 | 34.26 | 88.51 | 81.29 | 2.39 | 59.05 | 32.15 | 92.17 | 81.54 | 1.80 | 49.41 | 28.37 |
| **CPO** | | | | | | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 91.03 | 76.32 | 3.69 | 51.65 | 21.09 | 89.99 | 81.38 | 2.13 | 31.67 | 31.03 | 89.95 | 81.89 | 2.27 | 59.83 | 32.41 | 92.75 | 83.60 | 1.64 | 47.69 | 25.63 |
| No Base | 90.50 | 76.60 | 3.75 | 52.96 | 21.86 | 90.32 | 80.72 | 2.24 | 34.00 | 34.28 | 89.91 | 81.86 | 2.31 | 60.59 | 32.77 | 92.44 | 81.97 | 1.78 | 50.67 | 28.55 |
| No Ref | 90.80 | 75.98 | 3.79 | 51.96 | 21.67 | 90.09 | 81.40 | 2.09 | 33.20 | 33.49 | 89.59 | 81.73 | 2.30 | 59.94 | 32.82 | 92.74 | 83.13 | 1.65 | 48.54 | 27.17 |
| No GPT-4 | 91.27 | 75.23 | 3.79 | 49.49 | 19.91 | 88.41 | 80.44 | 2.26 | 32.65 | 33.47 | 88.81 | 81.35 | 2.35 | 57.91 | 30.94 | 92.22 | 83.16 | 1.67 | 46.82 | 24.53 |
| Chosen = Base | 89.59 | 74.75 | 4.11 | 51.34 | 21.65 | 89.06 | 79.83 | 2.27 | 32.40 | 32.22 | 87.94 | 81.02 | 2.44 | 58.62 | 31.94 | 91.75 | 81.06 | 1.85 | 48.56 | 27.60 |
| Chosen = Ref | 90.71 | 70.87 | 4.05 | 41.49 | 14.40 | 84.25 | 78.61 | 2.40 | 26.14 | 24.09 | 81.95 | 77.86 | 2.53 | 48.75 | 22.02 | 86.97 | 80.01 | 1.79 | 39.81 | 17.62 |
| Chosen = GPT-4 | 89.77 | 76.05 | 3.92 | 53.10 | 21.12 | 87.44 | 77.07 | 2.58 | 32.40 | 31.93 | 89.69 | 80.99 | 2.40 | 60.42 | 31.72 | 90.50 | 78.81 | 2.02 | 50.22 | 27.00 |
| → CometKiwi | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 91.03 | 76.69 | 3.61 | 52.68 | 21.40 | 89.35 | 81.27 | 2.23 | 31.90 | 30.95 | 89.82 | 82.04 | 2.24 | 60.22 | 32.58 | 92.19 | 83.64 | 1.71 | 48.11 | 25.35 |
| → chrF | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 89.56 | 68.89 | 4.50 | 38.41 | 12.48 | 85.29 | 79.47 | 2.43 | 28.37 | 26.33 | 78.51 | 75.62 | 2.83 | 45.32 | 19.41 | 88.89 | 80.99 | 1.84 | 42.50 | 19.33 |
| No Ref | 89.43 | 75.91 | 3.97 | 53.42 | 22.37 | 88.20 | 77.93 | 2.52 | 33.71 | 33.97 | 89.26 | 81.52 | 2.38 | 60.63 | 33.08 | 90.83 | 79.37 | 2.00 | 51.11 | 28.32 |
| *Mono-system* | | | | | | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | | | | | | |
| Vanilla | 90.82 | 76.01 | 3.77 | 51.92 | 21.44 | 90.14 | 81.75 | 2.15 | 33.22 | 33.80 | 89.35 | 81.80 | 2.29 | 59.52 | 32.26 | 92.69 | 82.91 | 1.69 | 49.02 | 27.74 |
| Optimized | 91.06 | 76.28 | 3.68 | 51.90 | 21.38 | 90.44 | 82.35 | 2.07 | 33.32 | 33.70 | 89.58 | 81.97 | 2.25 | 59.65 | 32.16 | 92.87 | 83.47 | 1.66 | 49.11 | 27.57 |

Table 9: Comprehensive downstream evaluation for the WMT'22 dataset, reporting xCOMET-QE, CometKiwi, Metric-X, chrF, and BLEU scores for all models and language pairs. Learning objectives are indicated in **bold** font, candidate settings in *italics*, and alignment metrics are preceded by an arrow (→).

**en-cs / de-en / en-de**

| | en-cs Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU | de-en Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU | en-de Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Base** | 85.90 | 73.23 | 1.91 | 52.57 | 27.45 | 84.79 | 76.57 | 3.73 | 66.64 | 39.38 | 84.97 | 71.97 | 3.04 | 61.69 | 32.78 |
| **SFT** | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 87.19 | 74.32 | 1.78 | 54.31 | 28.71 | 85.29 | 77.01 | 3.59 | 67.82 | 40.50 | 85.75 | 72.56 | 2.91 | 61.75 | 32.81 |
| No Base | 86.89 | 73.60 | 1.82 | 54.46 | 28.83 | 85.49 | 77.19 | 3.58 | 68.17 | 40.88 | 85.41 | 72.40 | 2.88 | 61.81 | 32.57 |
| No Ref | 87.44 | 74.13 | 1.83 | 54.43 | 29.20 | 85.46 | 77.19 | 3.64 | 67.90 | 40.46 | 85.78 | 72.71 | 2.91 | 62.05 | 33.15 |
| No GPT4 | 87.42 | 74.57 | 1.81 | 53.84 | 28.62 | 84.76 | 76.80 | 3.74 | 67.05 | 39.47 | 85.67 | 72.60 | 2.90 | 61.59 | 32.63 |
| Chosen = Base | 86.89 | 74.05 | 1.93 | 52.71 | 27.72 | 84.86 | 76.88 | 3.71 | 67.28 | 39.59 | 85.24 | 72.09 | 2.99 | 61.81 | 32.98 |
| Chosen = Ref | 87.68 | 74.94 | 1.81 | 53.56 | 28.18 | 85.08 | 76.78 | 3.62 | 66.26 | 38.70 | 86.14 | 72.93 | 2.86 | 61.36 | 32.16 |
| Chosen = GPT4 | 86.89 | 73.06 | 1.89 | 55.36 | 29.20 | 85.94 | 77.70 | 3.52 | 68.55 | 41.07 | 85.25 | 72.07 | 2.92 | 62.49 | 33.35 |
| → CometKiwi | | | | | | | | | | | | | | | |
| Vanilla | 87.19 | 74.65 | 1.77 | 54.27 | 28.71 | 85.42 | 77.16 | 3.58 | 68.00 | 40.54 | 85.90 | 72.93 | 2.83 | 62.09 | 32.99 |
| → chrF | | | | | | | | | | | | | | | |
| Vanilla | 87.15 | 74.18 | 1.85 | 53.91 | 28.57 | 84.93 | 76.61 | 3.69 | 65.86 | 38.47 | 85.68 | 72.47 | 2.90 | 61.50 | 32.37 |
| No Ref | 87.22 | 73.40 | 1.87 | 54.96 | 29.50 | 85.19 | 77.02 | 3.66 | 67.93 | 40.42 | 85.49 | 72.56 | 2.98 | 62.46 | 33.64 |
| *Mono-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 86.89 | 74.78 | 1.86 | 52.57 | 27.60 | 85.37 | 77.03 | 3.64 | 67.64 | 39.86 | 85.54 | 72.80 | 3.03 | 61.82 | 33.10 |
| **CPO** | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 86.53 | 76.83 | 1.65 | 49.84 | 23.64 | 86.18 | 77.55 | 3.45 | 67.15 | 39.97 | 87.61 | 73.39 | 2.57 | 58.32 | 27.66 |
| No Ref | 87.48 | 74.40 | 1.77 | 55.51 | 28.99 | 86.00 | 77.80 | 3.48 | 68.30 | 40.60 | 85.82 | 73.08 | 2.79 | 62.73 | 32.81 |
| No Base | 86.45 | 76.06 | 1.71 | 50.44 | 25.26 | 85.80 | 77.19 | 3.54 | 67.32 | 40.00 | 87.23 | 73.00 | 2.71 | 59.44 | 29.86 |
| No GPT4 | 86.27 | 76.71 | 1.64 | 49.10 | 23.01 | 85.96 | 77.19 | 3.57 | 65.40 | 37.11 | 87.39 | 72.99 | 2.53 | 57.36 | 26.45 |
| Chosen = Base | 86.71 | 74.16 | 1.92 | 51.74 | 26.79 | 84.95 | 76.91 | 3.74 | 66.64 | 38.79 | 85.67 | 72.46 | 2.98 | 61.20 | 32.37 |
| Chosen = Ref | 72.09 | 72.38 | 1.94 | 37.30 | 12.61 | 82.25 | 74.74 | 4.09 | 55.55 | 24.83 | 86.62 | 69.62 | 2.92 | 49.04 | 17.37 |
| Chosen = GPT4 | 85.67 | 70.85 | 2.06 | 55.59 | 27.58 | 85.41 | 77.18 | 3.53 | 68.10 | 39.75 | 83.74 | 71.58 | 2.92 | 63.11 | 32.22 |
| → CometKiwi | | | | | | | | | | | | | | | |
| Vanilla | 85.70 | 76.97 | 1.73 | 50.28 | 23.74 | 85.76 | 77.66 | 3.47 | 67.43 | 40.07 | 87.52 | 74.21 | 2.57 | 59.08 | 27.90 |
| → chrF | | | | | | | | | | | | | | | |
| Vanilla | 77.49 | 73.98 | 2.02 | 41.47 | 15.17 | 79.30 | 72.47 | 4.52 | 51.67 | 21.14 | 86.43 | 71.61 | 2.65 | 52.57 | 20.35 |
| No Ref | 85.71 | 71.61 | 2.09 | 56.21 | 28.93 | 85.13 | 77.15 | 3.57 | 68.15 | 40.49 | 84.20 | 71.56 | 3.12 | 63.31 | 33.63 |
| *Mono-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 87.92 | 76.62 | 1.76 | 52.03 | 26.54 | 85.72 | 77.14 | 3.53 | 67.22 | 39.30 | 87.09 | 74.07 | 2.70 | 61.01 | 31.70 |
| Optimized | 88.39 | 77.36 | 1.71 | 52.11 | 26.38 | 86.00 | 77.76 | 3.55 | 67.05 | 38.79 | 87.32 | 74.64 | 2.57 | 60.71 | 31.07 |

**ru-en / en-ru / zh-en**

| | ru-en Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU | en-ru Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU | zh-en Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Base** | 86.22 | 80.04 | 2.56 | 55.59 | 28.45 | 89.08 | 74.89 | 2.63 | 49.90 | 23.99 | 90.45 | 76.06 | 3.68 | 45.44 | 18.79 |
| **SFT** | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 87.66 | 80.62 | 2.53 | 56.62 | 29.52 | 89.48 | 75.83 | 2.50 | 50.74 | 24.80 | 91.11 | 76.98 | 3.61 | 46.37 | 19.72 |
| No Base | 87.96 | 80.55 | 2.55 | 56.80 | 29.50 | 89.53 | 75.86 | 2.52 | 50.72 | 24.66 | 91.04 | 77.15 | 3.57 | 46.70 | 19.93 |
| No Ref | 87.79 | 80.64 | 2.54 | 56.64 | 29.34 | 89.23 | 75.43 | 2.58 | 50.88 | 24.82 | 91.16 | 76.92 | 3.58 | 46.52 | 19.95 |
| No GPT4 | 86.62 | 80.13 | 2.56 | 55.99 | 29.03 | 89.58 | 75.85 | 2.47 | 50.36 | 24.36 | 90.72 | 76.52 | 3.70 | 45.69 | 19.28 |
| Chosen = Base | 86.39 | 80.10 | 2.57 | 55.72 | 28.79 | 88.94 | 74.68 | 2.68 | 50.17 | 24.24 | 90.43 | 76.33 | 3.72 | 44.82 | 18.99 |
| Chosen = Ref | 86.03 | 79.94 | 2.56 | 55.01 | 28.26 | 89.83 | 76.40 | 2.45 | 50.35 | 24.38 | 91.23 | 76.47 | 3.63 | 44.89 | 18.73 |
| Chosen = GPT4 | 88.30 | 80.79 | 2.57 | 57.15 | 29.47 | 88.86 | 75.45 | 2.60 | 51.26 | 24.72 | 91.68 | 77.92 | 3.46 | 47.38 | 20.22 |
| → CometKiwi | | | | | | | | | | | | | | | |
| Vanilla | 87.70 | 80.75 | 2.52 | 56.77 | 29.68 | 89.56 | 75.87 | 2.48 | 50.99 | 24.88 | 91.19 | 77.07 | 3.55 | 46.69 | 19.89 |
| → chrF | | | | | | | | | | | | | | | |
| Vanilla | 85.27 | 79.62 | 2.58 | 54.54 | 28.04 | 89.65 | 75.96 | 2.48 | 50.33 | 24.23 | 91.05 | 76.15 | 3.71 | 44.23 | 18.31 |
| No Ref | 87.63 | 80.64 | 2.56 | 56.67 | 29.43 | 88.88 | 75.00 | 2.61 | 51.14 | 24.83 | 91.02 | 76.99 | 3.61 | 46.74 | 19.95 |
| *Mono-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 86.94 | 80.33 | 2.56 | 55.86 | 28.64 | 89.16 | 75.09 | 2.64 | 50.02 | 24.10 | 90.84 | 76.61 | 3.64 | 45.89 | 18.97 |
| **CPO** | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 88.43 | 80.99 | 2.44 | 56.88 | 29.63 | 91.54 | 79.24 | 2.13 | 48.22 | 21.98 | 91.94 | 77.56 | 3.41 | 46.69 | 19.65 |
| No Ref | 88.50 | 81.12 | 2.50 | 57.26 | 29.55 | 89.33 | 76.05 | 2.47 | 50.98 | 24.10 | 91.50 | 78.04 | 3.45 | 47.58 | 20.12 |
| No Base | 87.98 | 80.69 | 2.46 | 57.06 | 29.77 | 90.90 | 77.72 | 2.26 | 48.94 | 23.02 | 91.89 | 77.09 | 3.45 | 46.42 | 19.67 |
| No GPT4 | 87.00 | 80.29 | 2.47 | 55.39 | 28.24 | 91.41 | 79.47 | 2.13 | 47.24 | 21.10 | 91.92 | 77.21 | 3.46 | 45.18 | 18.62 |
| Chosen = Base | 86.11 | 79.99 | 2.58 | 55.40 | 28.16 | 88.97 | 75.01 | 2.67 | 49.68 | 23.89 | 90.61 | 76.46 | 3.70 | 45.48 | 18.88 |
| Chosen = Ref | 77.69 | 76.31 | 2.70 | 46.94 | 20.32 | 89.79 | 76.43 | 2.20 | 39.65 | 14.17 | 90.11 | 72.25 | 3.88 | 39.36 | 14.29 |
| Chosen = GPT4 | 88.19 | 80.22 | 2.65 | 56.56 | 27.95 | 87.25 | 73.65 | 2.78 | 50.57 | 22.96 | 90.68 | 77.67 | 3.59 | 47.74 | 19.60 |
| → CometKiwi | | | | | | | | | | | | | | | |
| Vanilla | 88.34 | 81.08 | 2.44 | 57.08 | 29.59 | 91.08 | 79.22 | 2.21 | 48.53 | 21.83 | 91.71 | 77.91 | 3.39 | 47.51 | 19.89 |
| → chrF | | | | | | | | | | | | | | | |
| Vanilla | 74.05 | 74.06 | 3.00 | 43.60 | 17.59 | 90.38 | 77.64 | 2.22 | 42.21 | 15.79 | 89.24 | 76.90 | 4.22 | 36.75 | 12.42 |
| No Ref | 87.70 | 80.56 | 2.55 | 57.13 | 29.67 | 87.72 | 74.05 | 2.75 | 51.23 | 23.90 | 90.75 | 77.58 | 3.61 | 47.80 | 20.37 |
| *Mono-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 87.65 | 80.86 | 2.46 | 56.53 | 29.09 | 90.28 | 77.45 | 2.38 | 49.85 | 23.92 | 91.79 | 77.42 | 3.44 | 46.34 | 19.18 |
| Optimized | 87.63 | 80.84 | 2.45 | 56.56 | 29.03 | 90.64 | 78.18 | 2.28 | 49.94 | 23.82 | 91.79 | 77.58 | 3.39 | 46.59 | 19.31 |

**en-zh / xx-en / en-xx**

| | en-zh Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU | xx-en Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU | en-xx Neural xCOMET-QE | CometKiwi | Metric-X | Lexical chrF | BLEU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Base** | 83.92 | 71.42 | 2.56 | 35.58 | 35.94 | 88.00 | 77.74 | 3.23 | 52.30 | 25.37 | 86.19 | 73.08 | 2.42 | 47.31 | 29.43 |
| **SFT** | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 84.89 | 72.22 | 2.47 | 36.50 | 37.00 | 88.96 | 78.46 | 3.17 | 53.30 | 26.38 | 87.07 | 73.99 | 2.30 | 48.38 | 30.39 |
| No Base | 84.80 | 72.00 | 2.50 | 36.80 | 37.47 | 89.07 | 78.53 | 3.16 | 53.57 | 26.52 | 86.94 | 73.70 | 2.33 | 48.52 | 30.50 |
| No Ref | 84.78 | 71.50 | 2.53 | 37.05 | 37.73 | 89.05 | 78.47 | 3.17 | 53.39 | 26.41 | 87.04 | 73.60 | 2.36 | 48.65 | 30.79 |
| No GPT4 | 84.50 | 72.22 | 2.48 | 36.26 | 36.77 | 88.29 | 78.02 | 3.24 | 52.62 | 25.84 | 87.04 | 74.08 | 2.31 | 48.03 | 30.14 |
| Chosen = Base | 84.06 | 71.38 | 2.54 | 35.85 | 36.39 | 88.07 | 77.93 | 3.25 | 52.52 | 25.62 | 86.52 | 73.27 | 2.44 | 47.52 | 29.74 |
| Chosen = Ref | 84.66 | 72.81 | 2.40 | 36.06 | 36.28 | 88.33 | 77.92 | 3.19 | 51.75 | 25.18 | 87.29 | 74.57 | 2.27 | 47.86 | 29.82 |
| Chosen = GPT4 | 84.75 | 71.33 | 2.54 | 37.51 | 38.34 | 89.57 | 79.06 | 3.11 | 54.08 | 26.67 | 86.70 | 73.18 | 2.39 | 49.23 | 30.97 |
| → CometKiwi | | | | | | | | | | | | | | | |
| Vanilla | 84.90 | 72.44 | 2.48 | 36.43 | 36.92 | 89.03 | 78.57 | 3.13 | 53.53 | 26.53 | 87.11 | 74.21 | 2.29 | 48.45 | 30.40 |
| → chrF | | | | | | | | | | | | | | | |
| Vanilla | 84.40 | 72.14 | 2.44 | 36.59 | 36.94 | 87.91 | 77.62 | 3.25 | 51.20 | 24.86 | 86.95 | 73.96 | 2.31 | 48.14 | 30.11 |
| No Ref | 84.16 | 70.78 | 2.58 | 37.42 | 38.38 | 88.89 | 78.47 | 3.19 | 53.51 | 26.44 | 86.65 | 73.02 | 2.41 | 49.04 | 31.13 |
| *Mono-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 84.52 | 72.02 | 2.53 | 35.82 | 36.16 | 88.55 | 78.17 | 3.20 | 52.74 | 25.60 | 86.75 | 73.87 | 2.40 | 47.43 | 29.60 |
| **CPO** | | | | | | | | | | | | | | | |
| *Multi-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 86.32 | 75.09 | 2.36 | 31.02 | 29.15 | 89.77 | 78.95 | 3.02 | 53.47 | 26.32 | 88.09 | 76.75 | 2.09 | 44.29 | 25.15 |
| No Ref | 85.84 | 73.30 | 2.43 | 36.43 | 36.32 | 89.57 | 79.26 | 3.07 | 54.18 | 26.59 | 87.41 | 74.46 | 2.27 | 48.88 | 30.05 |
| No Base | 85.75 | 74.50 | 2.30 | 32.63 | 31.93 | 89.52 | 78.54 | 3.06 | 53.44 | 26.40 | 87.66 | 75.84 | 2.14 | 45.27 | 26.99 |
| No GPT4 | 84.69 | 74.65 | 2.45 | 29.78 | 27.30 | 89.16 | 78.46 | 3.07 | 51.94 | 24.91 | 87.45 | 76.62 | 2.11 | 43.30 | 24.02 |
| Chosen = Base | 84.64 | 72.42 | 2.50 | 34.26 | 34.58 | 88.05 | 77.95 | 3.25 | 52.24 | 25.22 | 86.68 | 73.75 | 2.41 | 46.54 | 28.75 |
| Chosen = Ref | 79.65 | 72.88 | 2.69 | 23.12 | 19.34 | 84.06 | 74.22 | 3.43 | 44.53 | 18.10 | 81.01 | 73.55 | 2.33 | 34.64 | 15.54 |
| Chosen = GPT4 | 82.88 | 69.66 | 2.81 | 36.00 | 35.56 | 88.99 | 78.64 | 3.20 | 53.95 | 25.59 | 85.14 | 71.40 | 2.58 | 48.68 | 28.99 |
| → CometKiwi | | | | | | | | | | | | | | | |
| Vanilla | 84.89 | 74.61 | 2.52 | 30.71 | 28.63 | 89.58 | 79.16 | 3.02 | 53.97 | 26.43 | 87.25 | 76.71 | 2.19 | 44.48 | 24.99 |
| → chrF | | | | | | | | | | | | | | | |
| Vanilla | 80.91 | 73.49 | 2.69 | 26.29 | 22.83 | 81.79 | 72.38 | 3.77 | 41.46 | 15.65 | 83.21 | 74.76 | 2.34 | 37.96 | 18.13 |
| No Ref | 83.08 | 69.85 | 2.81 | 37.64 | 37.80 | 88.79 | 78.73 | 3.17 | 54.21 | 26.74 | 85.40 | 71.82 | 2.60 | 49.59 | 30.49 |
| *Mono-system* | | | | | | | | | | | | | | | |
| → xCOMET-QE | | | | | | | | | | | | | | | |
| Vanilla | 85.85 | 74.47 | 2.33 | 34.22 | 34.00 | 89.33 | 78.78 | 3.06 | 53.17 | 25.80 | 87.94 | 76.01 | 2.20 | 46.65 | 28.45 |
| Optimized | 86.78 | 75.68 | 2.25 | 33.58 | 32.93 | 89.36 | 78.92 | 3.03 | 53.28 | 25.77 | 88.50 | 76.87 | 2.12 | 46.48 | 27.99 |

Table 10: Comprehensive downstream evaluation for the WMT'23 dataset. Metrics, notations and formatting are the same as in Table 9.

# B  Additional Data Details

## B.1  Building Preference Datasets in the Mono-System Setting

Following the experimental setup detailed in the main text (Section 3), we here provide further details on the method used to construct mono-system preference datasets. As a reminder, after generating the $K$ candidate translations for each source sentence, we have, for all $1 \leq i \leq N$,

$$\mathcal{Y}_i^{mono} = \left\{ y_i^1, \cdots, y_i^K \right\},$$

where $y_i^1 \preceq \cdots \preceq y_i^K$ are assumed to be sorted in increasing metric score order. For each sample, we evaluate $y_i^{Base}$ (the greedy-decoded translation) using metric $m$ and check its rank in the set of candidate translations. We denote it by $b_i$. Sorted in increasing quality order, we thereby have

$$y_i^1 \preceq \cdots \preceq y_i^{b_i-1} \preceq y_i^{Base} \preceq y_i^{b_i} \preceq \cdots \preceq y_i^K.$$

Finally, to determine the chosen and rejected hypotheses, we select two offset parameters $o^r, o^c \in \mathbb{N}$, such that the chosen and rejected options are respectively

$$\begin{cases} y_i^c = y_i^{\min(K, b_i + o^c)} \\ y_i^r = y_i^{\max(1, b_i - o^r)} \end{cases}.$$

Intuitively, $o^r$ and $o^c$ control the average quality of the chosen and rejected options in the resulting preference dataset and ensure that the chosen (resp. rejected) option always has a higher (resp. lower) quality than the base translation. Table 11 presents the average quality properties for mono-system preference datasets, and compares them to the multi-system setting.

| | | Neural | | Lexical |
|---|---|---|---|---|
| | Hyp. | xCOMET-QE | CometKiwi | chrF |
| **Multi-system** | | | | |
| Candidate systems | Base | 93.09 | 87.13 | 58.33 |
| | GPT-4 | 94.58 | 88.32 | 60.93 |
| | Reference | 91.84 | 86.72 | 100.00 |
| Vanilla preference dataset | Rejected | 87.86 | 84.15 | 78.48 |
| | Chosen | 97.24 | 89.81 | 75.95 |
| **Mono-system** | | | | |
| Multi-system replica | Rejected | 87.80 | 83.04 | 55.69 |
| | Chosen | 97.29 | 89.20 | 57.18 |
| Chosen = Low / Rejected = Low | Rejected | 75.36 | 75.46 | 52.95 |
| | Chosen | 93.60 | 87.04 | 57.14 |
| Chosen = Low / Rejected = Mid | Rejected | 84.54 | 81.02 | 54.93 |
| | Chosen | 93.60 | 87.04 | 57.14 |
| Chosen = Low / Rejected = High | Rejected | 92.15 | 85.54 | 55.86 |
| | Chosen | 93.60 | 87.04 | 57.14 |
| Chosen = Mid / Rejected = Low | Rejected | 75.36 | 75.46 | 52.95 |
| | Chosen | 95.77 | 88.40 | 57.43 |
| Chosen = Mid / Rejected = Mid | Rejected | 84.54 | 81.02 | 54.93 |
| | Chosen | 95.77 | 88.40 | 57.43 |
| Chosen = Mid / Rejected = High | Rejected | 92.15 | 85.54 | 55.86 |
| | Chosen | 95.77 | 88.40 | 57.43 |
| Chosen = High / Rejected = Low | Rejected | 75.36 | 75.46 | 52.95 |
| | Chosen | 98.16 | 89.84 | 57.56 |
| Chosen = High / Rejected = Mid | Rejected | 84.54 | 81.02 | 54.93 |
| | Chosen | 98.16 | 89.84 | 57.56 |
| Chosen = High / Rejected = High | Rejected | 92.15 | 85.54 | 55.86 |
| | Chosen | 98.16 | 89.84 | 57.56 |

Table 11: Average quality properties for xCOMET-QE-based mono-system preference datasets, compared to the multi-system setting. Multi-system replica is the mono-system dataset that matches the average chosen/rejected qualities of the multi-system preference data. Other mono-system datasets are represented by their relative average chosen/rejected qualities.

## B.2 Language Statistics



Figure 2: Language statistics for preference datasets. The y-axis represents the number of samples, corresponding percentages are displayed above each bar.

Figure 3: Language statistics for WMT'22 and WMT'23 test data. The y-axis represents the number of samples, corresponding percentages are displayed above each bar.

# Quality or Quantity? On Data Scale and Diversity in Adapting Large Language Models for Low-Resource Translation

**Vivek Iyer**    **Bhavitvya Malik**[*]    **Pavel Stepachev**[*]
**Pinzhen Chen**    **Barry Haddow**    **Alexandra Birch**
School of Informatics, University of Edinburgh
`vivek.iyer@ed.ac.uk`

## Abstract

Despite the recent popularity of Large Language Models (LLMs) in Machine Translation (MT), their performance in low-resource languages (LRLs) still lags significantly behind Neural Machine Translation (NMT) models. In this work, we explore what it would take to adapt LLMs for the low-resource setting. Particularly, we re-examine the role of two factors: a) the importance and application of parallel data, and b) diversity in Supervised Fine-Tuning (SFT). Recently, parallel data has seen reduced use in adapting LLMs for MT, while data diversity has been embraced to promote transfer across languages and tasks. However, for low-resource LLM-MT, we show that the opposite is true for both considerations: a) *parallel data* is critical during both pre-training and SFT; b) diversity tends to cause *interference* instead of transfer. Our experiments with three LLMs across two low-resourced language groups—Indigenous American and North-East Indian—reveal consistent trends, underscoring the generalizability of our findings. We believe these insights will be valuable for scaling to massively multilingual LLM-MT models that can effectively serve LRLs.

## 1 Introduction

Large Language Models (LLMs) have been dominating recent research in Machine Translation (MT), showing good few-shot prompting (Garcia et al., 2023; Hendy et al., 2023) and stronger instruction-tuning (Alves et al., 2024; Xu et al., 2024) performances—recently even outperforming commercial Neural Machine Translation (NMT) models (Kocmi et al., 2024). However, LLM translation for low-resource languages (LRLs) still lags significantly behind NMT models (Robinson et al., 2023; Zhu et al., 2024b). While the strong performance of LLMs on high-resource languages can be

---
[*]denotes equal contribution

| | Base LLM | #Tokens | #Langs |
|---|---|---|---|
| Zhang et al., 2024a | From scratch | 283B | 2 |
| Fujii et al., 2024 | Llama2 | 100B | 2 |
| Lu et al., 2024 | Llama{2,3} | ∼82B[†] | 101 |
| Alves et al., 2024 | Llama2 | 20B | 11 |
| Xu et al., 2024 | Llama2 | 20B | 6 |
| Ours | Mistral/Llama3 | 0.7B | 12 |

Table 1: Comparing data scales of previous works with ours, in terms of pre-training token counts, the base LLM (if pre-training continued from one) and how many languages this spanned. [†]Estimated from the reported sentence count assuming 100 tokens per sentence.

attributed to the skewed language distribution during pre-training and the unintentional consumption of parallel data at scale (Briakou et al., 2023), no such relief exists for LRLs. This leads to the main question motivating this paper: *What would it take to adapt LLMs for low-resource MT?*

Recent work on LRL translation with LLMs has explored using resources like multilingual lexicons (Lu et al., 2023), word alignments (Mao and Yu, 2024) or linguistic tools (Zhang et al., 2024b). While effective, reliance on such tools hinders ease of extensibility across languages. Instead, in this work, we take inspiration from research done for high-resource translation with LLMs, where a 2-stage training paradigm of Continued Pre-Training (CPT), followed by small-scale Supervised Fine-Tuning (SFT; (Xu et al., 2024; Alves et al., 2024)) has been successful. Aiming to adapt this framework for low-resource MT, we re-examine the role of two factors influencing the performance of translation LLMs: a) *how best to leverage parallel data*, and b) *the interplay between diversity and transfer* during SFT (also known as 'instruction tuning').

Recently, the role of **parallel data at scale**, long viewed as fundamental to the success of NMT models, has come into question in the era of LLM-MT systems. Motivated by the modest gains of training on 300M parallel sentences (Yang et al., 2023),

and the surprising benefits of scaling *down* during SFT (Zhou et al., 2023), subsequent works have used only tens of thousands of human-written bitext for LLM-MT (Zhang et al., 2023b; Alves et al., 2024; Xu et al., 2024)—with SFT scaling laws further showing the early plateau of LLM-MT performance (Zhang et al., 2024a). Even more surprisingly, Zhu et al. (2024a) showed MT abilities emerging with just 32 SFT examples! However, these explorations concern LLMs pre-trained on several billions of tokens in the languages in question. We revisit these notions for low-resource MT and work with languages having datasets that are 2-3 orders of magnitude smaller. In Table 1 we compare the scale of the datasets used in our work and related research. We discover that for low-resource MT, *parallel data is critical* not just during CPT, but even more so during SFT—in direct contrast with research on high-resource languages (HRLs).

Next, **diversity** in tasks, prompts, and datasets during SFT has been shown to significantly improve model performance across a range of tasks (Mishra et al., 2022; Chung et al., 2024). MT instructions have been shown to not just boost translation performance in unseen languages (Muennighoff et al., 2023), but also enhance LLM capabilities across diverse multilingual generation tasks (Ranaldi and Pucci, 2023; Zhu et al., 2023). Inspired by these findings, we study if SFT diversity could benefit low-resource LLM-MT systems too. By conducting experiments across a range of tasks and language pairs with SFT datasets of varying compositions, we establish that *diversity leads to negative interference* and fine-tuning on multilingual MT is the optimal strategy. Further, we observe that training for more epochs on MT data is more effective than curating and training on a diverse dataset of the same size.

Our contributions are thus as follows:

1. In contrast to findings for high-resource LLM-MT (Xu et al., 2024), we observe that for LRLs, LLMs benefit hugely from scale of parallel data, during both CPT and SFT stages

2. Linguistic and task diversity during SFT leads to negative interference for LRL LLM-MT, with focused multilingual MT fine-tuning for more epochs being the most effective recipe.

To ensure the generalizability of our findings, we conduct 2 sets of experiments training multilingual LLMs on different sets of languages: a) 11 Indigenous American and b) 4 North East Indian languages, wherein the former follows a Latin script and the latter includes languages that do not. Our focus in this work is on the eng/spa → X directions since generation in an LRL is known to be a much harder task than in an HRL like English or Spanish, and we are interested in studying the challenges involved. We experiment with 3 base LLMs of varying sizes—Gemma 2B (Gemma Team et al., 2024), Mistral 7B (Jiang et al., 2023), and Llama 3 8B (Dubey et al., 2024), and observe that findings are mostly consistent across these models.

By applying our findings to 2-stage training, our methods achieved a +16.5 average chrF++ improvement over few-shot prompting—with the largest gains coming from the 8 least-resourced American languages in our setup, all of which have about 10K-50K parallel sentences each. We hope that the findings of this work will be useful when scaling to LLMs that can effectively translate into lower-resource languages.

## 2 Related Work

**High-Resource Translation with LLMs** There has been considerable interest in using LLMs as MT systems recently. Following initial success in prompting LLMs for high-resourced pairs (Vilar et al., 2023; Garcia et al., 2023; Hendy et al., 2023; Zhang et al., 2023a; Iyer et al., 2023), subsequent works have attempted to train LLMs on parallel data at scale (Yang et al., 2023; Lu et al., 2024), but these yielded modest gains and underperformed smaller encoder-decoder baselines such as NLLB-200 (Costa-jussà et al., 2022). Zhang et al. (2024a) showed through scaling laws for SFT that LLMs pre-trained at the order of 50B-300B tokens saturate in MT performance with 20K-30K instructions. Concurrently, Xu et al. (2024) discovered excess parallel data washed out LLM knowledge, so they proposed a 2-stage paradigm called ALMA that involved pre-training on scaled-up monolingual data, followed by SFT on much smaller high-quality bitext (~60K lines). ALMA outperformed NLLB-200. Following their success, Alves et al. (2024) adopted the ALMA framework to train Tower 7B for 10 high-resourced languages, outperforming ALMA and also matching GPT-4.

It is worth noting that Xu et al. (2024) did not include parallel data during CPT. However, inspired by research showing LLMs unintentionally consume parallel data at scale (Briakou et al., 2023),

several works have included it at the order of several billions of tokens for the top 10-20 high-resource languages (Anil et al., 2023; Wei et al., 2023; Fujii et al., 2024; Alves et al., 2024). Concurrent work has explored pre-training on synthetic, code-switched parallel data for 101 languages (Lu et al., 2024) with a total of 400M sentences. This work explores the impact of parallel data exclusively for low-resource performance and experiments with 1M–13M parallel sentences (50M–750M tokens) during pre-training, two orders of magnitude smaller than prior work.

**Low-Resource Explorations in LLM-MT** LLMs have been shown to perform poorly in low-resource MT (Robinson et al., 2023; Zhu et al., 2024b). In response, there have been efforts to leverage external resources in the MT pipeline, including multilingual lexicons (Lu et al., 2023), rule-based linguistic tools (Zhang et al., 2024b), word alignments (Mao and Yu, 2024) and even entire grammar books (Reid et al., 2024). However, such approaches create dependencies on resources and hinder extensibility across languages. Instead, we focus on optimal data utilisation strategies during CPT and SFT, prioritizing extensibility.

**Cross-Lingual Instruction Tuning** There has been a body of work exploring multilingual instruction tuning that have touched on diversity and data scale, but most of it is limited to HRLs. MT examples were shown to improve cross-lingual generation (Ranaldi and Pucci, 2023; Zhu et al., 2023), while Chen et al. (2024) showed multilingual SFT on machine-translated Alpaca datasets matches or beats monolingual tuning. Kew et al. (2023) and Shaham et al. (2024) showed a small quantity of multilingual SFT data can improve cross-lingual generation capabilities in medium and high-resource languages, while Zhu et al. (2024a) showed only 32 examples in HRLs suffice to elicit MT capabilities from LLMs. In our work, we hypothesize that these findings have the common denominator of pre-training on HRLs at scale, and show that when moving to LRLs, these trends reverse and diversity is no longer beneficial.

## 3 Approach

We now describe our efforts to adapt the widely-used ALMA framework, originally designed for fine-tuning LLMs for HRL translation (Xu et al., 2024; Alves et al., 2024), for MT in LRLs.

### 3.1 Stage 1: Continued Pre-training (CPT)

**CPT on Monolingual Data** The objective of this stage is to 'teach' an LLM to model LRLs, which are scarce in the pre-training corpus. We conduct CPT on monolingual data with the standard Causal Language Modelling objective. We train with low-rank adaptation (LoRA) and attach rank 8 adapters to query and value matrices (Hu et al., 2022). We also fine-tune input and output embeddings.

**CPT on Parallel Data** In a scenario where monolingual data is scarce, it is crucial to investigate the most effective way to use parallel data—which, for our indigenous American languages, was found to surprisingly be more abundant than the former[1]. We investigate 3 methods of mixing all available parallel and monolingual data:

1. **All Mono**: Here, we merge monolingual data with only *the target side* of all available bitext—essentially using it as extra monolingual data

2. **Mono + parallel (concat)**: Here, we merge monolingual data with *concatenated source-target pairs* from parallel data. We prepend source and target language codes before concatenation and, following Guo et al. (2024), use a newline delimiter to separate them.

3. **Mono + parallel (separate)**: To ablate the impact of concatenation, we provide the source and target sides of parallel data as *separate* sentences, and shuffle with monolingual data.

We depict our approach in Figure 1. In the first technique, the motivation is that it might be hard for the LLM to learn to model concatenated sequences, given that they were likely scarce in the original pre-training corpus, with the added challenge of pre-training on 'new' low-resource languages. On the other hand, if the model is able to adapt to concatenated sequences, it could make the LLM more adjusted to the task of translation. Finally, the third method verifies whether the results of the 'concatenated' model are due to the concatenation itself, rather than simply being exposed to additional tokens in the source language. We use this terminology for all experiments in this work.

While these methods control how parallel data is incorporated in pre-training, we are also interested

---

[1]While monolingual online data is scarce, many translations of constitutions, articles etc. from Spanish do exist

Figure 1: Strategies explored for incorporating parallel data during Continued Pre-Training. We show a Spanish (es) to Aymara (aym) example from our parallel data.

in adjusting the ratio of parallel to monolingual data in the corpus, particularly for the *'concat'* method given pre-training on 100% parallel data can be suboptimal (Alves et al., 2024). To get a desired mixing ratio for *'concat'*, we include bitext only until it comprises a given percentage of the training corpus. Once this threshold is crossed, we use the *'all mono'* method to include it as monolingual data instead. Next, we use temperature sampling (Arivazhagan et al., 2019) to control the language-wise distribution in our monolingual and parallel pre-training data, since these are quite heterogeneous and certain languages are extremely low-resourced. We sample monolingual and parallel data independently if using the 'concat' method, else we just mix them all together and shuffle at the instance level to create our final pre-training corpus.

### 3.2 Stage 2: Supervised Fine-Tuning (SFT)

Next, we fine-tune with LoRA on supervised instruction data and detail the tasks explored below. Note that we convert all instructions to the standard Alpaca format, and compute loss on the target tokens only (Taori et al., 2023).

**Low-Resource MT**  Given our use case, the most intuitive task to include would be MT itself. The `instruction` for each example is chosen randomly from a set of translation prompts (Table 11), while the `input` and `output` fields are the source and target sentences respectively.

**High-Resource MT**  Apart from MT data in the LRLs, we also experiment with adding HRL MT data since it is more abundant and of higher quality, known to be important during SFT (Xu et al., 2024). To explore the impact of transfer learning, we work with HRL data that is in some way related to the source/target language, e.g. Spanish-English

data for experiments on Spanish-X. Instructions are formatted in the same way as the LRL data.

**General-Purpose Instruction Tuning**  Apart from MT data, we also explore adding widely used general-purpose instruction tuning datasets, such as Alpaca (Taori et al., 2023) and Aya (Singh et al., 2024), and use data from high-resource languages (comprising the source side in our tasks) to improve the model's overall instruction-following capabilities. However, for the most part, we are unable to find similar data in the LRLs we experiment on.

**Synthetic Cross-Lingual QA (XQA)**  We do not find any instruction tuning data for most LRLs, so we follow Iyer et al. (2024) to create synthetic Question Answering (QA) data. Starting from a parallel sentence pair $(X, Y)$, where $X$ is from an HRL (in our case, English/Spanish) and $Y$ is from an LRL, we prompt an LLM (Mixtral-8x7B-Instruct (Jiang et al., 2024) in this work) to generate a question $Q$ for which $X$ would be an answer. Since $X$ and $Y$ are semantically equivalent, $Y$ is treated as the answer to question $Q$. We add a requirement at the end of $Q$ to generate in the target language. Thus, we use $(Q, Y)$ as synthetic cross-lingual instruction data. We provide the templates used for generating XQA examples in Table 12.

## 4 Experiments and Discussions

We fine-tune 2 separate sets of multilingual LLMs for 2 different language groups to facilitate evaluation on test sets from 2 different low-resource MT shared tasks: AmericasNLP 2024 (Ebrahimi et al., 2024) and the Indic track of WMT 2023 (Pal et al., 2023). The former involves 11 Indigenous Central & South American languages, while the latter focuses on 4 North-East (NE) Indian languages. The first group includes Aymara (aym), Bribri (bzd), Asháninka (cni), Chatino (ctp), Guarani (grn), Huichol (hch), Nahuatl (nhe), Otomi (ote), Quechua (quy), Shipibo-Konibo (shp) and Tarahumara (tar). The second consists of Khasi (kha), Meitei (mni), Mizo (lus) and Assamese (asm). Our motivation in choosing these languages was to experiment with LRLs containing both Latin (American) and non-Latin (Indic) scripts; that also had widely used, high-quality test sets. We use the former for our main experiments and replicate the most interesting baselines in the Indic languages.

| Lang | #Tokens | Lang | #Tokens | Lang | #Tokens |
|------|---------|------|---------|------|---------|
| aym | 23.4M | bzd | 2.6M | cni | 2.2M |
| ctp | 5.4M | grn | 37.6M | hch | 3.5M |
| nhe | 32.1M | oto | 23.6M | quy | 45.1M |
| shp | 3.3M | tar | 2.3M | **Total** | **181.3M** |
| eng | 9.8M | spa | 27.9M | **Replay** | **37.6M** |

(a) Indigenous American Languages

| Lang | #Tokens | Lang | #Tokens | Lang | #Tokens |
|------|---------|------|---------|------|---------|
| asm | 1.1B | kha | 39.3M | lus | 165.1M |
| mni | 16.5M | eng | 7.3M | **Total** | **1.3B** |

(b) North-East Indian Languages

Table 2: Monolingual data statistics, with token counts calculated using the Llama3 8B tokenizer. English and Spanish are included as replay data. Note that LRL token counts overestimate data sizes, due to poor tokenization, and cannot be directly compared with HRLs.

## 4.1 Data

**Monolingual Data** Table 2 shows the token counts for the monolingual data collected for the 2 language groups. We note that the American languages are very low-resource, with 6 of 11 having 5M tokens or less. The Indic languages have relatively more data, with Assamese being medium-resourced but still likely low-resource in the original LLM pre-training corpus. Assamese and Meitei follow the Assamese-Bengali script, while Mizo and Khasi use the Latin script. Using the Llama3 tokenizer, we observe average fertilities of 2.87 and 3.83 for the American and Indic languages respectively, almost 3x that of high-resource languages, illustrating the under-representation of non-Latin scripts in SOTA LLMs. Finally, we include some data in English and Spanish as *replay data* to prevent catastrophic forgetting (Ibrahim et al., 2024).

**Parallel data** We curate parallel data from various sources, for use in both CPT and SFT. Tables 3 and 4 show the sizes for the American and Indian languages respectively. Given that our primary exploration is for the American languages with limited spa-X data, we also sample eng-X and por-X from OPUS (Tiedemann, 2012). Note the heavily skewed language distribution, with the 3 HRLs constituting 80% of spa-X data and 96% of the overall data. The lesser skew for spa-X is due to the efforts of AmericasNLP to collect data for these pairs and the prevalence of Spanish in Latin American countries. A similar skew also exists for the Indic languages, with English-Khasi being the least-resourced pair. We list sources for all curated

data, along with cleaning steps, in Appendix A.

## 4.2 Evaluation

To evaluate MT into the 11 American languages, we use AmericasNLP'23 validation sets (Ebrahimi et al., 2024) containing spa-X translation pairs. For the Indic pairs, we use the WMT 2023 test sets from the Indic track (Pal et al., 2023) which consist of eng-X pairs. Both evaluation datasets are multi-domain, as are the curated monolingual and parallel corpora. We show test set statistics in Table 5. Given the absence of neural metrics for these languages, we evaluate using ChrF++ (Popović, 2017), since both the American and Indic languages are morphologically rich wherein chrF++ is particularly effective (Popović, 2017). We use SacreBLEU (Post, 2018) for computing this. We also report confidence intervals with bootstrap resampling (Koehn, 2004), which we implement for the multilingual setting by computing the macro-average across all languages for each resample, and then computing the mean and variance across all resamples. We report the standard deviation as the confidence interval.

## 4.3 Experimental Settings

For temperature sampling our data, we use $\tau = 30$ for CPT and $\tau = 80$ for SFT. We used a batch size of 8 and gradient accumulation every 16 steps. We used a learning rate of 1e-4, with a cosine scheduler and a warmup ratio of 3%. We train all models on bf16 precision for 1 epoch. We use Llama-Factory (Zheng et al., 2024) for training and evaluating all models, with Deepspeed ZeRO3 (Rasley et al., 2020) for distributed training. For inference, we used a batch size of 16 with greedy decoding, since we found higher beam sizes yielded minimal gains.

## 4.4 Foundational Results

We first establish the importance of fine-tuning input and output embeddings in Table 6 and then show our foundational results for the American languages in Table 7, using 5-shot prompting as a baseline. Our findings are:

1. **Fine-tuning embeddings is critical.** Across the board, we observe that fine-tuning embeddings along with LoRA modules yields huge gains (Table 6), almost doubling chrF++ scores, indicating that this step is crucial to helping LLMs adapt to these new languages. Given this, we expect that full-weight fine-tuning would perform better, but stick to

| | Total | HRL | LRL | aym | bzd | cni | ctp | grn | hch | nhe | oto | quy | shp | tar |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| spa | 1M | 0.8M | 0.2M | 442K | 8K | 20K | 4K | 80K | 21K | 57K | 16K | 226K | 62K | 28K |
| por | 1.9M | 1.9M | 15K | 3.6K | 0 | 6.9K | 0 | 410K | 0 | 8K | 0 | 1520K | 0 | 0 |
| eng | 5.8M | 5.8M | 36K | 1053K | 0 | 13.9K | 0 | 2489K | 0 | 22K | 0 | 2271K | 0 | 0 |
| **Total** | **8.8M** | **8.5M** | **0.3M** | **1.5M** | **7.7K** | **41K** | **4.1K** | **3M** | **21K** | **87K** | **16K** | **4.0M** | **62K** | **28K** |

Table 3: Parallel data sentence counts for American languages, from source to each target language. HRL and LRL refer to the 3 high-resource languages (Aymara, Guarani, Quechua) and the other 8 low-resource ones respectively.

| Pair | eng-lus | eng-asm | eng-kha | eng-mni |
|---|---|---|---|---|
| **#Sents** | 6.5M | 5.0M | 25K | 443K |

Table 4: Indic Parallel Data Sizes (Sentence counts)

| Indigenous American | | | | | | NE Indic | |
|---|---|---|---|---|---|---|---|
| Pair | #Lines | Pair | #Lines | Pair | #Lines | Pair | #Lines |
| spa-aym | 996 | spa-nhe | 672 | spa-quy | 996 | eng-asm | 2000 |
| spa-bzd | 996 | spa-oto | 599 | spa-shp | 996 | eng-kha | 1000 |
| spa-cni | 883 | spa-gn | 995 | spa-tar | 995 | eng-mni | 1000 |
| spa-ctp | 499 | spa-hch | 994 | | | eng-lus | 2000 |

Table 5: Evaluation data statistics for the low-resourced American and Indic language experiments in this paper.

| Fine-Tuned Modules | Gemma 2B | Mistral 7B | Llama3 8B |
|---|---|---|---|
| LoRA only | 4.5 ± 0.1 | 8.8 ± 0.2 | 8.3 ± 0.2 |
| LoRA + embeddings | 9.6 ± 0.3 | 15.6 ± 0.4 | 15.6 ± 0.4 |

Table 6: Impact of fine-tuning input/output embeddings along with LoRA adapters is shown. Both models follow the "CPT *all mono*, SFT" recipe from Table 7.

| Method | Gemma 2B | Mistral 7B | Llama3 8B |
|---|---|---|---|
| 5-shot prompting | 2.8 ± 0.1 | 5.1 ± 0.1 | 3.9 ± 0.1 |
| SFT only | 8.7 ± 0.2 | 16.3 ± 0.4 | 14.8 ± 0.4 |
| CPT *all mono*, SFT | 9.6 ± 0.3 | 15.6 ± 0.4 | 15.6 ± 0.4 |
| CPT *mono+parallel*, SFT | **10.1 ± 0.3** | **16.7 ± 0.4** | **17.2 ± 0.4** |

Table 7: chrF++ scores for spa-X LLM-MT in the American languages, for LLMs of varying sizes. Both LoRA modules and embeddings are fine-tuned. For SFT, all models use 500K spa-X MT examples. Confidence estimates are computed using bootstrap resampling. *'mono+parallel'* uses concatenated bitext.

LoRA fine-tuning for cost-efficiency reasons. We fine-tune embeddings for all future results.

2. **Choosing larger base LLMs is crucial.** For under-represented (zero-resource) languages, is it more effective to train smaller LLMs with larger vocabularies (and thus, improved fertility), or vice versa?[2] We note that Gemma 2B has the largest vocabulary (256K tokens), followed by Llama3 (128K) and Mistral (32K tokens), resulting in improved fertility (2.36 vs 2.87 for Llama3/Mistral). Regardless, the larger models, Mistral 7B and Llama3 8B, vastly outperform the Gemma 2B model, suggesting fine-tuning smaller vocabulary LLMs like Mistral might be a better option from both cost and performance standpoints.

3. **SFT alone is effective, but CPT+SFT yields best results.** While SFT yields large gains over prompting, combining CPT and SFT seems optimal for both Gemma and Llama3, although the gap is smaller than that of high-resource MT (Xu et al., 2024), ostensibly due to the difference in scale. For Mistral, SFT alone seems to suffice—we hypothesize that

this might be due to the smaller vocabulary being effectively fine-tuned on SFT data alone.

4. **Pre-training on parallel data yields major gains consistently.** Lastly, we observe that pre-training on a mixture of *concatenated* monolingual and parallel data (*'mono + parallel'*) yields statistically significant gains over converting both as monolingual data (ie. *'all mono'*; refer Section 3.1). This trend is consistent for all 3 LLMs, with larger gains for the more effective models Mistral and Llama3. We note that for 'mono+parallel', we mix parallel and monolingual data in a 1:1 ratio, since we observed it worked best empirically, and show in Table 14 how increasing the ratio of parallel data during CPT monotonically improves performance. Given these gains, we explore the importance of bitext for low-resource LLM-MT further in Section 4.5.

## 4.5 Analysis: Importance of Parallel Data

**How important is concatenated parallel data at various scales of low-resource pre-training?**

---

[2]We found it prohibitively expensive to fine-tune Gemma 7B which has a larger vocabulary *and* a larger capacity.

Figure 2: Comparing Llama3 8B models pre-trained on monolingual data alone versus those included parallel data too—concatenated, or as separate texts at various scales. All models were pre-trained on 1M, 3M, 5M, 8M, and 13M sentences respectively, and markers denote the corresponding token counts. The y-axis shows chrF++ post SFT on 500K spa-X MT data for 1 epoch.

Having seen improvements in pre-training on the entire corpus (consisting of 13M 'mono+parallel' sentences, or 730M tokens) in Table 7, we now study the importance of parallel data as we scale *down*—an important consideration when one moves to even lower-resource settings. We pre-train on subsets of varying sizes and mix monolingual and parallel data in 3 ways: 'All mono', 'mono+parallel (concat)' and 'mono+parallel (separate)', as defined in Section 3.1). We fine-tune all these pre-trained models on the same SFT dataset: 500K spa-X MT instructions, and plot the resulting chrF++ scores in Figure 2, including error bars from bootstrap resampling. We note that 'all mono' has different markers than the others, as the token counts on including both source and target-side data in the corpus are obviously larger than only the latter. We find that a) starting at around 5M sentences (~300M tokens), it is *consistently advantageous* to include concatenated parallel data during pre-training. b) Given 'mono+parallel (separate)' severely underperforms, we establish that it is **concatenation** that adapts the LLM for the task of MT, not the extra data alone. Our findings complement those of Alves et al. (2024), who also observe gains from pre-training on parallel data in the 1B to 20B tokens range, using 10 high-resourced European languages. In contrast, our focus here is on investigating the *minimum data threshold* at which CPT leveraging parallel data becomes beneficial.

**How does scaling LRL parallel data during SFT impact performance?** We now turn our focus to scaling during SFT, which has not yielded gains in high-resource LLM-MT after 20K-30K sentences

(Zhang et al., 2024a). Motivated by our previous results, we re-examine this question for low-resource LLM-MT in Figure 3, wherein we evaluate at steps of 50K until ~1M SFT instructions, the point at which spa-X MT data runs out. Our findings are:

1. **Overall MT quality improves steadily with scale.** Unlike high-resource LLM-MT, we observe that, despite fluctuations, MT quality continues to grow until 1M sentences, particularly for LRLs, which have a steeper slope.

2. **The gains of CPT on parallel data carry over to SFT scaling, particularly for HRLs.** We observe that CPT on a mixture of concatenated monolingual and parallel data yields the largest initial gains, followed by CPT on monolingual-only data. We note that the former is particularly beneficial for HRLs, which have likely consumed a lot of parallel data during pre-training, and this helps them maintain a huge lead at lower SFT scales, suggesting that *concatenated parallel data teaches the LLM the task of translation and helps it adapt to MT more naturally during SFT*. Given the prevalence of bitext in LLM pre-training corpora (Briakou et al., 2023), the gains of concatenation would help explain why few-shot prompting (Garcia et al., 2023) and tiny-scale SFT (32 examples; c.f. Zhu et al. (2024a)) can elicit MT in the highest-resource languages!

3. **LRLs need SFT scaling much more than HRLs** So, while HRLs benefit hugely from CPT on bitext, the opposite is true for LRLs— the scarce amount of parallel data observed during CPT (see Table 3) is likely not enough to outperform 'all mono' pre-training. Instead, LLMs are far more responsive to the scale of LRL MT data during SFT, showing consistent performance improvement. Meanwhile, MT quality plateaus for the HRLs[3] (Figure 3b), suggesting that the 'less is more' (Zhou et al., 2023) trend popular in high-resource LLM-MT does *not* hold for LRLs, where *scale continues to remain the most effective option.*

For generalizability, we also provide scaling graphs for Mistral 7B in Figure 7 and report trends similar to Llama3. Moreover, we also compute

---

[3]although Aymara, Guarani and Quechua *are* LRLs for the original pre-training corpus of Mistral/Llama3

(a) Average (overall)     (b) Average (high-resource[†])     (c) Average (low-resource[η])

Figure 3: Scaling up Llama3 8B models with different CPT recipes (no CPT, CPT with monolingual data, and CPT with a mixture of monolingual and parallel data) on MT data for the American languages. [†] 'High-resource' refers to the relatively higher-resourced languages in our low-resource setup (Aymara, Guarani and Quechua) while the other 8 are grouped as low-resource[η].



Figure 4: Scaling up Llama 3 8B models with the 3 CPT recipes for the 4 Indic languages, until 5M sentences. We were forced to stop training 'No CPT' at 2.5M sentences, constrained by budget.



Figure 5: Epoch vs performance graph for low-resource LLM-MT. We use the entire 1M spa-X MT dataset, and plot average chrF++ for the Indigenous American languages, using Llama3 (Mono+Parallel (concat)) model.

scaling graphs for the 4 Indic languages in Figure 4 using the Llama3 model, until 5M eng-X sentences. Here, we observe that the gap between 'CPT: mono+parallel (concat)' and 'CPT (all mono)' is relatively lesser, and only significant until 500K sentences, at which point they start becoming comparable. This might be because the Indic parallel corpora used for CPT is relatively less diverse, consisting of only eng-X pairs in 4 languages, whereas our parallel data for the American languages is more heterogeneous (see Table 3)—which we observed to yield more gains in our preliminary experiments. Thus, we expect a more diverse corpus for the Indic languages to yield greater long-term gains, but we leave the verification of this to future work. In other respects, the trends are quite similar to the American languages: the gains of CPT carry over to SFT here too and interestingly, MT quality continues to improve until 5M SFT examples, once again in direct contrast to high-resourced LLM-MT research.

**More epochs: Scale through repetition** Finally, given the consistent gains from scale, we evaluate how effective data repetition could be during

SFT. We train for up to 10 epochs on the entire 1M spa-X MT dataset and plot our results in Figure 5. We note a strong monotonic gain of +3.3 average chrF++ until 5 epochs, with the largest coming from the 1st to the 2nd epoch (+2.0 chrF++). While the graph does plateau later (likely due to overfitting) this shows how simply training for more epochs can be an easy way to boost performance, given the data constraints for LRLs.

### 4.6 Analysis: Diversity in SFT is not *always* helpful

We now look at the importance of *diversity* in SFT, which plays a significant role in boosting LLM performance during general-purpose fine-tuning (Sanh et al., 2022; Longpre et al., 2023), but has also been favoured in training LLM-MT systems (Alves et al., 2024). Alves et al. (2024) observed improved performance on non-MT tasks by including non-MT instructions, but mixed results on MT quality—concluding that the "transfer/interference relations between tasks are complex". In this work, we try to tease apart these interference relations and study the effect of non-MT tasks on low-resource translation in depth. In particular, we look at 3 kinds of

| Method | All | HRL | LRL |
|---|---|---|---|
| Bilingual LLM-MT | 14.51 | 24.55 | 10.75 |
| Multilingual LLM-MT | 18.73 | 23.59 | 16.90 |

Table 8: Average chrF++ of SFT on a single language pair (bilingual) vs all language pairs (multilingual). As before, "HRL" includes Aymara, Guarani and Quechua, while the other eight are grouped under LRL.

diversity: *prompt*, *linguistic*, and *task* diversity, and examine their impact on MT quality when mixed in varying proportions in the SFT dataset.

**Bilingual or Multilingual fine-tuning?**   First, we ask whether given some multilingual low-resource data and a fixed FLOPS budget, would it be beneficial to do SFT on multiple, diverse low-resource MT pairs (multilingual LLM-MT) or only on a single pair (bilingual LLM-MT)? For this experiment, we simply concatenate all available data for the multilingual setting, while the bilingual setting consists of fine-tuning separate models for each pair, with the base LLM being constant in both cases: Llama3 (CPT with mono+parallel data). We show our results in Table 8. We observe that, on average, multilingual SFT outperforms the bilingual models by +4 chrF++ points. Looking closer, we notice that these gains mostly come from the lower-resourced languages (detailed results in Table 13). This is potentially because the bilingual setting has too little data for effective FT in these languages, while the multilingual option offers better transfer and scale. For higher-resourced pairs, the opposite is true: the performance is slightly lower, ostensibly due to negative interference. This is not unlike conventional multilingual NMT models, such as NLLB, wherein high-resource MT can often be worse than bilingual baselines.

**On Linguistic, Task and Prompt Diversity**   Inspired by the previous results showing target-side diversity in MT pairs boosts performance, we now broaden the scope of diversity during SFT. In Table 9a, we show varying mixtures of SFT data that all have the same size (500K examples) but are composed of different tasks. For *prompt diversity*, we ablate randomly sampling from a list of potential prompt templates (listed in Table 11) versus using a constant one (the first in Table 11). We observe a statistically significant gain in MT quality, similar to general-purpose SFT (Longpre et al., 2023).

Next, we study the more interesting question of *linguistic diversity*: can data in other MT pairs

| SFT Mixture | chrF++ |
|---|---|
| **Prompt Diversity** | |
| 500K spa-X MT (Same prompt) | 16.18±0.37 |
| 500K spa-X MT (Random prompts) | **17.22±0.40** |
| **Linguistic Diversity** | |
| 166K spa-X MT + 166K eng-X MT + 166K por-X MT | 14.26±0.37 |
| 250K spa-eng MT + 250K spa-X MT | 15.73±0.38 |
| 500K spa-X MT only | **17.22±0.40** |
| **Task Diversity** | |
| 250K spa-X MT + 250K XQA | 15.45±0.39 |
| 250K spa-X MT + 250K Aya (spa) | 15.68±0.39 |
| 250K spa-X MT (x 2 epochs) | 17.20±0.40 |
| 500K spa-X MT | **17.22±0.40** |

(a) Prompt, Linguistic and Task Diversity in American pairs

| SFT Mixture | chrF++ |
|---|---|
| 250K spa-X MT + 250K Aya (asm, mni) | 26.91 ± 0.42 |
| 250K spa-X MT + 250K Alpaca (eng) | 28.02 ± 0.41 |
| 500K spa-X MT only | **29.57 ± 0.43** |

(b) Task Diversity in Indic pairs

Table 9: Exploring interference due to diversity in SFT for our best Llama3 8B model (CPT on both monolingual and parallel data) on the American and Indic languages. The dataset size (example count) is prepended before each task. Scores shown are average chrF++.

transfer for test languages? In one baseline, we divide our SFT dataset equally into spa-X, eng-X, and por-X MT examples, into the American languages (statistics in Table 3). In another baseline, we use Spanish-English data from ParaCrawl (Bañón et al., 2020) and combine a 250K sentence subset with 250K pairs from our usual spa-X MT data. We find that this type of linguistic diversity leads to interference, and significantly underperforms the 500K spa-X MT baseline. This suggests that while *target diversity* in related languages might help performance, source diversity or unrelated high-resource languages like English may not.

Then, we look at *task diversity*: can non-MT tasks that elicit better instruction-following and general reasoning capabilities in the source or target language, benefit LLM-MT in LRLs? We mix Aya and XQA (Section 3.2) instructions with a 250K subset of spa-X MT examples. We also include an ablation that, in place of curating non-MT data, merely trains for 2 epochs on the 250K spa-X MT subset. Here we also discover that general-purpose tasks lead to interference and that training for 2 epochs on 250K MT examples is a more effective strategy comparable to 1 epoch on 500K.

We also look at the impact of task diversity on the Indic languages in Table 9b. Here, we are able

Figure 6: Mixing varying percentages of SFT tasks with spa-X MT, versus the impact on spa-X chrF++, for 500K examples. 'MT only' is a topline which *only* uses the spa-X MT data excluding all other tasks. Thus, it uses *lesser* data relatively and helps to estimate the interference of non-MT tasks on MT quality.

to find general-purpose instruction-tuning data in Assamese and Meitei, as mentioned in Section 3.2. Interestingly, including data in these *target languages*—the languages we want the model to be better at generating in—degrades performance the most! A similar amount of Alpaca data in English reduces performance far less, suggesting that reasoning in LRLs is such a hard task, that even data where the model *should* learn to generate in the languages of interest, is *worse* than generic English instruction-following data. Here, too, providing 500K spa-X examples is the most optimal strategy, underscoring the generalizability of our findings.

**Does non-MT SFT always cause interference?** Finally, we seek to establish conclusively if the negative results of diversity in Table 9a arise from factors like task mismatch or lack of quality; or if interference is simply *the norm* when fine-tuning on non-MT data for low-resourced LLMs. In Figure 6, we plot 5 graphs: the first 4 are graphs showing non-task-specific data (including XQA, Aya, {eng/por}-X MT and eng-spa MT) mixed in varying proportions with spa-X MT. Thus, for example, the second point of each plot represents 10% of that task mixed with 90% of spa-X MT, with the total always being 500K examples. We also include an ablation (dotted lines), which has the same quantity of spa-X MT data at each data point as other plots but without data from any other task listed.

We observe that *performance is always better using task-specific MT data*, and negative interference is indeed consistent. We conclude that transfer is challenging to achieve in low-resource settings, likely because at such scales, the LLM-MT model is still learning to generate translations, and reasoning is still a formidable challenge. Downstream

performance in such settings depends not on the diversity/composition of the SFT dataset, but only on the amount of LRL MT data provided—making optimizing for *quantity* the most effective strategy.

### 4.7 Discussion: LLMs vs NMT models

For the Indigenous American languages, the current SOTA systems are NordicAlps (Attieh et al., 2024) and DC-DMV (DeGenaro and Lupicki, 2024) that report average chrF++ scores of 26.73 and 23.76 respectively. Both are NMT models. NordicAlps trained a multilingual model from scratch on varying mixtures of eng-spa and spa-x data at scale, with their most major gains (+4 chrF++ points) coming from a novel redundancy-driven tokenization method. DeGenaro and Lupicki (2024) fine-tune NLLB-200 using a variety of parallel data sources, but unlike us, they generate a large amount of synthetic data. Our best-performing model, Llama3 with CPT on parallel data and 5 epochs of SFT, yields a score of 20.93 in the best setting. While this is a gain of +17.06 chrF++ from a 5-shot prompted Llama3, it is also almost 6 chrF++ points behind the SOTA. Since ours was not a shared task effort, we did not attempt explorations with tokenization or synthetic data, which are concurrent to our findings and would likely boost performance further. Overall, our results suggest that low-resourced LLM-MT systems, while promising, are still behind the curve compared to SOTA NMT models.

## 5 Conclusion

In this paper, we approach low-resource LLM-MT from a data-centric perspective and study performance along two axes: i) the size of parallel data and ii) diversity in SFT tasks. Through experiments and analyses, we conclude that quantity plays a dominant role in downstream performance. Specifically, parallel data drives performance significantly during both CPT and SFT, with HRLs and LRLs displaying different behaviours—making bitext a critical resource even in the modern era of LLM-MT. Moreover, we establish that diversity (on multiple fronts) consistently declines MT quality, with multilingual fine-tuning on task-specific data being the most effective option, reaffirming our previous findings on the value of scale. We hope these findings will be useful considerations when scaling to massively multilingual LLMs of the future.

## Limitations

One of the limitations of this work is that due to the unavailability of robust neural metrics for LRLs, we are forced to use string-based metrics like chrF++ as our primary evaluation metric. We note that while this is not optimal, it is also not unlike most other related works on low-resource translation, and chrF++ has been shown to align reasonably well with human assessment scores for morphologically rich languages. Secondly, we only focus on two specific language groups in this work. We expect that expanding the scope to a massively multilingual setting might yield even larger improvements, owing to scale and transfer learning.

## Acknowledgments

## References

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.

Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, et al. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges. *arXiv preprint arXiv:1907.05019*.

Joseph Attieh, Zachary Hopton, Yves Scherrer, and Tanja Samardzic. 2024. System description of the nordicsalps submission to the americasnlp 2024 machine translation shared task. In *Proceedings of the 4th Workshop on Natural Language Processing for Indigenous Languages of the Americas (Americas-NLP 2024)*, pages 150–158.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

Eleftheria Briakou, Colin Cherry, and George Foster. 2023. Searching for needles in a haystack: On the role of incidental bilingualism in PaLM's translation capability. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9432–9452, Toronto, Canada. Association for Computational Linguistics.

Pinzhen Chen, Shaoxiong Ji, Nikolay Bogoychev, Andrey Kutuzov, Barry Haddow, and Kenneth Heafield. 2024. Monolingual or multilingual instruction tuning: Which makes a better alpaca. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1347–1356, St. Julian's, Malta. Association for Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.

Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Ona De Gibert, Raúl Vázquez, Mikko Aulamo, Yves Scherrer, Sami Virpioja, and Jörg Tiedemann. 2023. Four approaches to low-resource multilingual NMT: The Helsinki submission to the AmericasNLP 2023 shared task. In *Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP)*, pages 177–191, Toronto, Canada. Association for Computational Linguistics.

Dan DeGenaro and Tom Lupicki. 2024. Experiments in mamba sequence modeling and nllb-200 fine-tuning

for low resource multilingual machine translation. In *Proceedings of the 4th Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP 2024)*, pages 188–194.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Abteen Ebrahimi, Ona de Gibert, Raul Vazquez, Rolando Coto-Solano, Pavel Denisov, Robert Pugh, Manuel Mager, Arturo Oncevay, Luis Chiruzzo, Katharina von der Wense, and Shruti Rijhwani. 2024. Findings of the AmericasNLP 2024 shared task on machine translation into indigenous languages. In *Proceedings of the 4th Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP 2024)*, pages 236–246, Mexico City, Mexico. Association for Computational Linguistics.

Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. 2024. Continual pre-training for cross-lingual llm adaptation: Enhancing japanese language capabilities. *arXiv preprint arXiv:2404.17790*.

Xavier Garcia, Yamini Bansal, Colin Cherry, George Foster, Maxim Krikun, Melvin Johnson, and Orhan Firat. 2023. The unreasonable effectiveness of few-shot learning for machine translation. In *International Conference on Machine Learning*, pages 10867–10878. PMLR.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

Jiaxin Guo, Hao Yang, Zongyao Li, Daimeng Wei, Hengchao Shang, and Xiaoyu Chen. 2024. A novel paradigm boosting translation capabilities of large language models. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 639–649, Mexico City, Mexico. Association for Computational Linguistics.

Kevin Heffernan, Onur Çelebi, and Holger Schwenk. 2022. Bitext mining using distilled sentence representations for low-resource languages. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2101–2112, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats Leon Richter, Quentin Gregory Anthony, Eugene Belilovsky, Timothée Lesort, and Irina Rish. 2024. Simple and scalable strategies to continually pre-train large language models. *Transactions on Machine Learning Research*.

Ayyoob ImaniGooghari, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora Kassner, Chunlan Ma, Helmut Schmid, André Martins, François Yvon, and Hinrich Schütze. 2023. Glot500: Scaling multilingual corpora and language models to 500 languages. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1082–1117, Toronto, Canada. Association for Computational Linguistics.

Vivek Iyer, Pinzhen Chen, and Alexandra Birch. 2023. Towards effective disambiguation for machine translation with large language models. In *Proceedings of the Eighth Conference on Machine Translation*, pages 482–495, Singapore. Association for Computational Linguistics.

Vivek Iyer, Bhavitvya Malik, Wenhao Zhu, Pavel Stepachev, Pinzhen Chen, Barry Haddow, and Alexandra Birch. 2024. Exploring very low-resource translation with LLMs: The University of Edinburgh's submission to AmericasNLP 2024 translation task. In *Proceedings of the 4th Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP 2024)*, pages 209–220, Mexico City, Mexico. Association for Computational Linguistics.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Tannon Kew, Florian Schottmann, and Rico Sennrich. 2023. Turning english-centric llms into polyglots: How much multilinguality is needed? *arXiv preprint arXiv:2312.12683*.

Mohammed Safi Ur Rahman Khan, Priyam Mehta, Ananth Sankar, Umashankar Kumaravelan, Sumanth Doddapaneni, Suriyaprasaad G, Varun Balan G, Sparsh Jain, Anoop Kunchukuttan, Pratyush Kumar,

Raj Dabre, and Mitesh M. Khapra. 2024. Indicllm-suite: A blueprint for creating pre-training and fine-tuning datasets for indian languages. *arXiv preprint arXiv: 2403.06350.*

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popovic, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. Preliminary wmt24 ranking of general mt systems and llms. *Preprint*, arXiv:2407.19884.

Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pages 388–395.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2024. Madlad-400: A multilingual and document-level large audited dataset. *Advances in Neural Information Processing Systems*, 36.

Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. The flan collection: Designing data and methods for effective instruction tuning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 22631–22648. PMLR.

Hongyuan Lu, Haoyang Huang, Dongdong Zhang, Haoran Yang, Wai Lam, and Furu Wei. 2023. Chain-of-dictionary prompting elicits translation in large language models. *arXiv preprint arXiv:2305.06575.*

Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. 2024. Llamax: Scaling linguistic horizons of llm by enhancing translation capabilities beyond 100 languages. *arXiv preprint arXiv:2407.05975.*

Zhuoyuan Mao and Yen Yu. 2024. Tuning LLMs with contrastive alignment instructions for machine translation in unseen, low-resource languages. In *Proceedings of the The Seventh Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2024)*, pages 1–25, Bangkok, Thailand. Association for Computational Linguistics.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao,

M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.

Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure, and Sandeep Kumar Dash. 2023. Findings of the wmt 2023 shared task on low-resource indic language translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 682–694.

Maja Popović. 2017. chrf++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Leonardo Ranaldi and Giulia Pucci. 2023. Does the English matter? elicit cross-lingual abilities of large language models. In *Proceedings of the 3rd Workshop on Multi-lingual Representation Learning (MRL)*, pages 173–183, Singapore. Association for Computational Linguistics.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '20, page 3505–3506, New York, NY, USA. Association for Computing Machinery.

Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530.*

Nathaniel Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. ChatGPT MT: Competitive for high- (but not low-) resource languages. In *Proceedings of the Eighth Conference on Machine Translation*, pages 392–418, Singapore. Association for Computational Linguistics.

Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti

Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*.

Uri Shaham, Jonathan Herzig, Roee Aharoni, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. Multilingual instruction tuning with just a pinch of multilinguality. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 2304–2317, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Shivalika Singh, Freddie Vargus, Daniel Dsouza, Börje F Karlsson, Abinaya Mahendiran, Wei-Yin Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura OMahony, et al. 2024. Aya dataset: An open-access collection for multilingual instruction tuning. *arXiv preprint arXiv:2402.06619*.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Jörg Tiedemann. 2009. News from opus-a collection of multilingual parallel corpora with tools and interfaces. In *Recent advances in natural language processing*, volume 5, pages 237–248.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in opus. In *Lrec*, volume 2012, pages 2214–2218. Citeseer.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for translation: Assessing strategies and performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15406–15427, Toronto, Canada. Association for Computational Linguistics.

Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li, Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhiwei Cao, Binbin Xie, et al. 2023. Polylm: An open source polyglot large language model. *arXiv preprint arXiv:2307.06018*.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024. A paradigm shift in machine translation: Boosting translation performance of large language models. In *The Twelfth International Conference on Learning Representations*.

Wen Yang, Chong Li, Jiajun Zhang, and Chengqing Zong. 2023. Bigtranslate: Augmenting large language models with multilingual translation capability over 100 languages. *arXiv preprint arXiv:2305.18098*.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023a. Prompting large language model for machine translation: A case study. In *International Conference on Machine Learning*, pages 41092–41110. PMLR.

Biao Zhang, Zhongtao Liu, Colin Cherry, and Orhan Firat. 2024a. When scaling meets LLM finetuning: The effect of data, model and finetuning method. In *The Twelfth International Conference on Learning Representations*.

Kexun Zhang, Yee Choi, Zhenqiao Song, Taiqi He, William Yang Wang, and Lei Li. 2024b. Hire a linguist!: Learning endangered languages in LLMs with in-context linguistic descriptions. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 15654–15669, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhengrui Ma, Yan Zhou, Langlin Huang, Mengyu Bu, Shangtong Gui, Yunji Chen, Xilin Chen, et al. 2023b. Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models. *arXiv preprint arXiv:2306.10968*.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, LILI YU, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. Lima: Less is more for alignment. In *Advances in Neural Information Processing Systems*, volume 36, pages 55006–55021. Curran Associates, Inc.

Dawei Zhu, Pinzhen Chen, Miaoran Zhang, Barry Haddow, Xiaoyu Shen, and Dietrich Klakow. 2024a. Fine-tuning large language models to translate: Will a touch of noisy data in misaligned languages suffice? *arXiv preprint arXiv:2404.14122*.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024b. Multilingual machine translation with large language models: Empirical results and analysis. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2765–2781, Mexico City, Mexico. Association for Computational Linguistics.

Wenhao Zhu, Yunzhe Lv, Qingxiu Dong, Fei Yuan, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Extrapolating large language models to non-english by aligning languages. *arXiv preprint arXiv:2308.04948*.

## A Data

### A.1 Monolingual Data

We provide the detailed monolingual data statistics for the indigenous American languages used for CPT in Table 10. We collect this data from various sources. For the indigenous American languages, we use MADLAD-400 (Kudugunta et al., 2024), GLOT 500 (ImaniGooghari et al., 2023), Wikipedia, data curated by the University of Helsinki (De Gibert et al., 2023) and OCR data collected by Iyer et al. (2024). The English and Spanish data used for replay comes from Wikipedia and MADLAD-400 too. We use MADLAD-400 and GLOT 500 again for the Indic languages again, along with the verified split of Sangraha (Khan et al., 2024)—a large corpus for Indian languages.

### A.2 Parallel Data

We provide the templates used for generating Machine Translation and XQA Instructions in Tables 11 and 12 respectively. For MT, we randomly use one of these prompts to create an instruction for doing SFT of our model. For XQA, the provided instructions are used to prompt Mixtral-8x-7B to generate a question, which is later used for creating the synthetic XQA datasets as described in Section 3.2.

We source the parallel data for generating these from the work of Iyer et al. (2024), which includes data sourced from AmericasNLP'23 official training datasets (Ebrahimi et al., 2024), NLLB-200 and FLORES-200 (Costa-jussà et al., 2022), OPUS (Tiedemann, 2009) and data curated by Helsinki in their 2023 submission (De Gibert et al., 2023).

#### A.2.1 Cleaning Parallel Data

**Rule-based Filtering** We clean parallel data by following standard filtering rules. First, we remove pairs with more digits or non-alphanumeric characters than alphabetic ones on either side. Then, we remove pairs where either the source or the target has less than 3 words or more than 120 words. To adjust for pairs containing partial or incomplete translations, we apply a combination of two filtering rules: a) for sentences with less than 25 words, the character length difference must be less than 65, and b) for those longer than 25 words, the character length ratio between either source and target, or vice versa, must be at most 1.55. Lastly, we filter our sentences with non-Latin characters in either sentence, as well as pairs with identical source and target sentences to create our final training dataset.

**Neural Quality Estimation** We also experimented with using neural quality estimation using models like LASER-3 (Heffernan et al., 2022). LASER-3 supports 3 of the 11 American languages—the HRLs (Aymara, Guarani and Quechua)—and we used it for scoring parallel sentences for these languages. We then sorted them in order of decreasing quality, calculated as the cosine similarity of LASER representations. For the other 8 LRLs, we used random scoring to simulate standard shuffling behaviour. We found that training on this 'sorted' data actually performed *worse* than our default baseline trained on unscored corpora, consistently at various SFT scales. As none of the authors spoke these languages, we could not substantively explain why, but we hypothesize that LASER representations might not be very reliable for these very LRLs, and might introduce certain undesirable biases in the sorted data, making the default baseline more robust to different kinds of bitext. In practice, we found that the standard rule-based filtering approaches worked best and thus, we stuck to them for cleaning our data.

Finally, we note that both our monolingual and parallel corpora span a variety of domains similar to our test data. Also, due to the paucity of data, we use all available sources.

## B Results

### B.1 Parallel vs Monolingual Data Ratio

In Table 14, we find that mixing higher ratios of parallel data with monolingual data either performs comparably or improves performance. We do not go higher than 50% since our parallel data runs out at this stage, and to ensure higher ratios we would have to oversample the existing dataset—which would not lead to a fair comparison with the other baselines. Nevertheless, given the monotonic trend, it would be interesting to explore if mixing higher ratios of parallel data continues to improve the performance even more. However, it is likely that there is some ceiling as to how much parallel data one should mix. For instance, Alves et al. (2024) show how a baseline trained on 100% parallel data underperforms compared to mixing it with monolingual data.

### B.2 Scaling up SFT: Mistral 7B

Figure 7 shows the scaling behaviour of the Mistral 7B model with different pre-training recipes (no

| Language | #Sentences | #Tokens | Language | #Sentences | #Tokens | Language | #Sentences | #Tokens |
|---|---|---|---|---|---|---|---|---|
| Aymara (aym) | 1M | 23.4M | Quechua (quy) | 1.9M | 45.1M | Guarani (grn) | 1.1M | 37.6M |
| Nahuatl (nhe) | 1.1M | 32.1M | Otomi (oto) | 0.6M | 23.6M | English (eng) | 0.4M | 9.8M |
| Spanish (spa) | 0.7M | 27.9M | Shipibo-Konibo (shp) | 0.1M | 3.3M | Bribri (bzd) | 0.1M | 2.6M |
| Asháninka (cni) | 0.2M | 2.2M | Chatino (ctp) | 0.3M | 5.4M | Huichol (hch) | 0.2M | 3.5M |
| Tarahumara (tar) | 0.2M | 2.3M | **Total (Replay)** | **1.1M** | **37.6M** | **Total** | **7.8M** | **218.9M** |

Table 10: Detailed monolingual data statistics for the American languages

## Translation Instructions

1. Translate the following sentence from {src_lang} to {tgt_lang}.
2. Can you convert the following sentence from {src_lang} to {tgt_lang}.
3. Kindly translate this sentence from {src_lang} into {tgt_lang}.
4. Could you translate the following from {src_lang} to {tgt_lang}?
5. Proceed to translate the subsequent sentence from {src_lang} to {tgt_lang}.
6. Change the following sentence from {src_lang} to {tgt_lang}.
7. Render the sentence below from {src_lang} into {tgt_lang}.
8. Switch the following sentence from {src_lang} into {tgt_lang} language.
9. Rephrase the following sentence into {tgt_lang} from {src_lang}.
10. Transform the following text from {src_lang} to {tgt_lang}.
11. Can you restate the following sentence from {src_lang} in {tgt_lang}?
12. Please provide a translation for this sentence from {src_lang} to {tgt_lang}.
13. Adapt the following into {tgt_lang} from the original {src_lang}.
14. Translate the subsequent text from {src_lang} into the {tgt_lang} language.

Table 11: MT Instruction Templates used during Supervised Fine-Tuning (SFT)

## XQA Instruction

"Consider this sentence: {input}\nWhat kind of specific instruction X could this be the unique answer to? Output ONLY the instruction, followed by a newline."

Table 12: Template used for generating XQA instructions

| Method | Avg | es-aym | es-bzd | es-cni | es-ctp | es-gn | es-hch | es-nhe | es-oto | es-quy | es-shp | es-tar |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bilingual LLM-MT | 14.5 | **21.2** | 6.1 | 12.3 | 8.4 | **27.9** | 15.0 | **14.1** | 9.5 | 24.6 | 14.5 | 6.1 |
| Multilingual LLM-MT | **18.7** | 20.2 | **15.8** | **18.0** | **26.0** | 25.8 | **21.6** | 18.2 | **11.2** | **24.9** | **14.8** | **9.5** |

Table 13: Comparison of Bilingual and Multilingual Es-X FT (MT only) Methods

CPT, CPT with monolingual data, and CPT with a mixture of monolingual and parallel data), where terminology is the same as that defined in Section 3.1. We observe that the trends for Mistral are largely similar to Llama3:

1. 'No CPT' underperforms the CPT baseline, but they become comparable at about ~500K sentences. As hypothesised in Section 4.4, due to a 4x larger vocabulary, Llama3 has 4x more parameters to fine-tune[4], meaning our

Llama3 models are likely more data hungry and retain the benefits of CPT over longer periods, but smaller models overfit sooner, leading to a shorter 'cross-over' threshold.

2. 'CPT (mono only)' consistently underperforms 'CPT (mono + parallel)', very similar to Llama3, lending further credence to our conclusion that concatenated parallel data adapts

---

[4]LoRA module parameters are negligible in comparison.

For fine-tuning input and output embeddings, Llama3 has $128K * 4096 * 2 \approx 1B$ parameters, whereas Mistral has $32K * 4096 * 2 \approx 250M$ parameters. LoRA parameters are on the order of 3M, for comparison.

| (a) Average (overall) | (b) Average (high-resource†) | (c) Average (low-resourceη) |

Figure 7: Scaling up Mistral 7B models with different CPT recipes (no CPT, CPT with monolingual data, and CPT with a mixture of monolingual and parallel data) on MT data for the American languages. We observe that the trends are very similar to Llama3, with the sole exception that the gains of CPT diminish faster at around 500K sentences. This might be because Mistral, due to a 4x smaller vocabulary, gets fine-tuned effectively with less data. † 'High-resource' refers to the relatively higher-resourced languages in our low-resource setup (Aymara, Guarani and Quechua) while the other 8 are grouped as low-resource$^\eta$.

| Ratio | Gemma 2B | Mistral 7B | Llama3 8B |
|-------|----------|------------|-----------|
| 0% | $9.6 \pm 0.3$ | $15.6 \pm 0.4$ | $15.6 \pm 0.4$ |
| 10% | $10.1 \pm 0.3$ | $15.9 \pm 0.4$ | $16.5 \pm 0.4$ |
| 25% | $10.1 \pm 0.3$ | $16.6 \pm 0.4$ | $16.6 \pm 0.4$ |
| 50% | $10.1 \pm 0.3$ | $16.7 \pm 0.4$ | $17.2 \pm 0.4$ |

Table 14: Model Performance vs. Parallel Data Ratio. While the gap between ratios is not always statistically significant, it is clear that the trend is monotonic and having 50% parallel data is consistently better than 0% (ie. the fully monolingual setting).

    the model to MT in a much better way

3. For HRLs, the performance plateaus quite quickly, while for LRLs SFT quality continues to grow with scale, once again following Llama3's trend

Our experiments with Mistral thus help in providing more evidence to support our claims with regard to the benefit of parallel data for low-resource LLM-MT systems.

### B.3 Bilingual vs Multilingual LLM-MT: Detailed Results

We provide detailed, language-specific results for bilingual vs multilingual LLM-MT in Table 13. We observe that multilingual LLM-MT mostly outperforms bilingual baselines with the exception of 3 relatively higher-resourced languages (Aymara, Nahuatl and Guarani) where there is a bit of a gap.

# Efficient Technical Term Translation: A Knowledge Distillation Approach for Parenthetical Terminology Translation

**Jiyoon Myung**[*], **Jihyeon Park**[*], **Jungki Son**[*], **Kyungro Lee**[*], **Joohyung Han**[*]

PrompTart LAB, MODULABS

{jiyoon0424, milhaud1201, aeolian83, lkr981147, ddang8jh}@gmail.com

[*] All authors contributed equally to this work.

## Abstract

This paper addresses the challenge of accurately translating technical terms, which are crucial for clear communication in specialized fields. We introduce the Parenthetical Terminology Translation (PTT) task, designed to mitigate potential inaccuracies by displaying the original term in parentheses alongside its translation. To implement this approach, we generated a representative PTT dataset using a collaborative approach with large language models and applied knowledge distillation to fine-tune traditional Neural Machine Translation (NMT) models and small-sized Large Language Models (sLMs). Additionally, we developed a novel evaluation metric to assess both overall translation accuracy and the correct parenthetical presentation of terms. Our findings indicate that sLMs did not consistently outperform NMT models, with fine-tuning proving more effective than few-shot prompting, particularly in models with continued pre-training in the target language. These insights contribute to the advancement of more reliable terminology translation methodologies.

## 1 Introduction

Terminology translation task is essential for understanding documents rich in technical terms, such as academic papers and technical reports. Traditionally, methods in the task have involved identifying term pairs in the source and target languages and using these pairs for training or post-editing purposes. However, challenges arise when there is no precise match for a term in the target language, or when new terms are used inconsistently. For instance, the term "fine-tuning" may be variably translated as "파인튜닝" or "미세조정" in Korean.

To address this, our research proposes a novel approach called Parenthetical Terminology Translation (PTT), which displays the original term in parentheses alongside its translation. This approach aims to mitigate reader confusion, especially when suitable translations are unavailable or translation accuracy is low. Although similar translation strategies using parenthetical form have been suggested in previous studies, effective technical solutions for this approach remain underexplored.

With the advent of advanced Large Language Models (LLMs), researchers have started exploring their potential for various tasks, including translation. LLMs can effectively support PTT through simple prompt usage, offering a promising solution for this approach. However, the practical application of LLMs is hindered by their high computational costs and latency, making them less feasible for real-time or large-scale deployment.

To mitigate these limitations, this study focuses on achieving the capabilities of LLMs using smaller, traditional Neural Machine Translation (NMT) models and small-sized Language Models (sLMs). We generated a high-quality PTT dataset using LLMs and distilled the knowledge by fine-tuning these smaller models with this dataset. This approach ensures that the benefits of LLMs can be harnessed without incurring high computational costs. Additionally, we evaluated the performance of various models and training methods to optimize model performance and efficiency, particularly for the specialized PTT task.

Our proposed task extends beyond mere translation accuracy; it also emphasizes the correct presentation of technical terms within parentheses, which is crucial for enhancing reader comprehension. To quantitatively evaluate this aspect, we introduced a novel metric specifically designed to assess the models' ability to accurately and effectively use parenthetical annotations. This metric not only evaluates translation quality but also ensures that technical terms are correctly presented, allowing for a robust comparison of model performance across different architectures and training techniques.

Thus, this paper makes three significant contributions to the field of terminology translation:

1410

1. **Synthetic Data Generation**: We propose a collaborative framework using Large Language Models (LLMs) to generate well-curated datasets specifically for the English-Korean Parenthetical Terminology Translation (PTT) task. This framework employs multiple agents to create high-quality sentence pairs, enabling smaller models to perform the PTT task with high accuracy. By leveraging robust data from LLMs, the framework ensures consistency and precision, making it effective for handling domain-specific terminology.

2. **Knowledge Distillation and Model Comparison**: Utilizing these high-quality datasets, we fine-tuned various Neural Machine Translation (NMT) models and small-sized Language Models (sLMs). We then conduct a comprehensive performance analysis from diverse perspectives, highlighting the strengths and limitations of each model. This analysis provides valuable insights for future research and development in the field.

3. **New Evaluation Metric**: We introduce a novel evaluation metric that quantitatively assesses the ability of models to present appropriate terms within parentheses. This metric ensures contextual accuracy and reader comprehension, offering a robust framework for evaluating model performance in the context of PTT.

These contributions aim to advance the domain of terminology translation by providing practical and efficient solutions that leverage the strengths of both large and small language models. Our approach addresses the inherent challenges of the terminology translation task and paves the way for more accessible translation methodologies in technical and specialized fields.

## 2 Related Work

Terminology translation plays a crucial role in ensuring consistency and accuracy in specialized domains like technical and academic documentation. Early approaches, such as rule-based and statistical machine translation, effectively leveraged predefined glossaries and translation memories (Melby et al., 1999). While these methods successfully maintain consistency within certain contexts, they often struggle with out-of-domain (OOD) words

and ambiguous terms (Och and Ney, 2003). Moreover, these approaches are less effective when dealing with domain-specific or emerging terms not covered by existing resources (Tiedemann, 2010; Tiedemann and Scherrer, 2017).

To maintain clarity and precision in academic and technical documentation, it is often necessary to preserve certain terms from the source language. This practice is particularly valuable in cases where the translated term may be unfamiliar to the reader or where retaining the original term is essential for legal or scientific accuracy (Moghadam and Far, 2015; Hasler et al., 2018; Michon et al., 2020). A further strategy to support this practice involves the strategic use of parentheses, where textual additions can help enhance translation quality and consistency through corpus-based improvements (Lin et al., 2008; Huang et al., 2017; Hawamdeh, 2018). Despite its potential benefits, the systematic implementation of this approach remains relatively underexplored in current research.

Recent studies have highlighted the effectiveness of knowledge distillation in transferring knowledge from large language models (LLMs) to smaller traditional translation models (Li et al., 2024; Enis and Hopkins, 2024). Through this process, datasets generated by a powerful teacher model are distilled into a student model, enabling small-sized models to perform specialized tasks like terminology translation. The multi-agent framework is particularly effective in generating targeted, high-quality datasets for specific tasks (Wu et al., 2023). Within this framework, different agents are assigned specialized roles, such as data generation and evaluation, collectively enhancing the quality and relevance of the resulting dataset. This collaborative process is essential for precise and context-aware data generation, which is crucial for training models to excel in specialized translation tasks.

## 3 Data Generation

To create a high-quality dataset for the Parenthetical Terminology Translation (PTT) task, we employed four collaborative agents—Writer, Translator, Evaluator, and Executor—utilizing the large language models (LLMs) GPT-4o-mini or GPT-4-turbo for each agent. Our goal was to generate English sentences containing technical terminologies alongside their Korean translations, with the original English terms included within parentheses. The overall framework is illustrated in Figure 1.

Figure 1: Multi-Agent Framework for generating a high-quality PTT dataset using four agents.

## 3.1 Writer

The Writer agent was responsible for creating academic English sentences that included the technical terms. To achieve this, we first compiled a comprehensive list of terms to be included in the dataset. Recognizing the rapid emergence of new terminology in the field of artificial intelligence (AI), we focused on terms frequently encountered in AI-related research. To ensure multiple terms could be incorporated into single sentences, we clustered similar domain-specific terms together.

Next, we utilized the arXiv API to find papers that contained all or some of the terms from each cluster. By including the summary of the most relevant paper in the prompt, we helped the Writer LLM understand the appropriate contexts in which these terms were used. This ensured that the generated sentences were contextually accurate and meaningful.

To enhance data diversity, the Writer were tasked with generating sentences where each term appeared either once or in conjunction with other terms. By combining these sentences post-generation, we facilitated the creation of sentences with various characteristics: sentences where terms appear only once, sentences containing different terms together, and sentences where the same term appears more than twice. This diversity allowed us to analyze the performance of PTT from multiple perspectives, ensuring a comprehensive evaluation of the models under different conditions. The complete prompt used for the Writer agent is provided in the Appendix (see Listing 1).

## 3.2 Translator

The Translator agent translated the English sentences into Korean, ensuring that each target term was followed by its original English term in parentheses to fulfill the PTT task requirements. To enhance accuracy, we employed the GPT-4 Turbo model, while other agents utilized GPT-4o-mini. Additionally, we applied one-shot prompting by providing a relevant example to guide the translation process. This approach helped maintain consistency and precision in the PTT task, ensuring that technical terms were accurately presented within parentheses. The complete prompt used for the Translator agent is provided in the Appendix (see Listing 2).

## 3.3 Evaluator/Executor

The Evaluator agent reviewed the translated sentences, scoring them from 0 to 10 based on the accuracy of term usage and overall translation quality. Then, The Executor agent transit the statement 'If the score is less than 8: Response "translator". If the score is 8 or greater: Response "final output".' If a sentence scored below 8, the Evaluator suggested corrections, prompting the Translator to revise the translation. The Translator would then repeat the translation task until the Executor rated the sentence as '"final output"', ensuring the highest quality and consistency in the dataset. The complete prompt used for the Evaluator agent is provided in the Appendix (see Listing 3).

After the automatic data generation process, human reviewers conducted a final quality check to

| Terms Set Index | 0 |
|---|---|
| **Domain** | AI (in-domain) |
| **Terms** | adversarial training, recurrent neural architectures, bayesian optimization |
| **Source** | By implementing adversarial training, researchers have demonstrated significant improvements in the performance of recurrent neural architectures against adversarial attacks. The integration of recurrent neural architectures with Bayesian optimization enhances the model's ability to learn from limited data while minimizing computational resources. |
| **Target** | 적대적 훈련(adversarial training) 도입으로 연구자들은 적대적 공격에 대응하는 순환 신경 구조(recurrent neural architectures)의 성능 향상을 크게 입증하였습니다. 순환 신경 구조(recurrent neural architectures)와 베이지안 최적화(bayesian optimization)를 결합함으로써 모델이 제한된 데이터로부터 학습하는 능력을 향상시키고, 계산적 자원 사용을 최소화합니다. |

Table 1: Sample of generated data from the AI domain (in-domain). Each entry includes the term set index, domain, source text, and the corresponding target translation. Red text highlights the targeted terms $T_{\text{Eng}}$, while blue text indicates the correct representation of terms $T_{\text{Kor}}$ in the Korean translation.

ensure the dataset's reliability. Following this review, we combined seven sentences from each cluster into three composite sentences as mentioned earlier, resulting in 1,398 English-Korean paired sentences encompassing 233 term clusters (a total of 699 distinct terms). We split these 1,398 sentences into 1,116 for training, 144 for validation, and 138 for testing the performance. We carefully ensured that sentences containing the same terms were allocated to the same dataset, maintaining consistency and preventing data leakage across the splits. The sample data is provided in Table 1. The entire dataset is available on Hugging Face at `https://huggingface.co/datasets/PrompTart/PTT_en_ko`.

## 3.4 Out-of-Domain Evaluation Dataset

To evaluate the generalization ability of models in the PTT task, we generated additional datasets in domains beyond artificial intelligence (AI), specifically targeting biology and physics. The data generation process followed the same methodology as the in-domain training dataset to ensure consistency. In total, we generated 171 paired sentences for biology (subcellular processes), 60 for nanoscale physics, and 168 for high-energy physics. Each domain-specific dataset was curated by referencing relevant academic papers, providing authentic and contextually accurate examples. These out-of-domain datasets allowed for a comprehensive assessment of the models' robustness and adaptability across different specialized fields.

## 4 Knowledge Distillation

In this study, we applied knowledge distillation to fine-tune both traditional neural machine translation (NMT) models and small-sized large language models (sLMs) using the synthetic Parenthetical Terminology Translation (PTT) dataset generated in the previous step. Our goal was to evaluate the effectiveness of distillation techniques across various model architectures, sizes, and training methodologies, offering insights into how distilled models perform in specialized translation tasks.

### 4.1 Fine-tuning Traditional Machine Translation Models

To evaluate the performance of knowledge distillation on traditional neural machine translation models, we employed several widely used opensource models. We focused on encoder-decoder Transformer-based models that support Korean. Specifically, we tested the following models:

- mBART50 (Liu et al., 2020) : This multilingual NMT model is pre-trained on monolingual corpora from 50 languages and is finetuned for translation tasks. It consists of 611 million parameters.

- M2M100 (Fan et al., 2020): A large-scale multilingual NMT model trained on 2200 translation directions, enabling many-to-many translation across 100 languages. We tested the base version with 418 million parameters.

- NLLB-200 (Koishekenov et al., 2023): Known for its extensive language coverage,

this model is particularly useful for low-resource languages and inclusive translation services. We tested the distilled version with 600 million parameters.

The fine-tuning parameters were provided in the Appendix (see Table 5).

### 4.2 Fine-tuning small-sized Large Language Models

To effectively compare performance with NMT models, we also fine-tuned open-source small-sized large language models. Our goal was to evaluate various models from multiple perspectives to gain comprehensive insights into the PTT task. To ensure reproducibility and a broad evaluation, we selected four well-known and high-performing open-source LLMs:

- Llama 3 (Touvron et al., 2023): The latest iteration of the Llama series, this model further refines the architecture introduced in earlier versions, enhancing its performance on large-scale datasets. We evaluated the 8B and 70B versions in our experiments.

- Gemma 2 (Team et al., 2024): A next-generation multilingual model, Gemma 2 is designed to deliver high performance across diverse natural language tasks with an emphasis on efficiency. We assess the model by testing three versions: the smallest (2B), a mid-sized variant (9B), and a larger configuration (27B).

- Qwen 2 (Yang et al., 2024): An updated version of the Qwen series, Qwen 2 is developed with a strong focus on flexibility and adaptability to domain-specific tasks. It offers improved performance and efficiency, particularly in handling complex language modeling challenges. In this study, we analyzed the 1.5B, 7B, and 72B versions.

- Mistral (Jiang et al., 2023): Mistral is known for its streamlined design and high efficiency in multilingual tasks. We specifically evaluate the 7B version to examine how its architecture balances performance with computational efficiency.

To compare pre-trained models with those that have been further fine-tuned specifically for the Korean language, we also tested models that underwent continual pre-training (Ke et al., 2023) in Korean. This approach allowed us to assess the impact of additional language-specific pre-training on the models' performance in the Parenthetical Terminology Translation (PTT) task.

- beomi/Llama-3-Open-Ko-8B[1]: A specialized version of Llama 3 focused on Korean language tasks. This open-source model is fine-tuned to excel in Korean linguistic applications.

- beomi/Llama-3-KoEn-8B[2]: A bilingual version of Llama 3 tailored for both Korean and English language tasks. This model is designed to maintain balanced performance across both languages, making it versatile for multilingual applications.

Furthermore, we explored instruction-tuned versions of the aforementioned models using different training techniques, such as LoRA (Hu et al., 2021) and QLoRA (Dettmers et al., 2023). In addition, we applied few-shot prompting (Brown et al., 2020) to both the instruction-tuned models and a commercial LLM (GPT-4o) to compare the effects of fine-tuning versus prompting. This comprehensive evaluation provides valuable insights into how knowledge distillation, combined with various tuning and prompting strategies, can enhance translation accuracy while maintaining efficiency across diverse model architectures.

The hyper-parameters for fine-tuning and LoRA are detailed in the Appendix (see Table 5), while the full prompt used for few-shot prompting is identical to the Translator agent's prompt (Listing 2), with the exception that we did not provide a list of terms in this case.

## 5 Custom Evaluation Metric

This section introduces a novel metric designed specifically for the Parenthetical Terminology Translation (PTT) task, aimed at evaluating not only the accuracy of overall translation but also the correct presentation of the technical terms within parentheses.

For each sentence in the dataset, let $T_{\text{Eng}}$ represent the list of all technical terms provided in the original English sentence, including duplicates

---

[1]https://huggingface.co/beomi/Llama-3-Open-Ko-8B

[2]https://huggingface.co/beomi/Llama-3-KoEn-8B

if the same term appears multiple times. Similarly, let $T_{\text{Kor}}$ represent the list of those terms that are correctly translated into Korean and accompanied by their original English terms in parentheses. We define $|T_{\text{Eng}}|$ as the total number of technical terms in the English sentence (including duplicates), and $|T_{\text{Kor}}|$ as the number of correctly translated terms from $T_{\text{Eng}}$ in the Korean sentence. The ratio of these terms is calculated as the weight $W_{\text{terms}} = \min\left(\frac{|T_{\text{Kor}}|}{|T_{\text{Eng}}|}, 1\right)$. This ratio is capped at 1 to ensure that no penalty is applied if more terms appear correctly in the Korean translation than in the original English sentence. The adjusted metric for the PTT task, $M_{\text{PTT}}$, is then computed by multiplying this clipped ratio with the original translation metric $M$, such that $M_{\text{PTT}} = W_{\text{terms}} \times M$. Finally, we average $M_{\text{PTT}}$ across all sentence pairs in the dataset to obtain the final evaluation metric.

We employed BLEU (Papineni et al., 2002), COMET (Rei et al., 2020), and BERTScore (Zhang et al., 2020) as $M$ to evaluate pure translation performance. The translation metrics are computed after removing the parenthetical terms, ensuring that we assess only the translation's accuracy and fluency. This approach allows us to maintain a focus on both the translation's quality and the correct handling of technical terms within parentheses.

## 6 Evaluation

### 6.1 Quantitative Analysis

The results presented in Table 2 provide a comprehensive overview of the quantitative performance of various models and training techniques on the in-domain Parenthetical Terminology Translation (PTT) dataset, while Table 3 presents results on the out-of-domain dataset. Key observations are summarized as follows:

1. **sLMs vs. NMT Models**: The performance comparison between small-sized Large Language Models (sLMs) and Neural Machine Translation (NMT) models reveals that sLMs do not consistently outperform NMT models, even though LLMs are often perceived as more advanced due to their architecture. For instance, mBART50 and M2M100, achieved weighted BLEU scores of 37.52 and 40.05, respectively, with corresponding weight indicate strong PTT performance. These scores were comparable or superior to those achieved by some sLMs, such as the Llama 3 8B and

70B models, which obtained similar weighted BLEU scores but required significantly larger model sizes.

2. **Instruction-Tuned vs. Base Models**: Within the same sLM families, base models generally slightly outperformed instruction-tuned models on the PTT task. For instance, the Llama 3 8B model with QLoRA achieved a weighted BLEU score of 38.88, while the instruction-tuned version (8B-it) with the same QLoRA technique scored slightly lower at 37.84. This trend suggests that instruction-tuned models, which are trained on a broad range of tasks, may not gain a specific advantage for the specialized requirements of the PTT task.

3. **Fine-Tuning vs. Prompt Engineering**: Applying prompt engineering instead of fine-tuning to instruction-tuned models, using a 1-shot prompting approach, resulted in very poor performance. For example, the Llama 3 8B-it scored only 0.523, and the Gemma 2-9B-it scored 0.342 on weight metric. Even the commercial LLM GPT-4o performed worse than other fine-tuned small models, underscoring the critical importance of fine-tuning for specialized tasks like PTT.

4. **Korean Continued Pre-trained Models**: Models that underwent continued pre-training in the target language (Korean) generally outperformed others, with the Llama-3-KoEn-8B-it achieving the highest score among all models. Although Llama-3-Open-Ko-8B, which was continued pre-trained exclusively in Korean, showed slightly lower performance with a weighted BLEU score of 39.869, it still performed well. This highlights the importance of bilingual proficiency in models for the PTT task, where handling both source and target languages effectively is crucial for success.

5. **Model Size and Out-of-Domain Performance**: In the in-domain dataset, model size had little impact on performance, with smaller models like Gemma 2 7B even outperforming larger ones like Gemma 2 27B. However, when tested on out-of-domain datasets, all models experienced significant performance drops, but larger models such as Gemma 2 27B or Llama 3 70B showed less decline, indicating better generalization capabilities.

| Model | #Params | Training Techniques | $W_{\text{terms}}$ | $M_{\text{PTT}}$ (BLEU) | $M_{\text{PTT}}$ (COMET) | $M_{\text{PTT}}$ (BERT) |
|---|---|---|---|---|---|---|
| **Open-source NMT systems** | | | | | | |
| mBART50 | 611M | Full Fine-Tuning | 0.931 | 37.519 | 0.831 | 0.863 |
| M2M100 | 418M | Full Fine-Tuning | 0.958 | 40.048 | 0.855 | 0.889 |
| NLLB-200 | 600M | Full Fine-Tuning | 0.685 | 24.544 | 0.606 | 0.630 |
| **Llama family sLMs** | | | | | | |
| Llama3 | 8B | LoRA | 0.959 | 37.632 | 0.856 | 0.887 |
| | 8B | QLoRA | 0.949 | 38.875 | 0.847 | 0.880 |
| | 70B | LoRA | 0.957 | 38.869 | 0.855 | 0.888 |
| Llama3-Instruct | 8B-it | QLoRA | 0.954 | 37.840 | 0.851 | 0.881 |
| | 8B-it | 1-shot prompting | 0.523 | 0.577 | 0.214 | 0.310 |
| **Gemma family sLMs** | | | | | | |
| Gemma2 | 2B | LoRA | 0.946 | 37.959 | 0.842 | 0.875 |
| | 9B | LoRA | 0.958 | 41.567 | 0.858 | 0.893 |
| | 9B | QLoRA | 0.935 | 38.955 | 0.835 | 0.869 |
| | 27B | LoRA | 0.966 | 40.856 | 0.865 | 0.899 |
| Gemma2-Instruct | 9B-it | QLoRA | 0.953 | 39.215 | 0.849 | 0.884 |
| | 9B-it | 1-shot prompting | 0.342 | 9.698 | 0.286 | 0.286 |
| **Qwen family sLMs** | | | | | | |
| Qwen2 | 1.5B | LoRA | 0.950 | 34.374 | 0.838 | 0.868 |
| | 7B | LoRA | 0.945 | 39.167 | 0.844 | 0.877 |
| | 7B | QLoRA | 0.951 | 38.014 | 0.846 | 0.879 |
| | 72B | LoRA | 0.956 | 40.837 | 0.855 | 0.889 |
| Qwen2-Instruct | 7B-it | QLoRA | 0.947 | 37.990 | 0.842 | 0.874 |
| **Mistral family sLMs** | | | | | | |
| Mistral | 7B | QLoRA | 0.931 | 37.646 | 0.830 | 0.862 |
| Mistral-Instruct | 7B-it | QLoRA | 0.927 | 37.990 | 0.826 | 0.857 |
| **Korean Continued Pre-trained sLMs** | | | | | | |
| Llama-3-KoEn | 8B-it | QLoRA | **0.974** | **41.789** | **0.873** | **0.907** |
| | 8B-it | 1-shot prompting | 0.614 | 0.333 | 0.080 | 0.110 |
| Llama-3-Open-Ko | 8B-it | QLoRA | 0.953 | 39.869 | 0.852 | 0.885 |
| **Commercial LLM** | | | | | | |
| GPT-4o | Unknown | 0-shot prompting | 0.616 | 20.596 | 0.547 | 0.564 |
| GPT-4o | Unknown | 1-shot prompting | 0.751 | 26.509 | 0.669 | 0.689 |

Table 2: Model performance metrics for in-domain test data. $W_{\text{terms}}$ represents the average ratio of correctly translated terms with original English terms in parentheses. $M_{\text{PTT}}$ (BLEU), $M_{\text{PTT}}$ (COMET), and $M_{\text{PTT}}$ (BERT) are the original tranlsation metrics adjusted using $W_{\text{terms}}$ and averaged over all data. The suffix '-it' indicates instruct-tuned models. The top scores for each metric are highlighted in bold.

This suggests that while smaller models can be highly effective in specialized tasks, larger models are more versatile and better suited for handling diverse and unfamiliar datasets. The larger models' ability to retain higher performance levels in out-of-domain tasks underscores their capacity to adapt to a wider range of terminologies and contexts, making them more versatile in applications where data variability is a key challenge.

## 6.2 Qualitative Analysis

1. **Progression of PTT and Translation Skills**: As illustrated in Table 4 , most of the models, including M2M100, initially demonstrated strong proficiency in the PTT task, particularly in incorporating original terms within parentheses, as indicated by the high weight metrics in the earlier epochs. Over successive training epochs, the model's overall translation skills improved gradually, leading to better performance across all weighted metrics. A detailed illustration of these improvements is provided in the Appendix (see Table 6).

2. **Challenges with Less Common Terms**: Our analysis highlights a persistent challenge among models in accurately translating less common terms, especially proper nouns. As demonstrated in Table 7, terms like "de Finetti's theorem" were inconsistently translated across different models, reflecting the difficulty these models face when dealing with less familiar terminology. This inconsistency

| Model | #Params | Training Techniques | $W_{\text{terms}}$ | $M_{\text{PTT}}$ (BLEU) | $M_{\text{PTT}}$ (COMET) | $M_{\text{PTT}}$ (BERT) |
|---|---|---|---|---|---|---|
| **Open-source NMT systems** | | | | | | |
| mBART50 | 611M | Full Fine-Tuning | 0.784 | 19.538 | 0.668 | 0.698 |
| M2M100 | 418M | Full Fine-Tuning | 0.763 | 20.472 | 0.650 | 0.680 |
| NLLB-200 | 600M | Full Fine-Tuning | 0.160 | 4.066 | 0.138 | 0.144 |
| **Llama family sLMs** | | | | | | |
| Llama3 | 8B | LoRA | 0.769 | 22.595 | 0.670 | 0.693 |
| | 8B | QLoRA | 0.849 | 23.792 | 0.736 | 0.760 |
| | 70B | LoRA | 0.854 | 33.321 | 0.762 | 0.788 |
| Llama3-Instruct | 8B-it | QLoRA | 0.864 | 23.498 | 0.750 | 0.775 |
| **Gemma family sLMs** | | | | | | |
| Gemma2 | 2B | LoRA | 0.824 | 24.203 | 0.711 | 0.736 |
| | 9B | LoRA | 0.886 | 38.639 | 0.793 | 0.822 |
| | 9B | QLoRA | **0.900** | 32.914 | 0.799 | 0.826 |
| | 27B | LoRA | 0.897 | **40.379** | **0.804** | **0.836** |
| Gemma2-Instruct | 9B-it | QLoRA | 0.883 | 34.861 | 0.785 | 0.813 |
| **Qwen family sLMs** | | | | | | |
| Qwen2 | 1.5B | LoRA | 0.762 | 9.831 | 0.598 | 0.635 |
| | 7B | LoRA | 0.750 | 18.531 | 0.637 | 0.662 |
| | 7B | QLoRA | 0.849 | 20.028 | 0.729 | 0.751 |
| | 72B | LoRA | 0.877 | 32.048 | 0.779 | 0.806 |
| Qwen2-Instruct | 7B-it | QLoRA | 0.864 | 23.498 | 0.750 | 0.775 |
| **Mistral family sLMs** | | | | | | |
| Mistral | 7B | QLoRA | 0.870 | 15.942 | 0.713 | 0.749 |
| Mistral-Instruct | 7B-it | QLoRA | 0.876 | 17.350 | 0.717 | 0.756 |
| **Korean Continued Pre-trained sLMs** | | | | | | |
| Llama-3-KoEn | 8B-it | QLoRA | 0.884 | 35.492 | 0.789 | 0.817 |
| Llama-3-Open-Ko | 8B-it | QLoRA | 0.887 | 35.409 | 0.790 | 0.813 |

Table 3: Model performance metrics for out-of-domain test data. $W_{\text{terms}}$ represents the average ratio of correctly translated terms with original English terms in parentheses. $M_{\text{PTT}}$ (BLEU), $M_{\text{PTT}}$ (COMET), and $M_{\text{PTT}}$ (BERT) are the original tranlsation metrics adjusted using $W_{\text{terms}}$ and averaged over all data. The suffix '-it' indicates instruct-tuned models. The top scores for each metric are highlighted in bold.

underscores the importance of the PTT task, which helps maintain translation accuracy by preserving original terms alongside their translations, thereby reducing the likelihood of incorrect interpretations.

3. **Out-of-Domain Translation Challenges**: Most models struggled with translating out-of-domain (OOD) sentences, as detailed in Table 8. They often failed to accurately translate OOD terms, frequently substituting them with unrelated or incorrectly adapted words, sometimes even drawing from other languages. These frequent mistranslations highlight the need for more robust training methods or supplementary mechanisms to improve the models' generalization ability for handling unseen datasets effectively.

## 7 Conclusion

In this study, we explored the Parenthetical Terminology Translation (PTT) task, a specialized translation problem that focuses on mitigating potential inaccuracies in term translation by displaying the original technical term in parentheses alongside its translation. To effectively evaluate this approach, we introduced a novel evaluation metric, $M_{\text{PTT}}$, designed to measure both the accuracy of overall translation and the proper parenthetical presentation, ensuring that technical terms are effectively communicated across languages.

To generate a high-quality dataset for this task, we utilized a collaborative approach involving Writer, Translator, Evaluator, and Executor agents, supported by large language models (GPT-4). This allowed us to create a diverse and contextually accurate dataset that reflects real-world usage of technical terms in artificial intelligence (AI), biology, and physics. We then applied knowledge distillation techniques to fine-tune both traditional Neural Machine Translation (NMT) models and small-sized Large Language Models (sLMs), comparing their performance across various model architectures, sizes, and training methods.

| Epoch | Weight | Weighted BLEU | Weighted COMET | Weighted BERT Score |
|-------|--------|---------------|----------------|---------------------|
| epoch 1 | 0.939 | 31.780 | 0.824 | 0.858 |
| epoch 3 | 0.924 | 35.668 | 0.821 | 0.853 |
| epoch 5 | 0.948 | 37.853 | 0.844 | 0.878 |
| epoch 7 | 0.956 | 38.685 | 0.851 | 0.886 |
| epoch 9 | 0.958 | 40.048 | 0.855 | 0.889 |

Table 4: Model performance metrics of the M2M100 model for test data across training epochs.

Our findings revealed that sLMs did not consistently outperform NMT models, challenging the assumption that more advanced architectures inherently lead to superior performance. Additionally, within the same sLM families, base models slightly outperformed instruction-tuned models, suggesting that broad task training may not offer advantages for specialized tasks like PTT. Fine-tuning proved crucial, as prompt engineering approaches like 1-shot prompting resulted in significantly poorer performance. Moreover, models with continued pretraining in Korean outperformed others, highlighting the importance of bilingual proficiency for the PTT task. While model size had little impact on in-domain performance, larger models demonstrated better generalization on out-of-domain datasets, suggesting they are more versatile and better suited for handling diverse and unfamiliar data. These insights contribute to optimizing models and training techniques for specialized translation tasks, offering practical guidance for future research and applications in terminology translation.

## Limitations

**Penalty Mechanism in Evaluation Metrics**: The current approach to evaluating PTT performance involves simply multiplying translation metrics by a weight that reflects the presence of correctly parenthesized terms. However, this straightforward multiplication can disproportionately affect the overall performance scores. A more sophisticated penalty mechanism, such as using an exponential function, could provide a more balanced assessment by reducing the impact on the metric scores. Additionally, the current metric does not penalize the model for excessively parenthesizing trivial or unintended terms, which could lead to over-parenthesization. Future work could incorporate penalties for such cases, potentially by introducing concepts of recall and precision to refine the evaluation.

**Potential Bias in the Dataset**: The PTT dataset was generated using GPT-4, and the performance metrics were assessed with this dataset as the ground truth. This approach may introduce biases inherent to the GPT-4 model into the dataset, potentially affecting the robustness and generalizability of the models trained on it. To mitigate this, future research should consider generating datasets using a variety of models, ensuring a broader representation of translation styles and reducing the potential for model-specific biases.

**Language Scope of the Study**: This study focused exclusively on translation into Korean, which limits the generalizability of the findings across different languages. PTT performance might vary significantly with other languages due to differences in linguistic structures and translation challenges. Expanding the study to include translations into multiple languages would enable a more comprehensive analysis of the PTT task and provide insights into how the models perform across different linguistic contexts.

## Acknowledgements

## References

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms.

Maxim Enis and Mark Hopkins. 2024. From llm to nmt: Advancing low-resource machine translation with claude.

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation.

Eva Hasler, Adrià de Gispert, Gonzalo Iglesias, and Bill Byrne. 2018. Neural machine translation decoding with terminology constraints. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 506–512, New Orleans, Louisiana. Association for Computational Linguistics.

Mohammad Amin Hawamdeh. 2018. Explicitation by textual addition in parentheses in translating the quranic text into english. International Journal of Applied Linguistics and English Literature.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.

Guoping Huang, Jiajun Zhang, Yu Zhou, and Chengqing Zong. 2017. Learning from parenthetical sentences for term translation in machine translation. In Proceedings of the 9th SIGHAN Workshop on Chinese Language Processing, pages 37–45, Taiwan. Association for Computational Linguistics.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. 2023. Continual pre-training of language models.

Yeskendir Koishekenov, Alexandre Berard, and Vassilina Nikoulina. 2023. Memory-efficient NLLB-200: Language-specific expert pruning of a massively multilingual machine translation model. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3567–3585, Toronto, Canada. Association for Computational Linguistics.

Jiahuan Li, Shanbo Cheng, Shujian Huang, and Jiajun Chen. 2024. Mt-patcher: Selective and extendable knowledge distillation from large language models for machine translation.

Dekang Lin, Shaojun Zhao, Benjamin Van Durme, and Marius Paşca. 2008. Mining parenthetical translations from the web by word alignment. In Proceedings of ACL-08: HLT, pages 994–1002, Columbus, Ohio. Association for Computational Linguistics.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. Transactions of the Association for Computational Linguistics, 8:726–742.

Alan K. Melby, Brigham Young, Sue Ellen Wright, and Kentucky State. 1999. Leveraging terminological data for use in conjunction with lexicographical resources 1 integration.

Elise Michon, Josep Maria Crego, and Jean Senellart. 2020. Integrating domain terminology into neural machine translation. In International Conference on Computational Linguistics.

Masoumeh Moghadam and Mansureh Far. 2015. Translation of technical terms: A case of law terms. Journal of Language Teaching and Research, 6:830.

Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. Computational Linguistics, 29(1):19–51.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 2685–2702, Online. Association for Computational Linguistics.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska,

Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. 2024. Gemma 2: Improving open language models at a practical size.

Jörg Tiedemann. 2010. Context adaptation in statistical machine translation using models with exponentially decaying cache. In Proceedings of the 2010 Workshop on Domain Adaptation for Natural Language Processing, pages 8–15, Uppsala, Sweden. Association for Computational Linguistics.

Jörg Tiedemann and Yves Scherrer. 2017. Neural machine translation with extended context. In Proceedings of the Third Workshop on Discourse in Machine Translation, pages 82–92, Copenhagen, Denmark. Association for Computational Linguistics.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. Qwen2 technical report.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert.

1420

## A Appendix

|  | Parameter | NMT | sLM (w/ LoRA) | sLM (w/ QLoRA) |
|---|---|---|---|---|
| Training Argument | Learning Rate | 3e-5 | 1e-4 | 2e-4 |
|  | Lr Scheduler Type | linear | cosine | cosine |
|  | Optimizer | AdamW | paged_adamw_32bit | paged_adamw_32bit |
|  | Warmup Ratio | N/A | 0.03 | 0.03 |
|  | Weight Decay | 0.01 | 0.001 | N/A |
|  | Max Grad Norm | 1.0 | 1.0 | 0.3 |
|  | Dtype | bfloat16 | bfloat16 | bfloat16 |
| LoRA Configure | LoRA R | N/A | 64 | 64 |
|  | LoRA Alpha | N/A | 16 | 16 |
|  | LoRA Dropout | N/A | 0.1 | 0.1 |
|  | Target Modules | N/A | all-linear | all-linear |

Table 5: Hyper-parameters for Fine-Tuning and LoRA Techinque

```
You are a professional paper writer.

[TERM1] = {terms[0]}
[TERM2] = {terms[1]}
[TERM3] = {terms[2]}

<reference>
{arxiv_summaries}
</reference>

<instruction>
- The request is to thoroughly review and cite the provided <reference> when writing
    theacademic paper.
- Write complex English sentences using the given technical terms.
- Use appropriate academic tone.
- Each sentence MUST be clear, accurate, and contextually appropriate for a
    scientific paper.
- Generate only in English.
</instruction>

## Output Format:
1.english: A sentence using terms [TERM1].
2.english: A sentence using terms [TERM2].
3.english: A sentence using terms [TERM3].
4.english: A sentence using terms [TERM1] and [TERM2].
5.english: A sentence using terms [TERM2] and [TERM3].
6.english: A sentence using terms [TERM1] and [TERM3].
7.english: A sentence using terms [TERM1], [TERM2], and [TERM3].

CAUTION: Ensure that exactly 7 sentences are generated.
```

Listing 1: Full Prompt of Writer

```
You are a professor specializing in AI, proficient in both Korean and English.

[TERM1] = {terms[0]}
[TERM2] = {terms[1]}
[TERM3] = {terms[2]}

<translation guideline>
- Translate while preserving the original term like 사전 훈련(pre-train).
- If there is an abbreviation, translate it like this Korean term(english term,
    abbreviation).
- Identify terms, acronyms, and concepts to keep in English.
- Maintain academic tone and technical accuracy in your translations.
- Ensure the translation is natural in Korean while accurately conveying the
    original meaning.
- Change all the letters within the parentheses in Korean sentences to lowercase.
- IMPORTANT: The terms corresponding to [TERM1], [TERM2], and [TERM3] MUST ALWAYS be
     enclosed in parentheses like this: Korean term(English term).
</translation guideline>

<example>
english: LLMs demonstrate new abilities such as in-context learning, instruction
    following, and multi-step reasoning, enabling them to learn new tasks, follow
    instructions, and effectively solve complex problems.
korean: LLM은 맥락 학습(in-context learning), 지시 사항 따르기(instruction following), 다단계 추론
(multi-step reasoning)과 같은 새로운 능력을 보여줌으로써 새로운 작업을 학습하고, 지시를 따르며,
복잡한 문제를 효과적으로 해결할 수 있습니다.
</example>

## Output Format:
1.korean: [Korean translation]
2.korean: [Korean translation]
...
( Continue this pattern for all 7 sentences )
```

Listing 2: Full Prompt of Translator

```
You're an expert evaluating English to Korean translations of research papers, with
    a specific focus on proper parenthetical translations of technical terms.

<criteria>
- The format for parenthetical translations should be: Korean term(English term).
- The specific terms {terms[0]}, {terms[1]} or {terms[2]} MUST ALWAYS be enclosed in
    parentheses in the Korean translation.
- Parentheses should be properly placed, ensuring consistency in parenthesizing
    across the entire sentence.
- Ensures the translation conveys the original meaning precisely and reads naturally
    and smoothly.
</criteria>

<instruction>
- Change all the letters within the parentheses in Korean sentences to lowercase.
- Evaluate the Korean translation of the provided English sentences.
- Check the consistency and correctness of parenthesization.
- Provide a score (0-10) based on the correctness and consistency of
    parenthesization as Korean term(English term).
- Offer specific improvement suggestions if the score is less than 10.
- DO NOT include any supplementary explanations.
- Check your output format again.
</instruction>

## Example Output:
english: The neural network uses backpropagation to optimize its weights.
korean: 신경망(neural network)은 역전파(backpropagation)를 사용하여 가중치(weight)를 최적화합니다.
score: 10/10
terms_check: [neural network: Yes, backpropagation: Yes, weight: Yes]
parentheses_count: 3
suggestions: No improvements needed / Suggest ensuring that "model compression" is
translated as "모델 압축(model compression)" and adjusting "모델 컴프레션" to "model compression"
for consistency and clarity.

## Example Format:
1.
english: [English text using term "{terms[0]}"]
korean: [Korean translation using parentheses]
score: [X/10]
terms_check: [{terms[0]}: Yes/No, {terms[1]}: Yes/No, {terms[2]}: Yes/No]
parentheses_count: [Number of parentheses pairs in the Korean translation]
suggestions: [Suggest capturing the original meaning and nuances in the translation
    while adjusting the structure for natural flow and grammar]
2.
english: [English text using terms "{terms[0]}" and "{terms[1]}"]
korean: [Korean translation]
3.
...
(Continue this pattern for all 7 sentences)
```

Listing 3: Full Prompt of Evaluator

| Domain | AI (in-domain) |
|---|---|
| **Terms** | heterogeneous graphs, neural probabilistic text segmentation |
| **Source** | Heterogeneous graphs provide a rich framework for the implementation of neural probabilistic text segmentation, facilitating improved contextual understanding through varied data representations. |
| **Target** | 이종 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 개선된 맥락적 이해를 촉진합니다. |

| Epoch | Model Output |
|---|---|
| epoch 1 | ~~해태로잰~~ 이종 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현에 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 맥락적 이해를 향상시킵니다. |
| epoch 2 | ~~동적~~ 이종 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 맥락 이해를 향상시킵니다. |
| epoch 3 | 이종 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 문맥 이해를 향상시킵니다. |
| epoch 4 | ~~동성애~~ 이종 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 개선된 맥락 이해를 촉진합니다. |
| epoch 5 | ~~동적~~ 이종 그래프(heterogeneous graphs)는 다양한 데이터 표현을 통해 개선된 맥락 이해를 촉진하여 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공합니다. |
| epoch 6 | ~~동성애~~ 이종 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 문맥 이해를 개선하는 데 도움을 줍니다. |
| epoch 7 | 이질적인 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 개선된 맥락 이해를 촉진합니다. |
| epoch 8 | 이질적인 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 개선된 맥락 이해를 촉진합니다. |
| epoch 9 | 이질적인 그래프(heterogeneous graphs)는 신경 확률적 텍스트 분할(neural probabilistic text segmentation)의 구현을 위한 풍부한 프레임워크를 제공하며, 다양한 데이터 표현을 통해 개선된 맥락 이해를 촉진합니다. |

Table 6: Illustration of M2M100 output for in-domain sample data across epochs. Red strikethrough text indicates incorrect translations, with blue text showing the corrected translations. Green text represents correct translations that were maintained or introduced during later epochs.

| Domain | AI (in-domain) |
|---|---|
| **Terms** | neural task-driven modeling, de Finetti's theorem |
| **Source** | Neural task-driven modeling, when applied to the constructs of de Finetti's theorem, unveils a sophisticated approach to managing uncertainty in predictive models within artificial intelligence. |
| **Target** | 신경 작업 중심 모델링(neural task-driven modeling)이 드 피네티의 정리(de finetti's theorem)의 구조에 적용될 때, 인공 지능 내에서 예측 모델의 불확실성을 관리하는 정교한 접을 드러냅니다. |

| Model | Model Output |
|---|---|
| mBART50 | 신경 작업 주도 모델링(neural task-driven modeling)이 데 페네터 이론의 구성 요소에 적용될 때, 인공지능 내 예측 모델에서의 불확실성을 관리하는 정교한 접근 방식을 밝혀냅니다. |
| M2M100 | 신경 작업 주도 모델링(neural task-driven modeling)이 디 피네티(de Finetti) 이론의 구조에 적용될 때, 인공 지능 내 예측 모델의 불확실성을 관리하는 정교한 접근 방식을 드러냅니다. |
| NLLB-200 | 데 피네티 정리의 적용은 신경 작업 주도 모델링(neural task-driven modeling)과 함께 사용함으로써 불확실성을 효과적으로 관리할 수 있는 더 탄력적인 모델을 가능하게 합니다. |
| Llama3-8B | 신경 작업 주도 모델링(neural task-driven modeling)을 데 피니의 정리의(de Finetti's theorem) 구축에 적용하면 인공 지능 내 예측 모델의 불확실성을 관리하는 정교한 접근 방식을 드러낼 수 있습니다. |
| Llama3-8B-Instruct | 신경 작업 주도 모델링(neural task-driven modeling)이 데 피니의 정리(de finetti's theorem)의 구축에 적용될 때, 인공지능 내에서 예측 모델의 불확실성을 관리하는 정교한 접근 방식을 드러냅니다. |
| Gemma2-9B | 신경 작업 주도 모델링(neural task-driven modeling)이 디 핀에티의 정리(de finetti's theorem)의 구조에 적용될 때, 인공지능 내 예측 모델에서 불확실성을 관리하는 정교한 접근 방식을 밝혀냅니다. |
| Gemma2-9B-Instruct | 신경 과제 주도 모델링(neural task-driven modeling)을 데 핀에티의 정리(de finetti's theorem)의 구조에 적용하면 인공지능 내 예측 모델에서 불확실성을 관리하는 정교한 접근 방식을 드러냅니다. |
| Qwen2-7B | 디-finetti의 정리(de finetti's theorem)의 구조에 적용된 신경 작업 주도 모델링(neural task-driven modeling)은 인공 지능에서 예측 모델의 불확실성을 관리하는 정교한 접근 방식을 밝혀냅니다. |
| Qwen2-7B-Instruct | 디 베니테이의 정리(de finetti's theorem)의 구조에 신경 작업 주도 모델링(neural task-driven modeling)을 적용하면 인공 지능에서 예측 모델의 불확실성을 관리하는 정교한 접근 방식을 밝혀낼 수 있습니다. |
| Mistrial-7B | 신경 작업 주도 모델링(neural task-driven modeling)이 디 피네티(de finetti)의 정리의 구조에 적용될 때, 인공지능 내에서 예측 모델에서 불확실성을 관리하는 정교한 접근 방식을 밝혀냅니다. |
| Mistrial-7B-Instruct | 신경 작업 주도 모델링(neural task-driven modeling)을 데 핀철리의 정리(de finetti's theorem)의 구조에 적용할 때, 인공 지능 내에서 예측 모델의 불확실성을 관리하는 정교한 접근 방식이 밝혀집니다. |
| Llama-3-Open-Ko-8B | 신경 과제 기반 모델링(neural task-driven modeling)이 데 핀에티의 정리(de finetti's theorem)의 구성 요소에 적용될 때, 인공 지능 내 예측 모델의 불확실성을 관리하는 정교한 접근 방식을 드러냅니다. |
| Llama-3-KoEn-Instruct | 신경 작업 주도 모델링(neural task-driven modeling)을 데 피니티의 정리(de finetti's theorem)의 구성에 적용할 때, 인공 지능 내에서 예측 모델의 불확실성을 관리하는 정교한 접근 방식을 드러냅니다. |

Table 7: Model output comparison for in-domain test data across various models, with sLMs fine-tuned using QLoRA. Blue text highlights the translations of specific terms that models frequently translate inconsistently, reflecting challenges in handling proper nouns.

| Domain | High-Energy Physics (out-of-domain) |
|---|---|
| **Terms** | jet quenching, quarkonia suppression |
| **Source** | The interplay between jet quenching and quarkonia suppression offers a comprehensive understanding of the dynamics governing the quark-gluon plasma and its evolution. |
| **Target** | 제트 퀀칭(jet quenching)과 쿼크오니아 억제(quarkonia suppression) 간의 상호작용은 쿼크-글루온 플라즈마의 동역학과 그 진화에 대한 포괄적인 이해를 제공합니다. |
| **Model** | **Model Output** |
| mBART50 | ject quenching(jet quenching)과 quarkonia suppression(quarkonia suppression) 간의 상호 작용은 쿼크- 글루언 플라즈마와 그 진화에 대한 종합적인 이해를 제공합니다. |
| M2M100 | jet quenching과 quarkonia suppression 간의 상호 작용은 quark-gluon plasma의 역학과 그 진화에 대한 포괄적인 이해를 제공합니다. |
| NLLB-200 | 제트 화와 쿼크니아 억제 간의 상호 작용은 쿼크-글루온 플라즈마와 그 진화에 대한 동적에 대한 포괄적인 이해를 제공합니다. |
| Llama 3-8B | 제트 .jet 퀀ン チ잉(jet quenching)과 쿼크오니아 억제(quarkonia suppression) 간의 상호 작용은 쿼크-글루온 플라즈마와 그 진화가 지배하는 역학에 대한 종합적인 이해를 제공합니다. |
| Gemma2-2B | 제트 격리(jet quenching)와 쿼크니아 억제(quarkonia suppression) 간의 상호 작용은 쿼크-글루온 플라즈마와 그 변화에 관한 역학을 종합적으로 이해하는 데 도움을 줍니다. |
| Gemma2-9B | 제트 콸닝(jet quenching)과 쿼크니아 억제(quarkonia suppression) 간의 상호 작용은 쿼크-글루온 플라즈마와 그 진화를 지배하는 역학에 대한 포괄적인 이해를 제공합니다. |
| Gemma2-27B | 제트 냉각(jet quenching)과 쿼크오니아 억제(quarkonia suppression) 간의 상호 작용은 쿼크-글루온 플라즈마와 그 진화를 지배하는 역학에 대한 종합적인 이해를 제공합니다. |
| Qwen2-1.5B | _jet quenching와_quetsquon suppression 간의 상호작용은 광자론 합성(quarkonia suppression)의 동적을 종합적으로 이해하는 데 도움을 줍니다. |
| Qwen2-7B | 점프 퀀터링(jet quenching)과 퀸코나임(quarkonia suppression) 간의 상호 작용은 퀸크-글루온 플라즈마와 그 진화를 지배하는 역학을 포괄적으로 이해하는 데 중요한 역할을 합니다. |
| Qwen2-72B | 제트 쿠enching(jet quenching)과 쿼크니아 억제(quarkonia suppression) 간의 상호 작용은 쿼크-글루온 플라스마와 그 진화를 지배하는 역학에 대한 종합적인 이해를 제공합니다. |

Table 8: Model output comparison for out-of-domain data across various models, with sLMs fine-tuned using LoRA. Red text highlights specific terms that are frequently mistranslated, indicating challenges in handling these out-of-domain terms.

| Model | # Params | Training Technique | BLEU | COMET | BERT |
|---|---|---|---|---|---|
| mBART50 | 611M | Full Fine-Tuning | 40.298 | 0.892 | 0.927 |
| M2M100 | 418M | Full Fine-Tuning | 41.789 | 0.892 | 0.928 |
| NLLB-200 | 600M | Full Fine-Tuning | 35.843 | 0.886 | 0.920 |
| Llama3 | 8B | LoRA | 39.243 | 0.892 | 0.924 |
| | 8B | QLoRA | 40.984 | 0.893 | 0.927 |
| | 70B | LoRA | 40.600 | 0.893 | 0.928 |
| Llama3-Instruct | 8B-it | QLoRA | 39.686 | 0.892 | 0.924 |
| | 8B-it | 1-shot prompting | 1.103 | 0.410 | 0.594 |
| Gemma2 | 2B | LoRA | 40.126 | 0.890 | 0.925 |
| | 9B | LoRA | 43.391 | 0.896 | 0.932 |
| | 9B | QLoRA | 41.620 | 0.893 | 0.929 |
| | 27B | LoRA | 42.313 | 0.896 | 0.931 |
| Gemma2-Instruct | 9B-it | QLoRA | 41.143 | 0.891 | 0.928 |
| | 9B-it | 1-shot prompting | 28.314 | 0.837 | 0.838 |
| Qwen2 | 1.5B | LoRA | 36.174 | 0.881 | 0.914 |
| | 7B | LoRA | 41.434 | 0.893 | 0.927 |
| | 7B | QLoRA | 39.975 | 0.890 | 0.924 |
| | 72B | LoRA | 42.704 | 0.894 | 0.929 |
| Qwen2-Instruct | 7B-it | QLoRA | 40.107 | 0.889 | 0.923 |
| Mistral | 7B | QLoRA | 40.424 | 0.891 | 0.925 |
| Mistral-Instruct | 7B-it | QLoRA | 39.368 | 0.891 | 0.924 |
| Llama-3-KO-EN | 8B-it | QLoRA | 42.862 | 0.896 | 0.931 |
| | 8B-it | 1-shot prompting | 2.031 | 0.490 | 0.673 |
| Llama-3-Open_Ko | 8B-it | QLoRA | 41.793 | 0.894 | 0.928 |
| GPT-4o | Unknown | 0-shot prompting | 33.406 | 0.889 | 0.915 |
| GPT-4o | Unknown | 1-shot prompting | 35.272 | 0.890 | 0.918 |

Table 9: Pure translation metrics $M$ for in-domain test data. The suffix '-it' indicates instruct-tuned models.

# Assessing the Role of Imagery in Multimodal Machine Translation

**Nicholas Kashani Motlagh[1], Jim Davis[1], Tim Anderson[2], Jeremy Gwinnup[2], Grant Erdmann[2]**

[1]Department of Computer Science and Engineering, Ohio State University

{kashanimotlagh.1, davis.1719}@osu.edu

[2]Air Force Research Laboratory

{timothy.anderson.20, jeremy.gwinnup.1, grant.erdmann}@us.af.mil

## Abstract

In Multimodal Machine Translation (MMT), the use of visual data has shown only marginal improvements compared to text-only models. Previously, the CoMMuTE dataset and associated metric were proposed to score models on tasks where the imagery is necessary to disambiguate between two possible translations for each ambiguous source sentence. In this work, we introduce new metrics within the CoMMuTE domain to provide deeper insights into image-aware translation models. Our proposed metrics differ from the previous CoMMuTE scoring method by 1) assessing the impact of multiple images on individual translations and 2) evaluating a model's ability to jointly select each translation for each image context. Our results challenge the conventional views of poor visual comprehension capabilities of MMT models and show that models can indeed meaningfully interpret visual information, though they may not leverage it sufficiently in the final decision.

## 1 Introduction

The use of multimodal data, combining visual and textual inputs, is becoming increasingly important in deep learning, especially in language modeling. Multimodal Machine Translation (MMT) presents a unique challenge in this area, as previous Machine Translation (MT) systems traditionally relied only on text. Despite the potential benefits of incorporating imagery, its efficacy in MMT remains controversial. Critics often view imagery as merely a regularizer rather than a core component of translation systems (Caglayan et al., 2016; Wu et al., 2021). This skepticism is fueled by results with the assumption that textual context alone suffices for most translation tasks (Caglayan et al., 2019).

To explore these concerns, the CoMMuTE dataset was developed to test MMT models on source sentences where visual context is essential for accurate selection between possible translations

(Futeral et al., 2023). Their proposed evaluation metric scores a model's *preference/choice between two reference translations*, diverging from traditional metrics such as BLEU (Papineni et al., 2002) and Meteor (Banerjee and Lavie, 2005) that instead compare a generated translation against a single reference. Initial analyses using the CoMMuTE dataset and metric indicate that current models show only slight, or no, improvement over using text-only models (Futeral et al., 2023).

Building on this recent foundation, we introduce a new complementary evaluative CoMMuTE metric that assesses a model's understanding of varying imagery on *a fixed reference translation* (as described above in (Futeral et al., 2023)). We additionally provide two group metrics designed to evaluate a model's ability to jointly choose each translation given their associated image contexts.

Results with our proposed metrics demonstrate that in many circumstances, models can indeed effectively understand and properly interpret the visual information, even if the final translation decisions are unaffected. This suggests the significant potential for improvements in model design to further leverage visual information.

## 2 Related Work

In this section, we present an overview of recent advancements and methodologies in two critical areas of related research. We first explore how imagery can enhance translation capabilities in MMT and subsequently shift our focus to contrastive evaluation methods, which represent a shift from traditional single-reference comparisons to more nuanced assessments using multiple contrasting references.

### 2.1 Multimodal Machine Translation

MMT typically trains with datasets such as Multi30k (Elliott et al., 2016) to enhance trans-

lation capabilities, yet results are not largely improved with sufficient textual context (Caglayan et al., 2019). Research such as Elliott (2018) demonstrates that the replacement of associated images with random counterparts often does not significantly impact translation quality, suggesting a predominant reliance on textual data. A later study further indicated that imagery typically serves merely as a form of regularization in training current models (Wu et al., 2021).

When imagery is available at inference time, approaches such as Graph-MMT (Yin et al., 2020), VTLM (Caglayan et al., 2021), Gated Fusion (Wu et al., 2021), and VGAMT (Futeral et al., 2023) are applicable. These methods leverage diverse global visual features from sources such as ResNet-50 (He et al., 2016) and CLIP (Radford et al., 2021), as well as visual semantic features through advanced object detectors like MDETR (Kamath et al., 2021).

In scenarios lacking visual data at inference time, innovative models such as CLIP-Trans (Gupta et al., 2023), UVR-NMT (Zhang et al., 2020), and ImagiT (Long et al., 2021) instead strategically leverage image-text datasets only during their training phase. These models employ sophisticated mechanisms to enhance their semantic understanding during training such as aligning image-text embedding spaces and synthesizing visual features. By pretraining on multimodal data, these models acquire a nuanced understanding of complex semantic relationships that text alone might not fully encapsulate. Some models, such as CLIP-Trans, can be modified to support the use of imagery at inference time by replacing CLIP text embeddings with CLIP image embeddings.

There has also been notable progress in adapting pretrained language models (LMs) such as BERT (Devlin et al., 2019) and XLM (Conneau and Lample, 2019) for multimodal use. Techniques such as visually-conditioned masked language modeling (VMLM) are explored in various architectures (Chen et al., 2020; Lu et al., 2019; Su et al., 2020; Li et al., 2020; Zhou et al., 2021; Ni et al., 2021; Futeral et al., 2023). Furthermore, the development of adapters and other lightweight modules can significantly enhance multimodal capabilities of LMs (Houlsby et al., 2019; Eichenberg et al., 2022; Yang et al., 2022; Tsimpoukelli et al., 2021; Sung et al., 2022; Futeral et al., 2023).

## 2.2 Contrastive Evaluation

Contrastive evaluation methodologies have become crucial for nuanced assessments of translation systems. These methodologies utilize contrastive test sets designed to challenge models to correctly rank pairs of translations, helping distinguish between correct and incorrect alternatives (Futeral et al., 2023). Contrastive datasets have been used to evaluate linguistic phenomena including grammaticality (Sennrich, 2017), pronoun translation (Müller et al., 2018; Bawden et al., 2018; Voita et al., 2019), and multi-sense word disambiguation (Rios Gonzales et al., 2017; Raganato et al., 2019; Futeral et al., 2023). Moreover, the coherence of lexical usage across translations has been thoroughly explored (Bawden et al., 2018; Voita et al., 2019).

## 3 CoMMuTE Dataset and Metric

The CoMMuTE dataset (Futeral et al., 2023) was recently introduced to score an MMT model's preference between two given translations for an ambiguous source based on the provided imagery. Specifically, CoMMuTE is comprised of 154 ambiguous English sentences, each paired with two contrasting images and their respective translations, where the two translations are available in French, German, and Czech. Each instance in the dataset is structured as a tuple $(s, i^a, t^a, i^b, t^b)$, where $s$ is an ambiguous source sentence and $(i^a, i^b)$ are images that disambiguate the sentence into two possible translations $(t^a, t^b)$, respectively. For example, in Fig. 1, the English source sentence "That's lots of bucks!" could refer to either deer or dollars, and the image is needed to determine the appropriate context.

To specifically score such disambiguation capabilities, the authors proposed a metric, which we refer to as TextCoMMuTE (TC), that compares the model's preference for the correct translation over the incorrect translation based on a single provided image context.

The model's uncertainty in a translation $t$ given a source $s$ and an image $i$ is quantified by perplexity, defined as

$$\mathcal{P}(s, i, t) = \exp\left(-\frac{1}{N}\sum_{k=1}^{N}\log p(t_k|s, i, t_{<k})\right)$$
(1)

Here, $N$ is the number of tokens in the translation, $t_k$ is the $k$-th token in the translation, and $p(t_k|s, i, t_{<k})$ denotes the conditional probability

of the $k$-th token given the source, image, and preceding tokens. In practice, this probability is approximated using the softmax of model outputs. Perplexity can be seen as a measure of uncertainty as it is the exponential of the negative mean log probability. Hence, *lower* perplexity is desired for a correct output versus an incorrect output.

The TC metric is then defined for a single image-translation triple $(i^m, t^m, t^n)$ as

$$TC^{m,n} = \mathbb{1}\{\mathcal{P}(s, i^m, t^m) < \mathcal{P}(s, i^m, t^n)\} \quad (2)$$

where $i^m$ and $t^m$ correspond to the matching image/translation and $t^n$ is the incorrect translation in the associated triple. Moreover, $\mathbb{1}$ is the indicator function that is 1 if the perplexity for the correct translation is less than that of the incorrect translation, and 0 otherwise.

Note that each of the 154 tuples in CoMMuTE yields 2 TC scores: $TC^{a,b}$ and $TC^{b,a}$. Hence, there are actually 308 individual TC scores for the dataset. An average is taken over the N=154 TC pairs as a summary statistic

$$TC = \frac{1}{2N} \sum_{j=1}^{N} \{TC^{a_j, b_j} + TC^{b_j, a_j}\} \quad (3)$$

Again, the TC score (Eqn. 3) views the two triples in each tuple *independently* even though both triples are associated with the same source sentence. TC scores range from 0-1 with 1 indicating correct disambiguation of all triples in the dataset. A text-only model scores a TC of 0.5 by definition (assuming no ties in perplexity) because for any tuple $j$ in the dataset, exactly one of $TC^{a_j, b_j}$ and $TC^{b_j, a_j}$ will be 1 while the other is 0 (*i.e.*, the image makes no contribution to the translation preference for a given source).

From an MMT perspective, this metric is insightful as translations with lower perplexities are typically more likely to be generated or appear higher in an n-best list.

## 4 Enhanced CoMMuTE Metrics

We now propose new complementary contrastive metrics to provide a more nuanced understanding of the interpretation of imagery for models with the CoMMuTE dataset.

### 4.1 ImageCoMMuTE

Rather than comparing two translations with the same image and source as is done with TC, we in-



(a) French Translation a: *Il y a beaucoup de cerfs !*

(b) French Translation b: *Cela fait beaucoup de dollars !*

Figure 1: English Source: *That's lots of bucks!*

stead examine the contribution of two *different* images to the *same* translation. From this perspective, we can directly assess whether the correctly associated image is appropriately affecting model uncertainty (reducing the perplexity of its corresponding translation). For a source $s$, images $(i^m, i^n)$, and a translation $t^m$, we define ImageCoMMuTE (IC) as

$$IC^{m,n} = \mathbb{1}\{\mathcal{P}(s, i^m, t^m) < \mathcal{P}(s, i^n, t^m)\} \quad (4)$$

where $i^m$ is the correctly associated image and $i^n$ is incorrectly associated image for translation $t^m$. Similar to TC, one can aggregate scores over a dataset by taking the mean of the N=154 pairs

$$IC = \frac{1}{2N} \sum_{j=1}^{N} \{IC^{a_j, b_j} + IC^{b_j, a_j}\} \quad (5)$$

Scores for IC range from 0-1, and a score of 0.5 indicates a random preference for the image context.

Our IC metric evaluates changes in model confidence for the *same* translation when presented with varying imagery. This approach directly assesses the interplay between imagery and text interpretation within the model. This differs from the work presented in Elliott (2018), where they assess average differences in model uncertainty, while we assess indicators of decisions. This IC metric also alleviates any possible concerns of the reliance on comparing perplexity averages and calibration *across* translations (as is done with TC). We will return to these potential issues in our discussion later. By maintaining a single reference translation across different visual contexts, our IC metric provides a more robust and precise measure of how imagery is understood by the model.

### 4.2 Group CoMMuTE

Though TC and IC are insightful metrics on their own, they both ignore the consistency desired for the underlying source-translation *pairs*. With TC,

the set of both *translations* is independently processed twice (each time with a different image context). Similarly with IC, the set of both *images* is independently processed twice (each time with a different translation target). What is truly desired is that the model consistently and correctly understands *both* cases for each set jointly to demonstrate true understanding.

Therefore, we propose a new group variant for TC and IC. To evaluate consistency across the paired nature of the task, we define Group TextCoMMuTE (GTC) as

$$GTC^{a,b} = TC^{a,b} \cdot TC^{b,a} \qquad (6)$$

and Group ImageCoMMuTE (GIC) as

$$GIC^{a,b} = IC^{a,b} \cdot IC^{b,a} \qquad (7)$$

These group metrics function with a logical "AND" between the two independent triple scores, ensuring that a score of 1 reflects consistent and correct interpretations for the tuple as a whole. As earlier, one can also aggregate group scores using a mean with

$$GTC = \frac{1}{N} \sum_{j=1}^{N} GTC^{a_j, b_j} \qquad (8)$$

$$GIC = \frac{1}{N} \sum_{j=1}^{N} GIC^{a_j, b_j} \qquad (9)$$

These scores also yield values between 0-1.

Our primary goal is to assess if the model properly interprets and understands imagery for the translations. Group scores such as GTC and GIC are crucial because they assess consistent model behavior with different text-image combinations, indicating true comprehension rather than coincidental correctness.

## 5 Experiments and Results

We present a comprehensive assessment of the previous and new CoMMuTE metrics on three pretrained English-to-French MMT models. Our evaluation is structured to elucidate how well these models understand the imagery with respect to resolving ambiguities in the CoMMuTE dataset. We begin by evaluating performance on the original CoMMuTE dataset, followed by an assessment using an extended set of imagery we collected for each CoMMuTE tuple to reveal further strengths and weaknesses across models.

### 5.1 Models

We employed three English-to-French MMT models, each chosen for its unique approach to integrating visual data with textual information. Across all models, we preprocessed imagery by resizing the smaller edge to 224px (maintaining the aspect ratio) and then taking a center crop of 224px $\times$ 224px.

**VGAMT.** The authors of CoMMuTE proposed VGAMT (Futeral et al., 2023), enhancing a pretrained mBART MT model (Liu et al., 2020) by incorporating CLIP ViT-B/32 image embeddings and fine-tuning adapters. While VGAMT included an object detector and a visually guided attention mechanism, our evaluation focused on its simplified variant from their ablation study (Futeral et al., 2023), which solely uses CLIP image embeddings. This model was trained using both visual masked language modeling and MMT objectives, having 1B total parameters. In our experiments, we employed three VGAMT models provided by the authors, each trained with a different random seed.

**CLIP-Trans.** The authors (Gupta et al., 2023) align the embedding spaces of a pretrained mBART MT model (Liu et al., 2020) with a multilingual M-CLIP model (Carlsson et al., 2022) via a mapping network. The model first trains on an image-captioning task using M-CLIP image embeddings followed by text-only MT training with M-CLIP text embeddings. They also suggest that imagery can be utilized at inference time, substituting M-CLIP text embeddings with image embeddings, even though it is not directly trained on MMT. We used a model following this approach with 1.3B total parameters. In the experiments, we evaluated one CLIP-Trans model provided by the authors.

**Gated Fusion.** This model introduces a dynamic gating mechanism that adaptively combines image and text representations, with gate values ranging from 0 to 1 for image components (Wu et al., 2021). The model leverages ResNet-50 (He et al., 2016) image features and a tiny transformer for a total of 32M parameters (substantially smaller than CLIP-Trans and VGAMT). We trained the model solely on the Multi30K dataset (Elliott et al., 2016), adhering to the authors' training protocol. We observed that the gating mechanism frequently assigns low values, often near 0, which tends to minimize the impact of visual data. To better incorporate image content into the translation process, we trained additional variants with fixed gate values of 0.25, 0.5,

Figure 2: Mixed imagery from Fig. 1 used for a pseudo-text-only baseline.

and 0.75. Each of these variants was trained and evaluated using three different random seeds.

## 5.2 Baseline Results

We first conducted a baseline evaluation on the CoMMuTE dataset. The second and third columns in Table 1 display the mean TC and GTC scores taken across models with random seeds (standard deviations were very low in all cases). For reference, a pure text-only MT model will have TC=0.5 and GTC=0, since the model will always choose one translation over the other for each tuple.

VGAMT scores highest in these two metrics, with the CLIP-Trans and Gated Fusion variants scoring near text-only in TC. This model also scores the highest in BLEU on Multi30k, as reported in previous work (Futeral et al., 2023; Gupta et al., 2023). The GTC scores of all models are above 0%, suggesting that all models can consistently disambiguate at least some tuples, though the scores are low. The gate values within the default Gated Fusion model were inspected and found to be near 0 (as expected). Interestingly, we see that TC for Gated Fusion improves slightly with a fixed larger gate value of 0.25 indicating that the strength of imagery does have the potential to change translations.

## 5.3 Comparison with Ambiguous Imagery

We next examined how much the imagery affected model decisions in comparison to the underlying textual bias. We compared the changes in TC scores using the original image context pairs (from CoMMuTE) versus an ambiguous mixed image.

As MMT models are trained with both imagery and text, one cannot properly obtain a pure text-only result through simple methods such as passing a zero image or removing the image context from the tokens. To obtain a pseudo-text-only baseline, we employed a 50/50% "mixup" (Zhang et al., 2018) of the two image contexts for each tuple to



Figure 3: Perplexities of the correct translations using the correct image, the incorrect image, and the mixed image.

create a single ambiguous image (see Fig. 2). Here, both image contexts are provided in a single image. However, there are other possible ways to create ambiguous imagery, such as arranging the images side-by-side. In Fig. 3, we see the perplexities of the correct translations using the mixed imagery typically fall between the perplexities using the correct and incorrect imagery, supporting the use of the mixed imagery as a baseline for comparison. We evaluated TC using this mixed image and also using the original images to get two competing TC scores for each image-translation triple. Note that the pseudo-text-only MMT model will score TC=0.5 (and GTC=0) by definition (we are using the same mixed image across two comparisons, and thus, preference does not change).

We measure changes in the score between the original images and the mixed image for each tuple using four consistency rates. The first two rates measure the percent of image-translation triples for which the original imagery and the mixed imagery gave different preferences for translations. That is, in these cases, the model's decision when using the original imagery was different from the model's decision when using the mixed imagery. The inconsistent positive rate (IPR) measures the percentage of image-translation triples that chose the right translation with the original imagery and the opposite/wrong translation with mixed imagery. The inconsistent negative rate (INR) measures the percentage of image-translation triples that chose the wrong translation with the original imagery and the opposite/right translation with mixed imagery. The performance of the remaining examples can be

| Model | Mean TC ↑ | Mean GTC ↑ | IPR ↑ | INR ↓ | CPR ↑ | CNR ↓ |
|---|---|---|---|---|---|---|
| VGAMT | 0.63 | 0.26 | 0.13 | 0.00 | 0.50 | 0.37 |
| CLIP-Trans | 0.51 | 0.03 | 0.01 | 0.00 | 0.50 | 0.49 |
| Gated Fusion | 0.50 | 0.02 | 0.01 | 0.01 | 0.49 | 0.49 |
| Gated Fusion$_{0.25}$ | 0.52 | 0.10 | 0.07 | 0.05 | 0.45 | 0.43 |
| Gated Fusion$_{0.5}$ | 0.50 | 0.07 | 0.05 | 0.05 | 0.45 | 0.45 |
| Gated Fusion$_{0.75}$ | 0.49 | 0.02 | 0.02 | 0.04 | 0.46 | 0.48 |

Table 1: Baseline TC and GTC scores on the original CoMMuTE dataset, and consistency rates compared to pseudo-text-only baseline.

quantified by a consistent positive rate (CPR) and a consistent negative rate (CNR), measuring the percentage of triples whose correct and incorrect preferences did not change when using the original or mixed imagery. Since the corpus is evenly split into 2 ambiguities, these rates are bounded in [0, 0.5] with IPR + CNR = INR + CPR = 0.5.

The last four columns in Table 1 display the consistency rates using the pseudo-text-only baseline for each of the models. The VGAMT model scores the highest IPR of 0.13 with an INR of 0, indicating that the model corrected 13% of translations without any negative impact when using the original imagery. In contrast, the CLIP-Trans and Gated Fusion variants show smaller IPR and INR rates, suggesting that imagery has a weaker yet still noticeable effect on these models. The higher INR rates for Gated Fusion models indicate that imagery can actually hurt their performance.

By examining the CPR and CNR rates in the table, we see that imagery may not be significantly impactful in the decisions across all models. These rates only measure the proportion of image-translation triples (with the original imagery) that *agree* with the pseudo-text-only baseline (with the mixed imagery). They do not describe if the model associates correct/incorrect imagery with translation confidence. The model still might correctly associate the original imagery, giving lower perplexity of the correct translation (desired), but this change may not be drastic enough to overturn the model's underlying textual preference. This highlights the need for a metric, such as the proposed IC, to measure how confidence in a translation changes with correct and incorrect imagery.

### 5.4 ImageCoMMuTE Results

We next conducted an evaluation of the CoMMuTE dataset using our proposed IC and GIC metrics. Table 2 displays the mean IC and GIC scores taken across the models with random seeds. Note that IC

| Model | Mean IC ↑ | Mean GIC ↑ |
|---|---|---|
| VGAMT | 0.81 | 0.66 |
| CLIP-Trans | 0.58 | 0.22 |
| Gated Fusion | 0.51 | 0.11 |
| Gated Fusion$_{0.25}$ | 0.51 | 0.12 |
| Gated Fusion$_{0.5}$ | 0.50 | 0.13 |
| Gated Fusion$_{0.75}$ | 0.50 | 0.11 |

Table 2: Baseline IC and GIC scores.

| Model | TC | IC |
|---|---|---|
| VGAMT vs CLIP-Trans | 0.39 | 0.18 |
| VGAMT vs Gated Fusion$_{0.25}$ | 0.25 | 0.16 |
| Gated Fusion$_{0.25}$ vs CLIP-Trans | 0.36 | 0.32 |

Table 3: Intersection-Over-Union of failures as determined by TC and IC.

and GIC metrics are undefined for a pure text-only MT model, and thus, we cannot compute the four consistency rates.

Our image-based metrics (IC and GIC) demonstrate that VGAMT interprets imagery most effectively, achieving 0.81 on IC and 0.66 on GIC, which are significantly higher than the TC of 0.63 and GTC of 0.26. Other models continue to score only slightly above 0.5. We find that of the models we tested, those that scored highest on MMT quality metrics also scored highest in our proposed metrics (as reported in (Futeral et al., 2023; Gupta et al., 2023)). These results demonstrate that VGAMT more appropriately adjusts uncertainty in a translation based on imagery.

We also investigated whether the different models made the same errors. We identified the image-translation triples where each model made errors in terms of TC and also for IC. We then calculated the intersection-over-union (IOU) between 2 models, which is a set similarity metric defined as the ratio of the number of image-translation triples common to both error sets for a given metric (intersection)

to the total number of unique image-translation triples in both error sets (union). This metric helps quantify the similarity in errors across models as a scalar bounded in [0,1] where 1 signifies exact similarity in errors. The results in Table 3 reveal that models do not strongly make the same mistakes yet do share some overlap.

### 5.5 Extended CoMMuTE

We next extended the CoMMuTE dataset by incorporating additional images per translation in each tuple. This extension allows for a broader assessment of model performance across diverse image inputs and enables a search for images that could either improve or degrade the scores.

For each ambiguous source $s$, we manually generated two distinct, <u>un</u>ambiguous captions, $c^a$ and $c^b$, which correspond directly to the translations $t^a$ and $t^b$, respectively. For example, the English sentence "That's lots of bucks!" is transformed to "a photo of deer" and "a photo of dollars".

Utilizing these unambiguous captions, we then sourced corresponding images from the DataComp-12.8M dataset (Gadre et al., 2023), which comprises 12.8 million image-text pairs harvested from the Common Crawl (Common Crawl). The DataComp dataset serves as a foundation dataset for enhancing the training of CLIP models. We employ a CLIP ViT-B/32 model, pretrained on the LAION-5B dataset (Schuhmann et al., 2022), to retrieve images most similar (cosine similarity) to our unambiguous captions.

From this candidate set of imagery, the top 15 images that most closely aligned with each caption, adhering to a minimum dimension of 64 pixels and a maximum aspect ratio of 2.5, were retrieved automatically. We manually selected the four most representative images from this set (due to potentially noisy images retrieved). If fewer than 4 suitable images were found, additional images were sourced from Google Images. This method resulted in a total of 1540 images, providing 5 images (instead of just 1) for each unambiguous translation. Consequently, this extended CoMMuTE dataset includes the original source $s$, translations $t^a$ and $t^b$, and now 5 images each for $i^a$ and $i^b$.

With this extended CoMMuTE dataset, we examined if there existed subsets of imagery that could significantly increase or decrease the GIC score (as we deem GIC the most important metric for each model). For each tuple in our extended dataset, we identified the image pair (one image taken from each image set) that maximizes or minimizes the GIC score. As multiple pairs can meet the criteria, we select the pair that optimizes

$$
\begin{aligned}
\{\mathcal{P}(s, i^a, t^a) - \mathcal{P}(s, i^b, t^a)\} \quad & + \\
\{\mathcal{P}(s, i^b, t^b) - \mathcal{P}(s, i^a, t^b)\} &
\end{aligned} \tag{10}
$$

This expression reflects the confidence gaps for the translations. Given that a lower perplexity indicates a better result and considering the ordering of differences in Eqn. 10, we minimize (or maximize) this equation to maximize (or minimize) the GIC score accordingly. When seeking images to maximize the GIC score, we break ties by finding the image pair that *minimizes* Eqn. 10 (can be negative). When seeking images to minimize the GIC score, we break ties with the image pair that *maximizes* Eqn. 10. We refer to the image subset specifically tailored to maximize GIC as Image-Oracle. We also tracked the replacement rate (RR) of the number of images replaced from the original dataset.

As shown in Table 4, the maximal GIC image subsets show high effectiveness, with VGAMT scoring a Max IC of 0.96 and a Max GIC of 0.92. This suggests that the model can accurately interpret the intended visual signals in these particular image pairs for nearly all translations. This is further supported by the notably higher Max IC and GIC scores in the CLIP-Trans and Gated Fusion variants. Conversely, we see that sets of images can be found to hurt performance, especially in CLIP-Trans and Gated Fusion. Examples of replaced imagery can be seen in Fig. 4. Therefore, it is possible to have imagery that drastically improves or degrades the scores. We see that replacement rates are high, indicating that the original dataset is not prominent in these maximal/minimal subsets. The results with maximal/minimal GIC show that the model does indeed have an internal understanding of the imagery with respect to the translation task.

We would expect the Image-Oracle images that maximized GIC to similarly improve TC and GTC scores. However, Table 5 shows only minor improvements in TC and GTC across models. Thus, even though the IC and GIC metrics strongly indicate the image interpretability of the models, the TC and GTC metrics fail to highlight the potential contribution of imagery.

### 6 Discussion

This study introduced image-based and group metrics for CoMMuTE to better evaluate if models do

| Model | Min IC ↑ | Min GIC ↑ | RR | Max IC ↑ | Max GIC ↑ | RR |
|---|---|---|---|---|---|---|
| VGAMT | 0.59 | 0.33 | 0.80 | 0.96 | 0.92 | 0.71 |
| CLIP-Trans | 0.46 | 0.01 | 0.77 | 0.89 | 0.77 | 0.77 |
| Gated Fusion | 0.40 | 0.00 | 0.77 | 0.73 | 0.48 | 0.78 |
| Gated Fusion$_{0.25}$ | 0.38 | 0.00 | 0.80 | 0.86 | 0.71 | 0.80 |
| Gated Fusion$_{0.5}$ | 0.35 | 0.00 | 0.81 | 0.88 | 0.76 | 0.77 |
| Gated Fusion$_{0.75}$ | 0.37 | 0.00 | 0.79 | 0.85 | 0.71 | 0.80 |

Table 4: Minimum and maximum IC and GIC scores along with replacement rates.



Figure 4: Examples from the CoMMuTE dataset with original imagery (top row), oracle best replacements (middle row), and oracle worst replacements (bottom row) as determined by VGAMT.

| Model | Mean TC ↑ | Mean GTC ↑ |
|---|---|---|
| VGAMT | 0.67 | 0.34 |
| CLIP-Trans | 0.52 | 0.05 |
| Gated Fusion | 0.51 | 0.02 |
| Gated Fusion$_{0.25}$ | 0.64 | 0.28 |
| Gated Fusion$_{0.5}$ | 0.59 | 0.18 |
| Gated Fusion$_{0.75}$ | 0.56 | 0.12 |

Table 5: Image-Oracle TC and GTC scores.

| Model | Mean TC ↑ | Mean GTC ↑ |
|---|---|---|
| VGAMT | 0.66 | 0.32 |
| CLIP-Trans | 0.52 | 0.03 |
| Gated Fusion | 0.51 | 0.01 |
| Gated Fusion$_{0.25}$ | 0.60 | 0.21 |
| Gated Fusion$_{0.5}$ | 0.58 | 0.15 |
| Gated Fusion$_{0.75}$ | 0.53 | 0.07 |

Table 6: Image-Oracle TC scores with the shared prefix removed in perplexity computation.

understand imagery in MMT. In this section, we explore possible reasons why TC scores are so much lower than IC and discuss future directions on how to further leverage the imagery to improve MMT.

There are two potential issues related to perplexity and calibration that may affect the TC/GTC scores. First, there is an assumption that perplexity is indeed an appropriate uncertainty metric to compare *two* translations. Perplexity is a transform

of the mean log probability and, therefore, relies on averages where all tokens are weighted equally (Ueda et al., 2024). There may indeed be other better measures of uncertainty (Kauf and Ivanova, 2023). It is also assumed that the model is well calibrated to properly compare *across* translations.

One method to examine the effects of averages across sequences of different lengths in the perplex-

Figure 5: Calibration results using temperature scaling.

ity computation is to remove any shared prefix in $t^a, t^b$ before computing perplexity and then compare to the results without prefix removal (original method). Ignoring common prefixes (while still weighting the remaining tokens equally) actually shows a slight degradation in scores (as illustrated in Table 6). These results suggest perplexity (a transform of mean log probability) does have some issues as a comparison method. However, this does not fully explain the low TC/GTC scores.

We also investigated the effects of model calibration using a simple global temperature scaling method (Guo et al., 2017) across a range of temperature values from 0.25 to 2. As shown in Fig. 5, the TC scores appear unaffected, indicating potential miscalibration, while IC scores suggest that models are relatively well-calibrated (at T=1). We also examined higher temperatures, which did not change the results, suggesting calibration does not appear to be primarily responsible for the TC/GTC degradation.

Therefore, given the stronger results from IC/GIC, we believe the main overall issue with TC/GTC is that the underlying textual preference/bias in these models is too strong and does not allow much influence from the imagery (which we have shown to be interpreted well by the models).

## 7 Recommendations for Future Work

One future area of work is the integration of imagery *earlier* in the model's architecture rather than appending them at the end of the processing chain (Wu et al., 2021; Gupta et al., 2023). Integrating image features earlier in the model's architecture could enhance the model's ability to better leverage the rich contextual cues provided by the imagery. This approach may result in translations that are more contextually nuanced, with increased attention to specific words critical for disambiguation.

Additionally, enhancing the impact of visual sig-

nals *within* the model could also prove beneficial. This could be achieved by adjusting the gate values in models that use gating mechanisms, such as Gated Fusion (Wu et al., 2021), to strengthen the influence of visual data. As demonstrated, setting a fixed gate value that prioritizes visual information could help in situations where visual context is crucial for disambiguating textual content. Even though the non-gated VGAMT was the top performer, there is still room for improvement by strengthening the role of imagery in the processing using some method of gating or amplification.

Earlier we have shown that the IOU of errors between model pairs did not have strong alignment. This diversity implies that ensembling different models could potentially mitigate individual weaknesses and enhance overall performance.

## 8 Conclusion

Our study challenges the widespread belief that visual cues are not generally very helpful to MMT. By employing our proposed IC and Group CoMMuTE metrics within an expanded CoMMuTE dataset, we have established a robust framework for assessing if visual information is understood in MMT systems. Our results reveal that while visual data does indeed support translation preferences, it is not leveraged significantly to enhance the outcomes over the underlying textual bias. Our findings mark a promising direction for future research in MMT, suggesting that further exploration could uncover ways to amplify this positive impact.

## Limitations

Firstly, we evaluated English-French translations in CoMMuTE. It remains to be seen whether the results generalize to other languages. Additionally, our evaluations were conducted on an extended set of 5 images, whereas larger sets (e.g., 100 images) would provide more robust insights. Furthermore, we relied on the default single reference translation for each image. Having additional translations for each image context would enable a more comprehensive evaluation.

# References

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Rachel Bawden, Rico Sennrich, Alexandra Birch, and Barry Haddow. 2018. Evaluating Discourse Phenomena in Neural Machine Translation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1304–1313, New Orleans, Louisiana. Association for Computational Linguistics.

Ozan Caglayan, Walid Aransa, Yaxing Wang, Marc Masana, Mercedes García-Martínez, Fethi Bougares, Loïc Barrault, and Joost van de Weijer. 2016. Does Multimodality Help Human and Machine for Translation and Image Captioning? In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 627–633, Berlin, Germany. Association for Computational Linguistics.

Ozan Caglayan, Menekse Kuyu, Mustafa Sercan Amac, Pranava Madhyastha, Erkut Erdem, Aykut Erdem, and Lucia Specia. 2021. Cross-lingual Visual Pretraining for Multimodal Machine Translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1317–1324, Online. Association for Computational Linguistics.

Ozan Caglayan, Pranava Madhyastha, Lucia Specia, and Loïc Barrault. 2019. Probing the Need for Visual Context in Multimodal Machine Translation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4159–4170, Minneapolis, Minnesota. Association for Computational Linguistics.

Fredrik Carlsson, Philipp Eisen, Faton Rekathati, and Magnus Sahlgren. 2022. Cross-lingual and Multilingual CLIP. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6848–6854, Marseille, France. European Language Resources Association.

Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. UNITER: UNiversal Image-TExt Representation Learning. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, volume 12375, pages 104–120. Springer International Publishing, Cham. Series Title: Lecture Notes in Computer Science.

Common Crawl. Common crawl. https://commoncrawl.org.

Alexis Conneau and Guillaume Lample. 2019. Crosslingual Language Model Pretraining. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Constantin Eichenberg, Sidney Black, Samuel Weinbach, Letitia Parcalabescu, and Anette Frank. 2022. MAGMA – Multimodal Augmentation of Generative Models through Adapter-based Finetuning. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2416–2428, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Desmond Elliott. 2018. Adversarial Evaluation of Multimodal Machine Translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2974–2978, Brussels, Belgium. Association for Computational Linguistics.

Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30K: Multilingual English-German Image Descriptions. In *Proceedings of the 5th Workshop on Vision and Language*, pages 70–74, Berlin, Germany. Association for Computational Linguistics.

Matthieu Futeral, Cordelia Schmid, Ivan Laptev, Benoît Sagot, and Rachel Bawden. 2023. Tackling Ambiguity with Images: Improved Multimodal Machine Translation and Contrastive Evaluation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5394–5413, Toronto, Canada. Association for Computational Linguistics.

Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, Eyal Orgad, Rahim Entezari, Giannis Daras, Sarah Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei W Koh, Olga Saukh, Alexander J Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. 2023. DataComp: In Search of the Next Generation of Multimodal Datasets. In *Advances in Neural Information Processing Systems*, volume 36, pages 27092–27112. Curran Associates, Inc.

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. On Calibration of Modern Neural Networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1321–1330. PMLR. ISSN: 2640-3498.

Devaansh Gupta, Siddhant Kharbanda, Jiawei Zhou, Wanhua Li, Hanspeter Pfister, and Donglai Wei. 2023. CLIPTrans: Transferring Visual Knowledge with Pretrained Models for Multimodal Machine Translation. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2863–2874, Paris, France. IEEE.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, Las Vegas, NV, USA. IEEE.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-Efficient Transfer Learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2790–2799. PMLR. ISSN: 2640-3498.

Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. 2021. MDETR - Modulated Detection for End-to-End Multi-Modal Understanding. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1760–1770, Montreal, QC, Canada. IEEE.

Carina Kauf and Anna Ivanova. 2023. A Better Way to Do Masked Language Model Scoring. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 925–935, Toronto, Canada. Association for Computational Linguistics.

Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. 2020. Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, volume 12375, pages 121–137. Springer International Publishing, Cham. Series Title: Lecture Notes in Computer Science.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual Denoising Pre-training for Neural Machine Translation. *Transactions of the Association for Computational Linguistics*, 8:726–742. Place: Cambridge, MA Publisher: MIT Press.

Quanyu Long, Mingxuan Wang, and Lei Li. 2021. Generative Imagination Elevates Machine Translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5738–5748, Online. Association for Computational Linguistics.

Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Mathias Müller, Annette Rios, Elena Voita, and Rico Sennrich. 2018. A Large-Scale Test Set for the Evaluation of Context-Aware Pronoun Translation in Neural Machine Translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 61–72, Brussels, Belgium. Association for Computational Linguistics.

Minheng Ni, Haoyang Huang, Lin Su, Edward Cui, Taroon Bharti, Lijuan Wang, Dongdong Zhang, and Nan Duan. 2021. M3P: Learning Universal Representations via Multitask Multilingual Multimodal Pre-Training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3977–3986.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: A Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning*, pages 8748–8763. PMLR. ISSN: 2640-3498.

Alessandro Raganato, Yves Scherrer, and Jörg Tiedemann. 2019. The MuCoW Test Suite at WMT 2019: Automatically Harvested Multilingual Contrastive Word Sense Disambiguation Test Sets for Machine Translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 470–480, Florence, Italy. Association for Computational Linguistics.

Annette Rios Gonzales, Laura Mascarell, and Rico Sennrich. 2017. Improving Word Sense Disambiguation in Neural Machine Translation with Sense Embeddings. In *Proceedings of the Second Conference on Machine Translation*, pages 11–19, Copenhagen, Denmark. Association for Computational Linguistics.

Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski,

Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. LAION-5B: An Open Large-scale Dataset for Training Next Generation Image-Text Models. *Advances in Neural Information Processing Systems*, 35:25278–25294.

Rico Sennrich. 2017. How Grammatical is Character-level Neural Machine Translation? Assessing MT Quality with Contrastive Translation Pairs. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 376–382, Valencia, Spain. Association for Computational Linguistics.

Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *International Conference on Learning Representations*.

Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. 2022. VL-ADAPTER: Parameter-Efficient Transfer Learning for Vision-and-Language Tasks. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5217–5227, New Orleans, LA, USA. IEEE.

Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, S. M. Ali Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal Few-Shot Learning with Frozen Language Models. In *Advances in Neural Information Processing Systems*, volume 34, pages 200–212. Curran Associates, Inc.

Naoya Ueda, Masato Mita, Teruaki Oka, and Mamoru Komachi. 2024. Token-length Bias in Minimal-pair Paradigm Datasets. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 16224–16236, Torino, Italia. ELRA and ICCL.

Elena Voita, Rico Sennrich, and Ivan Titov. 2019. When a Good Translation is Wrong in Context: Context-Aware Machine Translation Improves on Deixis, Ellipsis, and Lexical Cohesion. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1198–1212, Florence, Italy. Association for Computational Linguistics.

Zhiyong Wu, Lingpeng Kong, Wei Bi, Xiang Li, and Ben Kao. 2021. Good for Misconceived Reasons: An Empirical Revisiting on the Need for Visual Context in Multimodal Machine Translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6153–6166, Online. Association for Computational Linguistics.

Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Zero-Shot Video Question Answering via Frozen Bidirectional Language Models. In *Advances in Neural Information Processing Systems*.

Yongjing Yin, Fandong Meng, Jinsong Su, Chulun Zhou, Zhengyuan Yang, Jie Zhou, and Jiebo Luo. 2020. A Novel Graph-based Multi-modal Fusion Encoder for Neural Machine Translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3025–3035, Online. Association for Computational Linguistics.

Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. 2018. mixup: Beyond Empirical Risk Minimization. In *International Conference on Learning Representations*.

Zhuosheng Zhang, Kehai Chen, Rui Wang, Masao Utiyama, Eiichiro Sumita, Zuchao Li, and Hai Zhao. 2020. Neural Machine Translation with Universal Visual Representation. In *International Conference on Learning Representations*.

Mingyang Zhou, Luowei Zhou, Shuohang Wang, Yu Cheng, Linjie Li, Zhou Yu, and Jingjing Liu. 2021. UC2: Universal Cross-lingual Cross-modal Vision-and-Language Pre-training. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4153–4163.

# Error Span Annotation:
# A Balanced Approach for Human Evaluation of Machine Translation

**Tom Kocmi**[★][1]    **Vilém Zouhar**[★][2]    **Eleftherios Avramidis**[3]    **Roman Grundkiewicz**[1]
**Marzena Karpinska**[4]    **Maja Popović**[5]    **Mrinmaya Sachan**[2]    **Mariya Shmatova**[6]

[1]Microsoft    [2]ETH Zurich    [3]DFKI
[4]UMass Amherst    [5]DCU & IU    [6]Dubformer

tomkocmi@microsoft.com    vzouhar@inf.ethz.ch

## Abstract

High-quality Machine Translation (MT) evaluation relies heavily on human judgments. Comprehensive error classification methods, such as Multidimensional Quality Metrics (MQM), are expensive as they are time-consuming and can only be done by experts, whose availability may be limited especially for low-resource languages. On the other hand, just assigning overall scores, like Direct Assessment (DA), is simpler and faster and can be done by translators of any level, but is less reliable. In this paper, we introduce Error Span Annotation (ESA), a human evaluation protocol which combines the continuous rating of DA with the high-level error severity span marking of MQM. We validate ESA by comparing it to MQM and DA for 12 MT systems and one human reference translation (English to German) from WMT23. The results show that ESA offers faster and cheaper annotations than MQM at the same quality level, without the requirement of expensive MQM experts.

## 1 Introduction

While automatic evaluation metrics are important and invaluable tools for rapid development of Machine Translation (MT) systems, human assessment remains the gold standard of translation quality (Kocmi et al., 2023; Freitag et al., 2023). The translation quality is conceptually measured through adequacy (preservation of the original meaning) and fluency (grammaticality of the translated text; Koehn and Monz, 2006), and sometimes through comprehension (how readable or understandable the translation is; White et al., 1994).

Annotators are usually asked to assign a score on a particular quality aspect. Likert and 0–100 scale are often used for discrete and continuous scales.



Figure 1: Stylized annotation user interface with Error Span Annotation (ESA). The annotator first marks errors with minor and major severity and then assigns a final score. This is more robust than asking for score directly.[1]

The most popular scoring method in machine translation field in recent years is Direct Assessment (DA; Graham et al., 2013), which is used to portray a human assessment of MT quality in the WMT shared tasks since 2016. Since 2022, the DA+SQM metric is used, namely direct assessment enriched with more objective Scalar Quality Metrics (SQM) guidelines (Kocmi et al., 2022).

Translation scores indicate the overall quality of a translation, but they can be subjective and do not provide details about the translation errors. The usual way to overcome this drawback is error classification: asking the evaluators to mark each translation error and assign an error tag from a set of predefined categories, such as *terminology* or *style*.

---

★Equal contributions. Others alphabetically.
[0]Code & collected data:
  ⌂ github.com/wmt-conference/ErrorSpanAnnotation

---

[1]Our experiments are on English→German. In Figure 1, Spanish→English is only an illustration for English-speaking readers. The first example, based on our data, omits *Rainy days!* and incorrectly translates *about field items* as *im Feld* (=*in the field*). The second example, for illustrative purposes, does not capitalize *H* and mistakenly adds extra *week*.

In recent years, the dominant error classification protocol is the Multidimensional Quality Metrics (MQM; Lommel et al., 2014; Freitag et al., 2021a). MQM error classification is the standard human metric in the WMT Metrics shared task since 2021 (Freitag et al., 2021b). While error classification provides interesting insights into the distribution of different types of errors, it requires much more time and effort, both for annotators and task organizers, who need to prepare the error taxonomy, annotation guidelines and training examples.

We present a new evaluation protocol based on highlighting errors and followed by assigning scores, Error Span Annotation (ESA), and compare it to the MQM error classification and DA+SQM scores. We compare MQM, DA+SQM and ESA annotations in parallel on a subset of English→German machine translation outputs from WMT23. We find that the proposed ESA protocol is faster and cheaper than MQM whilst providing the same usefulness in ranking MT systems.

## 2 Related work

Assigning overall scores was the very first method of manual MT evaluation (ALPAC, 1966; White et al., 1994), where the evaluators assessed some or all of the translation quality criteria at once: adequacy, comprehensibility, and fluency. The first WMT (Workshop/Conference on Machine Translation) shared task in 2006 and the subsequent task in 2007 adopted this technique and used adequacy and fluency scores as official metrics (Koehn and Monz, 2006; Callison-Burch et al., 2007). Later, Vilar et al. (2007) proposed binary ranking of two or more MT outputs, which became the official metric at WMT 2008 (Callison-Burch et al., 2008). It required less effort and showed better inter-annotator agreement than adequacy and fluency scores. This remained the official WMT metric until 2017 when it was replaced by continuous Direct Assessment (DA). DA (Graham et al., 2013) does not use discrete scales, but a continuous one between 0 and 100. Bojar et al. (2016) scrutinized the quality criteria and recommended to focus on adequacy and use fluency to break ties only. DA replaced ranking methods in 2017 (Bojar et al., 2017) and since 2022 (Kocmi et al., 2022) it is used with SQM guidelines (Freitag et al., 2021a) in a slightly modified version with more descriptive scale labels, which increased the inter-annotator agreement.

None of the described methods provides information about the erroneous or problematic parts of the translation. An early work of Vilar et al. (2006) analyzes errors in translation outputs assigning them to error classes from a predefined error typology. Most popular error typology recently is Multidimensional Quality Metrics (MQM; Lommel et al., 2014; Klubička et al., 2018; Freitag et al., 2021a), which is used in WMT metrics task since 2021 (Freitag et al., 2021b).

Several error span marking methods have been proposed recently (Kreutzer et al., 2020; Popović, 2020) as a less demanding error annotation approach than error classification. While it does not provide the fine-grained details about different error classes, it still gives the information about the position and amount of errors, and also enables further fine-grained analysis on the annotated data, if necessary (e.g. classification of already marked errors, identifying linguistic phenomena causing the errors, or focusing on particular error type). While the previously reported findings on this method are promising, no systematic comparison for the purposes of evaluating machine translation systems has been carried out so far. Furthermore, the simplified error marking method does not solve the challenges in determining how to appropriately weight individual errors to obtain segment-level scores, a problem that becomes particularly pronounced when extending the evaluation to document level.

This work combines the advantages of error annotation (like MQM) and assigning direct scores (like DA). The annotators are first asked to identify and mark all errors, and afterwards to assign an overall score. When deciding about the score, they are *primed* by the preceding error annotation and see all the marked errors that can be taken into consideration for the final score.

## 3 Comparison with DA+SQM and MQM

Our proposed method lies between DA+SQM and MQM protocols, so we provide a detailed comparison between the two before describing ESA in details in the next section.

While both DA+SQM and MQM generally exhibit low inter-annotator agreement (Knowles and Lo, 2024; Freitag et al., 2021a), DA+SQM scores have high variance, which needs to be compensated with higher number of annotations per system (Wei et al., 2022). On the other hand, MQM requires human experts trained with the MQM protocol and

error classification. Trained experts can be twice as expensive as translators or bilingual speakers evaluating DA+SQM. The required expertise is a hard constraint which makes evaluation on some languages prohibitively expensive or not possible at all, especially low-resource languages. Furthermore, assigning a DA+SQM numerical score to a segment is anecdotally and intuitively much faster than MQM, where the evaluators need to mark each error span, classify it and assign severity. This altogether can make each MQM annotated segment up to approximately $10\times$ more expensive.

DA+SQM is usually based on sentence-level scores, and the paragraph-level score is computed as the average of all sentence scores in the paragraph. Paragraph-level DA+SQM evaluation is possible, but evaluating an entire paragraph takes more time, substantially decreases the total number of collected scores, and is more demanding cognitively, which negatively impacts the inter-annotator agreement (Castilho, 2020). On the other hand, MQM as an error classification annotation is agnostic to the choice of annotation unit.

As an example, in App. Figure 6 we show the system ranking based on two approaches DA+SQM and MQM on Chinese-English (sentence-level evaluation) and English-German (paragraph-level evaluation) language pairs. Although both techniques reach same order of system clusters, DA+SQM produces much fewer clusters in paragraph-level setup, thus putting many systems within a single cluster. On the other hand, MQM is better able to distinguish different systems. To increase statistical power of DA+SQM, we would have to collect much more DA+SQM samples (Wei et al., 2022), which would further drive the cost up. In addition, DA is much more skewed towards fluency as opposed to adequacy (Martindale and Carpuat, 2018).

The cost difference was one of the main reason behind the high usage of DA+SQM at the WMT General MT shared task (Kocmi et al., 2023, 2022). For all these reasons, our hypothesis is that a new annotation protocol, ESA, which is between DA+SQM and MQM can provide better annotations than DA+SQM at a lower cost than MQM.

## 4 Error Span Annotation

**Annotation process.** In Error Span Annotation (ESA), the evaluators first mark all problematic parts (characters, words, phrases, sentences) in the translated text. For each marked span, they

| Source/Translation+ESA | Score |
|---|---|
| **SRC**: ... I've entered the burrata dimension. <br> **TGT**: ... ich habe die Burrata-Dimension eingegeben. <br> **gloss**: *habe eingegeben(=I put in)* should be *bin eingetreten* | 70% |
| **SRC**: Not like other tomb raider games <br> **TGT**: Nicht wie andere Gräberüberfäller Spiele <br> **gloss**: *Gräberüberfäller* overtranslates *Tomb Raider* | 35% |
| **SRC**: (PERSON2) Yeah, so just know like- <br> **TGT**: Ja, also weißt du einfach... [missing] <br> **gloss**: *PERSON2* is missing | 86% |
| **SRC**: All collards, kale, chard is transplanted. <br> **TGT**: Alle Kohlköpfe, Grünkohl, Schmalz sind verpflanzt. <br> **gloss**: *Kohlköpfe(=cabbages)* and *Schmalz(=lard)* are incorrect translations | 17% |

Example 1: ESA-annotated examples with associated manual score. The error severity distinction is between minor and major. The first example has a single error (confusion *eingeben(=put in)* as *enter* meaning *go in*) which also affects the auxiliary verb habe/bin. Thus, the same error is marked twice.

also provide one of the two severity levels: **major** (e.g. changed meaning) or **minor** (e.g. incorrect grammar, style; see Example 1). Because all error spans are marked in the translation, not the source text, we include a special tag for marking omission errors. This was an intentional choice over annotating the source text to make the annotation protocol forward-compatible with other translation modalities, such as audio and video translations.

After the annotators mark all the error spans, they are asked to provide an overall score for the entire segment, on the scale from 0 to 100, reminiscent of DA+SQM. We implement the annotation interface in Appraise (Federmann, 2018) and show a screenshot in Figure 2. The full guidelines displayed to annotators are shown in Appendix A.

**Segment-level scores.** To rank systems, we need scalar values. There are two evident ways to extract them from ESA: (1) using the annotator's overall segment-level score directly, like DA+SQM, or (2) converting error span severity levels into a segment-level score, like MQM. By instructing the annotators to identify and mark all errors first, we *prime* them to be more accurate when assessing the overall quality of the segment—when making the decision about the score, they have already marked all errors in the segment and can see them, therefore they can take them into consideration.

MQM is primarily error diagnostics protocol which has been repurposed for translation segment scoring. The transition from MQM error spans

into a single segment score has been proposed to be done with the formula based on error severity counts (Freitag et al., 2021a):

$$\text{MQM-like} = -5 \cdot \#\text{MAJOR} - 1 \cdot \#\text{MINOR}$$

Notice that this does not scale with different text sizes.[2] As an example, if translation of a segment has two major errors, it receives the score of $-10$. However, if the source is repeated twice and the corresponding translation as well, the score would be further decreased to $-20$. This is especially problematic for paragraph-level evaluation which features segments of different length. Additionally, the segment-level MQM score might not correspond with the segment-level translation quality, such as when marking one error affecting several places in the segment, as in Example 1 (top). To avoid such issues, we use the annotators' direct scores as the main scoring approach for the ESA protocol unless specified differently. The system-level scores for all human evaluation methods in this work are calculated as the average of all segment-level scores for particular system. We further revisit the score computation in App. Appendix B.1. In some analysis, we use ESA error spans to calculate MQM-like score, we refer to such scores as $\text{ESA}_{\text{spans}}$.

**Advantages.** Assigning overall scores is guided by errors in the translation, and through error marking, the annotator can first focus on direct highlighting of these issues instead of determining the overall score directly. The advantage of ESA over error classification is that it is less demanding, while still informative—the annotations can be further refined in subsequent analyses. Furthermore, the evaluators are not limited to any pre-defined annotation protocol and can highlight a larger range of errors. The error marking approach can be seen as *descriptive* (encouraging annotator subjectivity and capturing their individual beliefs) and the error classification as a *prescriptive* (discouraging annotator subjectivity and asking annotators to align with one specific belief, in this case the pre-defined error protocol), as per Rottger et al. (2022).

## 5 Experimental setup

We conduct experiments comparing the ESA protocol with MQM and DA+SQM protocols. We design them in a way which makes it comparable with the previously collected annotations for

---

|  | $\text{ESA}_1$ | $\text{ESA}_2$ | MQM | $\text{MQM}^{\text{WMT}}$ |
|---|---|---|---|---|
| # error spans | 0.45 | 1.00 | 0.53 | 3.37 |
| % minor | 63% | 68% | 67% | 67% |
| % major | 37% | 32% | 33% | 33% |
| Score (MQM-like) | 81.8 (-1.1) | 84.5 (-2.2) | (-1.2) | (-7.1) |

Table 1: Average number of error spans per segment, ratio between minor and major errors, and scores across different annotation protocols.

$\text{MQM}^{\text{WMT}}$ and $\text{DA+SQM}^{\text{WMT}}$. For this reason, we reproduce the human evaluation campaign for the WMT23 English to German systems (Kocmi et al., 2023; Freitag et al., 2023) with ESA and our reimplementation of MQM.

The original campaign featured 13 translations of 557 source segments. To facilitate running multiple campaigns for proper comparison, we had to scale down and subsampled 207 segments per system (74 documents), which yields 2,691 segments in total. To keep the ESA annotation comparable to other protocols, we subsample by selecting a subset of documents evaluated by Freitag et al. (2023) keeping the entire documents. This differs from Freitag et al. (2023), who removed some paragraphs from the ends of long documents. In our analysis we only consider segments overlapping with both previous annotation collections, thus obtaining 2,027 annotations evaluated across all protocols. We note that the subsampling makes all annotations protocols statistically less powerful, but keeps them fair in terms of statistical power per evaluated segment. Therefore, clustering and final system ranking in our analysis differs from Kocmi et al. (2023); Freitag et al. (2023). To keep the study comparable with Kocmi et al. (2023), we use the Wilcoxon rank-sum test with $p<0.05$ when producing system clusters. However, as all systems are evaluated on the same set of segments, we advice to use Wilcoxon signed-rank test when employing ESA as proposed by Kocmi et al. (2021).

To analyze inter annotator agreement, we run ESA protocol twice with different sets of annotators. We hired 28 annotators to evaluate our protocols and each protocol was evaluated by different sets of experts to avoid bias. Specifically, we had 8 bilingual annotators for the initial run of $\text{ESA}_1$, 10 translators for $\text{ESA}_2$ (different vendor), and 10 annotators to evaluate MQM protocol. For MQM, we hired professionals already experienced with MQM annotation protocol, while for ESA, we hired translators or bilingual speakers. All of them were native speakers of the target language, German.

**Highlighting errors:**

　Highlight the text fragment where you have identified a translation error (drag or click start & end).

　Click repeatedly on the highlighted fragment to increase its severity level or to remove the selection.

　**Minor Severity:** Style/grammar/lexical choice could be better/more natural.

　**Major Severity:** Seriously changed meaning, difficult to read, decreases usability.

　If something is missing from the text, mark it as an error on the **[MISSING]** word.

　The highlights do not have to have character-level precision. It's sufficient if you highlight the word or rough area where the error appears.

　Each error should have a separate highlight.

**Score:** After highlighting all errors, please set the overall segment translation scores. The quality levels associated with numerical scores on the slider:

　**0: No meaning preserved**: Nearly all information is lost in the translation.

　**33%: Some meaning preserved**: Some of the meaning is preserved but significant parts are missing. The narrative is hard to follow due to errors. Grammar may be poor.

　**66%: Most meaning preserved and few grammar mistakes**: The translation retains most of the meaning. It may have some grammar mistakes or minor inconsistencies.

　**100%: Perfect meaning and grammar**: The meaning and grammar of the translation is completely consistent with the source.

---

*Getting my eyes checked, because there was a decent deal for it here, and my left eye's been weird for some time now, so maybe they can tell me something what's happening.*

Ich lasse meine Augen untersuchen, denn es gab hier ein gutes Angebot dafür, und mein linkes Auge ist seit einiger Zeit komisch, also vielleicht können sie mir sagen, was los ist. **[MISSING]**

| 0%: No meaning preserved | 33%: Some meaning preserved | 66%: Most meaning preserved | 100%: Perfect |

---

*I'm splurging on a new set of frames, these red ones I reeeeally like.*

Ich verschwende mein Geld für eine neue Brillenfassung, diese rote mag ich seeehr. **[MISSING]**

| 0%: No meaning preserved | 33%: Some meaning preserved | 66%: Most meaning preserved | 100%: Perfect |

---

*So, apparently the ghost image my left eye sees isn't too much of a concern. But I need basically varifocals too. So the lenses have a combo of those. Not the full on varifocals but something in between those and "normal" glasses (no idea of the terminology in English).*

Also ist das Geistsrbild, das mein linkes Auge sieht, nicht so sehr Anlass zum Sorge. Aber grundsätzlich brauche ich auch Gleitsichtgläser. Die Linsen sind also eine Modifikation dieser. Nicht die starken Gleitsichtgläser, sondern etwa zwischen diesen und „normalen" Gläsern (keine Ahnung der Terminologie in Englisch). **[MISSING]**

| 0%: No meaning preserved | 33%: Some meaning preserved | 66%: Most meaning preserved | 100%: Perfect |

---

Figure 2: Screenshot of the beginning of one annotated document in the ESA interface (following segments are not shown). By showing and annotating whole documents at the segment-level, the annotators see all the relevant context. Segment *reset* button and *completed* labels removed for brevity. See the interactive tutorial shown to all annotators in Appendix Figure 10.

# 6   Analysis

We first analyze the collected data, system ranking, agreement with other protocols, quality assurance, and finally the annotation time. The findings reveal that the ESA quality is comparable, if not better than MQM, takes less time, and does not require highly trained annotators.

## 6.1   Score distribution

As per Table 1, on average the $ESA_1$ annotators mark 0.45 error spans per segments, which is close to MQM's 0.53 error spans per segment. The second run of $ESA_2$ has more than double of error spans per segment, which could be the result of different characteristics of annotators group, specifically experienced annotators in $ESA_1$ versus trans-

lators in $ESA_2$. On the other hand, $MQM^{WMT}$ contains 7x more errors per segment than our re-run of MQM. We attribute this difference primarily to the differences in annotation crowds, which further motivates our own evaluation of both MQM and ESA so that the annotations differ only in the annotation protocol and past MQM training. The severity levels are distributed similarly across campaigns. Important insights are in the score range distribution presented in Figure 3: the MQM-like score computation creates more skewed distribution around 0, which is in addition unbounded and can go to -infinity the longer the evaluated segment is, complicating modeling and comparisons. In contrast, the manual scores from annotators are spread out and guaranteed to be in [0, 100].

Figure 3: Distribution of scores for one annotation campaign. For ESA, we either use the manual score or ESA$_{spans}$ computation based on error severities. For MQMs, the distribution is clipped $\geq -15$ for higher resolution.

## 6.2 System ranking capabilities

We now investigate how well different annotation protocols (ESA and MQM) can rank MT systems. For purpose of this experiment, we consider MQM$^{WMT}$ which comes from an independent very high-quality annotation crowd and implementation, as the gold standard. Note that this creates a positive bias towards our implementation of MQM as opposed to ESA. In Appendix B.2 we show an evaluation without this gold standard assumption.

When comparing two protocols, we ideally want them to rank all systems in the same order. This is not always possible as some systems are very similar and cannot be significantly distinguished with the evaluated sample size. Another problem is that different protocols may weight different phenomena (e.g. fluency or adequacy) differently. To compare different protocols in the task of ranking systems, we use pairwise accuracy (Kocmi et al., 2021), which is also used in WMT Metrics shared task when comparing different automatic metrics (Freitag et al., 2023). Pairwise accuracy measures how many system pairs does a protocol rank the same way as MQM$^{WMT}$. As we have only 78 system pairs, any wrong system pair will change pairwise accuracy by 1.28%. Therefore, we also calculate Spearman's correlations as we mainly want protocols to have monotonic ranking.

In Figure 4, each subplot compares system-level scores between one protocol on x-axis and MQM$^{WMT}$ on y-axis. Our repeated MQM experiment and ESA protocol rank systems identically

(94.9%), while ESA has slightly higher Spearman's correlation with MQM$^{WMT}$. On the other hand, DA+SQM$^{WMT}$ significantly lacks behind both protocols. This suggest that our ESA protocol has comparable system ranking capabilities to MQM and is superior to DA+SQM$^{WMT}$.

The Figure 4 also shows, that relying on error spans only is not optimal, as ESA$_{spans}$ has lower accuracy and Spearman's correlation than ESA. We can notice that this is even lower than our rerun of MQM. This can be attributed to the evaluation crowd, where we used professional MQM annotators for MQM protocol, while we used bilingual speakers and translators for the ESA protocol.

Further focusing on the clustering, MQM$^{WMT}$ significantly differentiates the top system from others (highlighted in orange), while DA+SQM$^{WMT}$ strongly puts this system into the second cluster. This system is human reference, which we assume should be of highest quality. The reduced number of clusters in contrast to Figure 6 is due to the lower sample size. This may be result of DA+SQM$^{WMT}$ higher sensitivity to fluency and style errors, which contribute to 60% of all errors in human reference as marked by MQM$^{WMT}$. This conflict in clustering is one of the critiques of DA+SQM$^{WMT}$ if we assume that human reference should be the highest scoring translation. For example, in Kocmi et al. (2023), human reference was the best translation only in 2 out of 8 language pairs.

|            | MQM$^{WMT}$ |
|------------|-------------|
| ESA$_1$             | 0.227 |
| ESA$_{1\,spans}$    | 0.170 |
| ESA$_2$             | 0.250 |
| ESA$_{2\,spans}$    | 0.236 |
| MQM                 | 0.189 |
| DA+SQM$^{WMT}$      | 0.209 |

Table 2: Kendall $\tau$ segment-level correlations between evaluation protocols.

## 6.3 Agreement with other protocols

We now compare how different protocols correspond on the segment-level. We use MQM$^{WMT}$ as the gold standard to compare against because it was done independently outside of our setup and with high quality assurance. We analyze two aspects: (1) segment scores, and (2) spans, where we consider spans overlapping even with a single

Figure 4: Each point represents a system, with the original MQM$^{\text{WMT}}$ scores on the *y*-axis plotted against our rerun of DA+SQM$^{\text{WMT}}$ (first plot), ESA (second plot), ESA$_{\text{spans}}$ (third plot), and MQM (forth plot). Stripped lines indicate cluster separations determined by each method with alpha threshold 0.05. We compute Spearman correlation $\rho$ and pairwise accuracy Acc.

character as a match. For this evaluation, we use Kendall $\tau$ variant C, which is more suitable for data with different underlying scales and many ties.

In Table 2, we see that although all protocols correlate similarly with MQM$^{\text{WMT}}$, our protocols obtain the highest $\tau$ for both runs. The segment-level correlation also confirms that relying only on the error spans is not optimal and ESA$_{\text{spans}}$ obtains lower Kendall score.

In Table 3 we focus on the annotated spans. We consider any spans that overlap as matching, irrespective of severity. Because different protocols have different average number of error spans, presenting just the size of the intersection would be misleading. Instead, we show normalized set similarity that is not symmetric. It answers the questions: *What proportion of samples in B were covered by A?* Both ESA and MQM cover MQM$^{\text{WMT}}$ similarly, with 29% and 32% respectively. At the same time, MQM$^{\text{WMT}}$ covers ESA and MQM with 93% and 85%.

| ↓A  B→ | ESA | MQM | MQM$^{\text{WMT}}$ |
|---|---|---|---|
| ESA | | 77% | 29% |
| MQM | 85% | | 32% |
| MQM$^{\text{WMT}}$ | 93% | 89% | |
| $|A \cap B|/|B|$ | | (100% for $B = \emptyset$) | |

Table 3: Similarity between spans of different annotation protocols (one campaign) computed as percentage how much of $B$ does $A$ contain. For example, 93% of ESA spans were also in MQM. Any span that overlaps with another one is considered a hit, even if only with as single character.

## 6.4 Quality of annotations

**Intra annotator agreement.** We want the human evaluation protocol to consistently assign similar scores for the same translations over time. A good

indicator of the annotation quality is how noisy and subjective it is, which is reflected by how much annotators agree on the same segments (inter annotator agreement) as well as how a single annotator agrees with themselves (intra annotator agreement). To measure intra-AA, we ask the same annotators to again annotate the same documents two months later. We prepare an identical campaign with the same distribution of systems in the same order as originally, asking the same annotators to redo it again for both ESA and MQM.

It is not obvious how to measure the agreements for protocols that have different features. One issue is the frequency of ties, where MQM has more ties than rating from ESA or DA+SQM$^{\text{WMT}}$. For example, MQM$^{\text{WMT}}$ contains 30.8% no-errors, while DA+SQM$^{\text{WMT}}$ contains only 5.2% of score 100. Secondly, each protocol uses different range and distribution of scores, which makes calculation of agreement complicated. See App. Figure 7 to understand different distribution of scores.

Previous works comparing different protocols discretize the scale into bins (Graham et al., 2013; Freitag et al., 2021a), however, this approach is sensitive to subjective selection of bin sizes and benefits already discrete protocols. Instead, we propose to use Kendall's Tau-c correlation to measure inter-annotator agreement. Secondly, we also want to take into consideration that small changes in scores are less damaging than large shifts, therefore, we want inter annotator's scores to correlate linearly, ideally having identical score each time. To measure this, we use Pearson's correlation. Lastly, we measure recall of how often annotator mark *any* error in the same segment in contrast to leaving the segment without marked errors.

Table 4 shows that ESA has all scores higher than MQM. Higher Kendall and Pearson suggest that the task is easier for annotators to agree on the

score. We expect the recall to be comparable for both techniques as the task is similar, we hypothesize that the drop in MQM could be explained by annotators saving time and skipping minor errors as the annotation is more complicated for MQM than ESA, which can be confirmed when looking at minor error's recall only.

Lastly, evaluating inter annotator agreement is heavily affected by the strategy of annotators, where different annotation strategy does not mean different performance of the task as Riley et al. (2024) showed. Secondly, the MQM$^{WMT}$ was collected with a different tooling and the documents have been presented to annotators in different order, which could also impact the inter annotator agreement.

| | Intra AA | | Inter AA | |
|---|---|---|---|---|
| | ESA | MQM | ESA | MQM |
| Kendall's Tau-c | 0.149 | 0.109 | 0.254 | 0.116 |
| Pearson | 0.403 | 0.189 | 0.482 | 0.281 |
| Error recall | 69.6% | 61.9% | 66.6% | 40.1% |
| Minor e. recall | 70.7% | 66.2% | 67.7% | 44.4% |
| Major e. recall | 82.6% | 82.1% | 84.8% | 62.9% |

Table 4: Intra- and inter-annotator agreement on segment-level.

**Inter annotator agreement.** To measure how different annotators agree between themselves on the same protocol, we compare $ESA_1$ to $ESA_2$, where each protocol was evaluated by different group of annotators (bilingual annotators vs. translators). To calculate MQM's inter-annotator agreement, we compare our MQM run with MQM$^{WMT}$. However, the comparison is not as 1:1 as for ESA protocol. Our MQM was collected with different interface and system outputs have been shown to annotators in a different order. Results in Table 4 suggest, that ESA has also higher inter-annotator agreement than MQM. Unfortunately, we could not rerun DA+SQM protocol to calculate agreements, which needs to be reevaluated in future work.

**Agreement on the error span.** We now investigate how much MQM annotators agree on the error spans, error categories and severity levels. We evaluate from two angles: *intra*, where we check if the same annotator marks the same error spans or have overlapping parts, and *inter*, where we compare our MQM error spans to MQM$^{WMT}$. Table 5 shows that only in 30% of cases, the same annotator marked at least part of the same segment as an

error regardless error severity and category. If we look at cases preserving severity and category, this number drops to 8.4%.

When comparing the inter annotator agreement, only 50% of errors are overlapping. This number is higher as the total number of errors in MQM$^{WMT}$ is $7\times$ higher than in MQM, therefore it is more likely an error will have overlap.

| | Intra AA | Inter AA |
|---|---|---|
| Any errors | 29.3% | 50.2% |
| Same severity | 16.9% | 23.7% |
| Same category | 18.9% | 24.1% |
| Same sev. + categ. | 11.6% | 10.0% |
| Same sev. + subcateg. | 8.4% | - |

Table 5: ESA intra- and inter-annotator agreement (frequency) on marking overlapping errors with same severities, categories or subcategories.

**Quality control.** To measure the quality of annotations, we added "attention checks" in the form of segments for which we can reliably check whether the annotator annotated them correctly or not. In random documents, we perturbed the translation by replacing part of it with random sequence of words of the same length introducing a major translation error. Within 100 segments annotated by the annotator, they see both original and perturbed versions of that document. We can control the annotation quality by checking if the perturbed documents received more error spans. We show a worked-out Example 2.

**SRC***: Sie haben gestern das Treffen wieder verschoben.*
**TGT***: He postponed the meeting again yesterday.*
**TGT$^P$***: He postponed the meeting squirrels tense.*

Example 2: An example of a perturbed translation **TGT$^P$** based on the original system translation **TGT**. The **TGT** has one error (*He* should be *You* or *They*) and **TGT$^P$** introduces one more errors (*squirrels tense*).

The MQM and ESA setups used the same perturbations and we show their results in Table 6. The scores comparing original and perturbed segments for ESA and MQM are vastly different,[3] showing that annotators paid attention to the quality control items. For ESA the original segment had a higher score than the perturbed one in 86% of cases, while for MQM in 78% of cases. When investigating

---

[3]The protocols use different scales, for example, one major error under MQM is $-5$ points, while 25 points in ESA represents quarter of the full scale.

whether an annotator marked the error span or not, MQM has higher recall than ESA, however, this is less crucial for ESA as annotator can adjust the ranking without marking the error.

|  |  | Original | Perturbed | OK |
|---|---|---|---|---|
| **ESA** | Score | 79.5 | 52.6 | 86% |
|  | Span count | 0.85 | 1.86 | 54% |
|  | Perturbation marked |  |  | 56% |
| **MQM** | Score | -1.87 | -6.49 | 78% |
|  | Span count | 0.66 | 1.68 | 70% |
|  | Perturbation marked |  |  | 76% |

Table 6: Annotations assigned to perturbed attention check items (either scores or number of spans). **OK** is percentage in how many cases the non-perturbed item received a higher score or had fewer error spans, and how often the perturbed span was marked by the annotator.

## 6.5 Annotation time

One of the reasons behind development of the new protocol is reducing the time requirements of the evaluation. In this section, we analyze times on our experiments only, therefore our rerun of MQM in contrast to ESA. We do not include the times for DA+SQM and MQM$^{WMT}$ because they were done independently outside of our study and the time data are not available. Assessing the speed of annotations is challenging as annotators took breaks during the annotation (from short breaks taking several minutes up to several hours). This makes the evaluation of time problematic and we therefore investigate the time estimate in several ways. Counting all annotators together, the median time for annotation per single paragraph for MQM is 38 seconds and for ESA is 29 seconds, a reduction of 23%. However, as median time for each annotator fluctuates, we look at the average median time across annotators. **For MQM, the median is 49 seconds and for ESA it is 34 seconds, a reduction of 32%**. This can be contributed mainly to the less demanding error span annotation approach used in ESA.

We note that further speedups could be made by instructing the ESA annotators to not spend extra time e.g. marking multiple error-span annotations of a single grammatical phenomenon.

**Speedup during annotations.** Naturally, the reported total annotation time does not distinguish between the duration of the first and last annotations. In practice, annotators *learn* to perform the

annotation task more effectively. In Figure 5 we show time per segment depending on how many segments the annotator already processed. For both MQM and ESA, there is a small learning effect. For MQM, with each segment, the annotator becomes 0.20s faster, while for ESA this is 0.17s.



Figure 5: Time per segment with respect to progression in the annotation. The faint gray lines represent individual annotators, while the bold black line shows the average time. The lines are smoothed with a window of size 15 segments. We also compute the average speed at the beginning and at the end, which yields the *learned speedup*. This is how much the annotator speeds up after working on one segment.

## 7 Conclusion

Existing annotation protocols for machine translation evaluation are either expensive because they require expert labor (MQM), or they are noisy and less reliable (DA+SQM). To this end, we propose, describe, and analyze **Error Span Annotation (ESA)**, which builds on top of previous protocols to enable economic evaluation at scale. It works by asking the annotators to mark error spans, but with only the error severities and not types. In contrast to MQM, we also solicit final translation score. This is more reliable than DA+SQM, as the annotators are primed and informed about the translation errors to assess quality of longer documents.

We showed that our protocol has the higher inter and intra annotator agreement than MQM while being 32% faster. In addition, the protocol does not require annotators trained in MQM categorizations.

Lastly, we showed that relying only on error spans and not using the ranking score as we did in ESA$_{spans}$ produces suboptimal scoring, therefore the combination of error spans and ranking seems to produce the best results.

## Limitations

A possible limitation in contrast to MQM is that now the system evaluation does not provide breakdown of error types, which could help practitioners in improving their systems. ESA does not provide this because the goal is *evaluation* and not *diagnosis*. Because of this and the costs of scaling expert labor, we are convinced that this is not a true shortcoming of ESA. Furthermore, the annotated errors can be further classified and analysed, if necessary.

Our experiments, for monetary reasons, were done only on one language pair, English→German. Nevertheless, it is unlikely that the results would be vastly different for other languages. The most difficult setup could be with Chinese, Japanese, and Korean texts that do not use spaces. However, we made a deliberate decision to allow highlighting of individual characters, as opposed to words, so that the user experience is unified across all languages. This was done in spite of speed improvements (selecting on the word level is easier than selecting individual character boundaries) in order to make the tool scalable to a large range of languages.

## Ethics Statement

The annotators were paid a standard commercial translator wage in the respective country. The experts in the MQM annotation has been paid double the hourly wage. No personal data was collected and the showed data was screened for potentially disturbing content.

We follow up with a questionnaire asking annotators on their feedback. Almost all annotators specified that the annotation experience was positive and instructions were clear. The main concern they mentioned was that some documents have been too long to evaluate.

## References

ALPAC. 1966. Language and machines. Computers in translation and linguistics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. Findings of the 2017 conference on machine translation (WMT17). In *Proceedings of the Second Conference on Machine Translation*, 169–214. Association for Computational Linguistics.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, 131–198. Association for Computational Linguistics.

Aljoscha Burchardt. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*. Aslib.

Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2007. (meta-) evaluation of machine translation. In *Proceedings of the Second Workshop on Statistical Machine Translation*, 136–158. Association for Computational Linguistics.

Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. 2008. Further meta-evaluation of machine translation. In *Proceedings of the Third Workshop on Statistical Machine Translation*, 70–106. Association for Computational Linguistics.

Sheila Castilho. 2020. On the same page? comparing inter-annotator agreement in sentence and document level human machine translation evaluation. In *Proceedings of the Fifth Conference on Machine Translation*, 1150–1159. Association for Computational Linguistics.

Christian Federmann. 2018. Appraise evaluation framework for machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, 86–88. Association for Computational Linguistics.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, 578–628. Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021b. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human

evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, 733–774. Association for Computational Linguistics.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, 33–41. Association for Computational Linguistics.

Filip Klubička, Antonio Toral, and Víctor M. Sánchez-Cartagena. 2018. Quantitative Fine-grained Human Evaluation of Machine Translation Systems: A Case Study on English to Croatian. *Machine Translation*, 32(3):195–215.

Rebecca Knowles and Chi-kiu Lo. 2024. Calibration and context in human evaluation of machine translation. *Natural Language Processing*, 1–25.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, 1–42. Association for Computational Linguistics.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, 1–45. Association for Computational Linguistics.

Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In *Proceedings of the Sixth Conference on Machine Translation*, 478–494. Association for Computational Linguistics.

Philipp Koehn and Christof Monz. 2006. Manual and automatic evaluation of machine translation between European languages. In *Proceedings on the Workshop on Statistical Machine Translation*, 102–121. Association for Computational Linguistics.

Julia Kreutzer, Nathaniel Berger, and Stefan Riezler. 2020. Correct me if you can: Learning from error corrections and markings. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, 135–144. European Association for Machine Translation.

Arle Lommel, Aljoscha Burchardt, Maja Popović, Kim Harris, Eleftherios Avramidis, and Hans Uszkoreit. 2014. Using a new analytic measure for the annotation and analysis of MT errors on real data. In *Proceedings of the 17th Annual Conference of the European Association for Machine Translation*, 165–172. European Association for Machine Translation.

Marianna Martindale and Marine Carpuat. 2018. Fluency over adequacy: A pilot study in measuring user trust in imperfect MT. In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, 13–25. Association for Machine Translation in the Americas.

Maja Popović. 2020. Informative manual evaluation of machine translation output. In *Proceedings of the 28th International Conference on Computational Linguistics*, 5059–5069. International Committee on Computational Linguistics.

Parker Riley, Daniel Deutsch, George Foster, Viresh Ratnakar, Ali Dabirmoghaddam, and Markus Freitag. 2024. Finding replicable human evaluations via stable ranking probability. *arXiv preprint arXiv:2404.01474*.

Paul Rottger, Bertie Vidgen, Dirk Hovy, and Janet Pierrehumbert. 2022. Two contrasting data annotation paradigms for subjective NLP tasks. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 175–190. Association for Computational Linguistics.

David Vilar, Gregor Leusch, Hermann Ney, and Rafael E. Banchs. 2007. Human evaluation of machine translation through binary system comparisons. In *Proceedings of the Second Workshop on Statistical Machine Translation*, 96–103. Association for Computational Linguistics.

David Vilar, Jia Xu, Luis Fernando D'Haro, and Hermann Ney. 2006. Error analysis of statistical machine translation output. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC'06)*. European Language Resources Association (ELRA).

Johnny Wei, Tom Kocmi, and Christian Federmann. 2022. Searching for a higher power in the human evaluation of MT. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, 129–139. Association for Computational Linguistics.

John S. White, Theresa A. O'Connell, and Francis E. O'Mara. 1994. The ARPA MT evaluation methodologies: Evolution, lessons, and future approaches. In *Proceedings of the First Conference of the Association for Machine Translation in the Americas*.

Vilém Zouhar, Tom Kocmi, and Mrinmaya Sachan. 2024. AI-assisted human evaluation of machine translation.

# A  User Guidelines

The following are annotation guidelines for our local ESA and MQM campaigns.

## A.1  ESA (Error Span Annotations)

**Higlighting errors:**  Highlight the text fragment where you have identified a translation error (drag or click start & end). Click repeatedly on the highlighted fragment to increase its severity level or to remove the selection.

- **Minor Severity:** Style/grammar/lexical choice could be better/more natural.
- **Major Severity:** Seriously changed meaning, difficult to read, decreases usability.

If something is missing from the text, mark it as an error on the **[MISSING]** word. The highlights do not have to have character-level precision. It's sufficient if you highlight the word or rough area where the error appears. Each error should have a separate highlight.

**Score:**  After highlighting all errors, please set the overall segment translation scores. The quality levels associated with numerical scores on the slider:

- **0%**: No meaning preserved: Nearly all information is lost in the translation.
- **33%**: Some meaning preserved: Some of the meaning is preserved but significant parts are missing. The narrative is hard to follow due to errors. Grammar may be poor.
- **66%**: Most meaning preserved and few grammar mistakes: The translation retains most of the meaning. It may have some grammar mistakes or minor inconsistencies.
- **100%**: Perfect meaning and grammar: The meaning and grammar of the translation is completely consistent with the source.

## A.2  MQM (Multidimensional Quality Metrics)

**Higlighting errors:**  Highlight the text fragment where you have identified a translation error (drag or click start & end). Click repeatedly on the highlighted fragment to increase its severity level or to remove the selection.

- **Minor Severity**: Style/grammar/lexical choice could be better/more natural.
- **Major Severity**: Seriously changed meaning, difficult to read, decreases usability.

If something is missing from the text, mark it as an error on the **[MISSING]** word. The highlights do not have to have character-level precision. It's sufficient if you highlight the word or rough area where the error appears. Each error should have a separate highlight.

**Error types:**  After highlighting an error fragment, you will be asked to select the specific error type (main category and subcategory). If you are unsure about which errors fall under which categories, please consult the typology definitions.



Figure 6: System scores illustrating differences between DA+SQM and MQM. Each point is a single system and dashed lines mark clusters. DA+SQM produces fewer clusters and groups many systems into one single cluster, while MQM better distinguishes different systems. Scores and clusters are from Kocmi et al. (2023).



Figure 7: Intra annotator agreement; changes in scoring by the same annotator when evaluated again. Each point represents an annotated segment with x-axis being annotator's score assigned in March and y-axis their score assigned in May.

| Feature | Corr. with ESA score ($\rho$) |
|---|---|
| Source token count | -0.16 |
| Target token count | -0.06 |
| Minor error count | -0.20 |
| Major error count | -0.52 |
| Missing error count | -0.45 |
| Minor error count (normalized) | -0.13 |
| Major error count (normalized) | -0.37 |
| Missing error count (normalized) | -0.31 |

Table 7: Segment-level Pearson correlation of individual features with the ESA score.

## B   Additional results

### B.1   From error spans to final score

To find out what influences the score, we show correlation between individual segment-level features in Table 7. On average, longer segments have lower translation quality. Importantly, the error counts normalized by segment length correlate less than the non-normalized counterparts. However, we note that the normalized scores are more continuous that the non-normalized MQM computation, as per Figure 3.

The MQM formula was crafted with respect to preserving system ranking and not segment-level matching from the same annotator (Freitag et al., 2021a). However, the construction of ESA allows us to revisit this problem as each annotator gives both the error spans and the final score. We scan for multiple minor/major ratios of error weights and show the results in Figure 8. We find that the optimal formula that optimizes the correlation between the direct score and the score from the spans has the following form: $\text{SEG.SCORE} = -1 \cdot \#\text{MINOR} - 4.8 \cdot \#\text{MAJOR}$, which is very close to the originally proposed 1:5 ratio. From Figure 8, the weight of the major error class seems to have a much bigger effect on the final translation score, suggesting that minor errors play a lesser role.



Figure 8: Correlation between direct ESA score and scores computed from error spans with minor errors having weight $-1$ and major errors $-x$.

### B.2   Protocol evaluation without a gold standard

Evaluating the quality of a protocol without a target to compare to is difficult. In previous sections we assumed that $\text{MQM}^{\text{WMT}}$ is the gold standard, which might bias the evaluation in favor of MQM. Even though the results showed higher correlations of ESA with $\text{MQM}^{\text{WMT}}$. For completeness, we consider the methodology of Zouhar et al. (2024)

which does not require target gold standard to compare the quality of annotation protocols.

The assumption is that annotation protocols have various levels of noise, but are unbiased in what they measure. Because MQM and DA+SQM might measure different things, we want to compare each to the perfect ranking of the particular thing they aim to measure. The linking hypothesis is that even noisy and low-quality annotation protocols would lead to the final system ranking with large enough data. Vice-versa, only robust annotation protocols would arrive at the final ranking with only a few data points per each system. This is formalized by the subset consistency accuracy. It is the system ranking accuracy on a subset of annotations with respect to the ranking induced by the full data from one annotation protocol.

We show the results in Figure 9. Out of the comparable lines, ESA (with direct scoring) achieves the highest subset consistency accuracy. In practice, this translates to needing fewer annotated examples to achieve the final system ranking. This directly corresponds to lower annotation costs.



Figure 9: Subset consistency accuracy (Zouhar et al., 2024) of annotation protocols. E.g. with just 60 annotated segments, $\text{MQM}^{\text{WMT}}$ achieves 95% system ranking accuracy with its final ranking based on 160 annotated segments. Values in the legend are averages, corresponding to normalized area under the curve. The only comparable lines are ESA, $\text{ESA}_{\text{spans}}$, and $\text{MQM}_{\text{spans}}$ because they were run in the same setting with similar crowds.

Figure 10: Tutorial to ESA annotations shown at the beginning of the campaign. All tutorial segments need to be annotated correctly before continuing.

# Neural Methods for Aligning Large-Scale Parallel Corpora from the Web for South and East Asian Languages

**Philipp Koehn**

Center for Language and Speech Processing
Johns Hopkins University
`phi@jhu.edu`

## Abstract

We introduce neural methods and a toxicity filtering step to the hierarchical web mining approach of Paracrawl (Bañón et al., 2020), showing large improvements. We apply these methods to web-scale parallel corpus mining for 9 South and East Asian national languages, creating training resources for machine translation that yield better translation quality for most of these languages than existing publicly available datasets in OPUS. Our methods also generally lead to better results than the global mining approach of Schwenk et al. (2021).

## 1 Introduction

The goal of this work is to apply neural methods to the task of parallel corpus mining from the web and to create large useful parallel corpora for languages that have not received much attention. We demonstrate when applying these methods at scale, they yield better data resources than the two main existing approaches Paracrawl (Bañón et al., 2020) and CC-Matrix (Schwenk et al., 2021).

In addition to six Southeast Asian national languages (Burmese, Thai, Lao, Khmer, Vietnamese, Indonesian), we also included the South Asian languages Hindi and Nepali and the East Asian language Korean. These are mostly mid-resource languages, they have millions of speakers, mostly significant presence on the web, but have not received as much attention in the research community as European languages (Bañón et al., 2020), Indian languages (except, we also include Hindi) (Siripragada et al., 2020), Chinese (Ziemski et al., 2016; Zhai et al., 2020), and Japanese (Morishita et al., 2022).

Building on the work of the Paracrawl project (Bañón et al., 2020), we follow the same general sequence of steps: targeted web crawling, document alignment, sentence alignment, and parallel corpus filtering. Note that compared to the European-



Figure 1: National languages covered: Hindi, Nepali, Burmese, Thai, Lao, Khmer, Vietnamese, Indonesian, Korean. We build parallel corpora for these languages paired with English.

focused Paracrawl project, we deal with languages with fewer existing resources, mostly non-Latin scripts, and challenges such as lack of explicit word segmentation and even sentence boundary marking (in the case of Thai).

In contrast to Paracrawl, we deploy neural methods in three steps: document alignment with an efficient Marian (Junczys-Dowmunt et al., 2018) neural machine translation model distilled from the multilingual NLLB (NLLB Team et al., 2022) model, sentence alignment with Vecalign (Thompson and Koehn, 2019), and using LASER for parallel corpus filtering (Chaudhary et al., 2019). We also added a novel toxity filtering step.

We obtain large parallel corpora of 1.5–7.7 million sentence pairs per language. We validate the usefulness of these corpora by showing better machine translation quality of up to +18.2 BLEU compared to CC-Matrix (Schwenk et al., 2021) for 7 languages and up to +13.0 BLEU compared to other existing parallel corpora on OPUS[1] (Tiede-

---

[1] `https://opus.nlpl.eu/`

mann, 2009) for 6 languages (tied for another language). While this required significant computational resources, the effort was carried out using only CPUs and consumer-grade GPUs (GTX 1080ti).

## 2 Related Work

While the idea of mining the web for parallel data has been already pursued in the 20th century (Resnik, 1999), the initial large-scale efforts were limited to large companies such as Google (Uszkoreit et al., 2010) and Microsoft (Rarrick et al., 2011), or targeted efforts on specific domains such as the Canadian Hansards and Europarl (Koehn, 2005). More recently, large corpora have been released by broad web mining efforts, such as Paracrawl (Bañón et al., 2020) and CC-Matrix (Schwenk et al., 2021). A recent effort to assemble large-scale monolingual and parallel corpora is the EU Project High Performance Language Technologies (Aulamo et al., 2023).

Currently, there are two main approaches to extract parallel sentence pairs from web documents: hierarchical and global mining. In *hierarchical mining* (as in Paracrawl), the task is broken up into the steps of identifying websites with parallel text, document alignment within websites, sentence alignment within document pairs, and sentence pair filtering.

In contrast, in *global mining* (as in CC-Matrix), all content is split up into sentences, each sentence represented by a cross-lingual sentence embedding and stored in one index per language. Then, sentences in one language are used to query the index of sentences in another language, using nearest neighbor search. There are also efforts that lie in-between these two extremes, such as local mining in CC-Align (El-Kishky et al., 2020) where the hierarchical mining is followed up to the step of document alignment, and then sentences for each document are stored in an index and then queried regardless of the order of sentences in the document.

We follow the hierarchical mining approach. We believe that it leads to cleaner parallel corpora since it matches alignment with the underlying structure of the data. There has been varying amount of work on the steps in hierarchical mining. Matching documents pairs uses some similarity measure to compare the content of documents across languages. A common approach is to translate the non-English document into English and perform monolingual matching of words (Buck and Koehn, 2016) or n-grams (Dara and Lin, 2016; Uszkoreit et al., 2010). There have been some attempts to use document embeddings (Guo et al., 2019). Besides matching the URL (Le et al., 2016; El-Kishky et al., 2020) — e.g., `example.com/en/page.html` and `example.com/fr/page.html` — other structural information such the DOM-tree (Shi et al., 2006), links to the same images, links between pages, etc. have been rarely used.

Sentence alignment has been a rich field of research dating back to the 1990s (Brown et al., 1991; Gale and Church, 1993). This also requires a similarity measure, defined over sentences or sequences of sentences. Typical features are sentence length and matches in a bilingual dictionary (Moore, 2002; Varga et al., 2005). Sennrich and Volk (2010) translate the non-English sentence and match the translation against the English sentence using the BLEU score. Vecalign (Thompson and Koehn, 2019) is a sentence alignment method that relies on bilingual sentence embeddings and achieves linear run time with a coarse-to-fine dynamic programming algorithm.

Finally, a lot of effort has been spent on developing methods for filtering noisy parallel corpora which are particularly harmful for neural models (Khayrallah and Koehn, 2018). Four shared tasks were dedicated to this problem (Koehn et al., 2018, 2019, 2020; Sloto et al., 2023). Besides basic simple filtering rules based on sentence or token length and their ratios (Kurfalı and Östling, 2019; Soares and Costa-jussà, 2019), typically a scoring function is used. Popular methods are based on the scores obtained by force-decoding the sentence pair with a machine translation model (Junczys-Dowmunt, 2018), and the cosine distance between cross-lingual sentence embeddings (Chaudhary et al., 2019). Recently, the most successful approach are classifiers that distinguish between genuine parallel sentence pair and misalignments, typically based on neural sentence representations (Açarçiçek et al., 2020; Esplà-Gomis et al., 2020; Xu et al., 2020; Tan et al., 2023).

Filtering has been focused on impact on machine translation quality using traditional metrics. There has not been much published work on toxicity filtering (NLLB Team et al., 2022) — a task that is also hard to delineate and evaluate.

| Model | Vietnamese | | | Nepali | | | Thai | | |
|---|---|---|---|---|---|---|---|---|---|
| | time | chrF | BLEU | time | chrF | BLEU | time | chrF | BLEU |
| MoE 54b official | - | 62.3 | 43.8 | - | 66.9 | 48.1 | - | 57.8 | 36.9 |
| Dense 3b official | - | 61.5 | | - | 65.9 | | - | 56.8 | |
| Dense 1b official | - | 59.8 | | - | 64.5 | | - | 54.9 | |
| Dense distilled 1b official | - | 60.4 | | - | 65.1 | | - | 54.9 | |
| Dense distilled 600m official | - | 62.3 | | - | 62.5 | | - | 52.7 | |
| Dense 3b quantized | 207s | 60.7 | 41.1 | 202s | 62.2 | 41.1 | 238s | 55.4 | 33.4 |
| Dense distilled 1b quantized | 61s | 59.5 | 39.6 | 71s | 63.2 | 42.2 | 74s | 53.8 | 31.2 |
| **Dense distilled 1b** | 45s | 59.8 | 39.2 | 44s | 63.7 | 42.7 | 50s | 54.1 | 31.5 |
| Dense 1b | 45s | 58.9 | 38.6 | 44s | 62.4 | 41.5 | 51s | 54.2 | 31.6 |

Table 1: Speed/Quality trade-offs for different versions of NLLB, the model we distill. Translation time to translate the 1012 sentences of the Flores devtest set into English on a single GTX-1080 GPU (bottom). Official NLLB evaluations are in the top of the table. Based on these findings, we use the dense distilled 1 billion parameter model.

## 3 Methods

### 3.1 Targeted Crawling

We follow the Paracrawl approach of crawling a list of targeted web sites. The crawl list has been mainly obtained by using meta-data from CommonCrawl but also opportunistically extended over several years, e.g., by web searches for language-specific terms. Based on Commoncrawl statistics, any website that has pages in English and any of the targeted languages and somewhat balanced ratio was selected and crawled with httrack[2], an open source web copying tool. We only follow links to web pages on the same webdomain. We stop crawling after crawling 50,000 pages for each website, both to avoid downloading duplicate webpages and due to computational limitations of subsequent processing steps.

### 3.2 Distilling Machine Translation Models

Our document alignment approach requires the translation of all non-English web pages for a targeted language into English. Since this implies the translation of a massive volume of text, we need an efficient but still sufficiently high-quality machine translation model.

The multilingual machine translation model NLLB (NLLB Team et al., 2022) covers 200 languages, including all the languages we target here. It comes in versions with 600 million to 54 billion parameters. However, using even the smallest model would be computationally prohibitive given the scale of our effort and the limitations of our technical means. Hence, we decided to distill these



Figure 2: Amount of synthesized training data from the NLLB model and BLEU scores of distilled Marian models. For Lao, Khmer, and Burmese, we exhausted the monolingual data in mC4.

models into an efficient model that can be run on CPU via data distillation. Specifically, we use the NLLB model to translate monolingual text and then use the resulting synthetic parallel corpus to train a faster model. The monolingual text for distillation is drawn from mC4[3] (Xue et al., 2021).

Table 1 shows machine translation quality scores and the time it takes to translate the 1012 sentences of the Flores-200 devtest set for three of our languages (Vietnamese, Nepali, and Thai) into English given different NLLB models. We explored the use of quantized parameters. However, we observed worse speed/quality trade-offs. We settled on using the dense distilled 1 billion parameter model. It

---

[2]available at https://www.httrack.com/

[3]available at https://huggingface.co/datasets/mc4

| Language | Forw. | Backw. | Both | +OPUS |
|----------|-------|--------|------|-------|
| Hindi | 36.4 | 31.1 | **36.2** | 36.3 |
| Nepali | 30.8 | 30.1 | **33.6** | 33.4 |
| Burmese | 21.3 | 18.3 | **22.7** | 21.6 |
| Thai | 23.9 | 18.4 | **25.3** | 23.5 |
| Lao | 24.9 | 21.8 | **28.4** | 27.6 |
| Khmer | 25.5 | 13.3 | 26.3 | **26.7** |
| Vietnamese | 30.9 | 27.5 | 32.4 | **34.8** |
| Indonesian | 41.0 | 37.6 | **41.3** | – |
| Korean | 26.0 | 22.6 | 26.0 | **26.5** |

Table 2: BLEU scores for different data types for distillation: synthetic corpus generated by forward translation (X→English) or back translation (English→X). Forward translation fares better than backward translation, but combination of both is typically best.
We also checked if we can better system by adding OPUS data. This is the case for Khmer, Vietnamese, and Korean, so we use these system in our mining pipeline.

gives reasonable performance at translation speeds of about 500 words per second on GPU.

We explored how much data we need to distill to get a reasonable Marian system. As illustrated in Figure 2, system quality plateaus at around 1 billion words of distilled data. Note that we exhausted all monolingual data in mC4 for Lao, Khmer, and Burmese, so we distilled less data for these.

We generate synthetic parallel corpora by translating both from English and into English. The forward direction (X→English) is motivated by the idea of data distillation while backward translation (English→X) is well-established in the field of machine translation since it builds on authentic text on the target side. As shown in Table 2, we find that forward translation gives better results, but combining both forward and backward translation fares generally best.

We filter the synthesized corpus with LASER using a threshold of 1.05 (1.00 for Burmese and Lao, unfiltered for Hindi). See Section 3.5 for more details on this method. We also added all of OPUS to the training of Khmer, Vietnamese, and Korean distilled models. As shown in Table 2, adding OPUS data yielded better translation quality.

The configuration of Marian (Junczys-Dowmunt et al., 2018) is given in Appendix A. The model is trained with guided alignment training and a vocabulary shortlist. The translation model uses quantized parameters for efficient vector integer computations supported by Intel CPUs (8 bit, avx512). When translating web content, we observe transla-

tion speeds of about 1000 words per second in a single Intel Xeon Silver 4110 CPU core. Contrast that to 500 words per second on a GPU for the NLLB model: a roughly thousand-fold increase in translation speed when measured by sentences per compute core.

### 3.3 Document Alignment

Our document aligner follows the method by Buck and Koehn (2016). For a website where we found web pages in English and in the targeted language, we translate all the latter web pages into English and represent each document (i.e., web page) in form of word counts. Document similarity is measured by tf/idf-weighted cosine distance between these representations. A greedy algorithm iteratively finds the best matching document pair and removes them from the pool of documents. The process terminates if documents in either language are exhausted.

The main difference to the Paracrawl approach is the use of a very efficient neural translation model instead of a statistical Moses model. The neural model has higher translation quality and is faster.

### 3.4 Sentence Alignment

We used Vecalign (Thompson and Koehn, 2019) as sentence aligner. It uses the cosine-distance between LASER embeddings with modified CSLS scoring (normalizing by distance to randomly chosen neighbors). It is also constrained by the order of the sentences in the pair of documents. Just like other sentence aligners (Hunalign, Bleualign, etc.), it may skip and merge sentences but it is not allowed to reorder them. Hence, it combines a powerful sentence matching method with the structural bias coming from the fact that documents are in almost all cases translated in sequence.

Documents were split into sentences with NLTK's sentence tokenizer. Thai required special treatment due to its lack of marking of sentence boundaries. We used the library pythainlp (Phatthiyaphaibun et al., 2023) for sentence splitting. We use LASER3, the latest version (Heffernan et al., 2022), that supports all our languages.

### 3.5 Noise Filtering

The previous processing steps are geared towards high recall instead of high precision. In other words, we try to retain as much data as possible. This requires a final filtering step that removes

| Language | 1.00 | | 1.05 | |
|---|---|---|---|---|
| | **Size** | **BLEU** | **Size** | **BLEU** |
| Hindi | 136.5m | 31.7 | 75.0m | **31.8** |
| Nepali | 32.3m | **26.0** | 18.3m | 25.0 |
| Burmese | 12.5m | **11.5** | 5.2m | 10.1 |
| Thai | 21.9m | **19.3** | 12.9m | 18.9 |
| Lao | 152.4m | **23.3** | 110.0m | 22.3 |
| Khmer | 22.6m | **13.7** | 6.4m | 9.4 |
| Vietnamese | 134.4m | 29.7 | 94.5m | **30.1** |
| Indonesian | 27.0m | 37.2 | 13.5m | **37.6** |
| Korean | 228.1m | 22.2 | 118.4m | **23.4** |

Table 3: Impact of different thresholds in LASER-based filtering: Corpus size in million words and BLEU score.

noisy data, an open problem that has received much research attention.

We use LASER-based filtering (Chaudhary et al., 2019), using LASER3 (Heffernan et al., 2022). This method embeds sentences in a cross-lingual embedding space, so that an English sentence and its translation should have identical representations. Hence, the distance between an English sentence embedding and a non-English sentence embedding is a measure for their meaning similarity. The exact formula to compute similarity between the two embedding vectors is the cosine distance, normalized by how similar each vector is to its closest neighbors in the embedding space.

We carried out limited experiments with the filtering threshold and chose a value of 1.00 for Nepali, Burmese, Thai, Lao, and Khmer and 1.05 for Hindi, Vietnamese, Indonesian, and Korean. We note that the more permissive threshold (1.00) worked better for the smaller corpora (see Table 3). For some languages we tried even lower thresholds but that led to worse results.

### 3.6 Toxicity Filtering

While we are aiming to collect parallel data across the entire web, we do want to exclude toxic content, so that machine translation systems are not trained to produce offensive language. We narrow down the concept of excluded toxic content to pornographic web sites which not only feature derogatory and offensive language but are also often machine translated.

Toxicity filtering may be carried at several levels. We argue that filtering on the level of web sites will lead to the most robust results. Simple key word filtering on the sentence level has to contend with the fact that many words are ambiguous, and excluding all sentences that have, say, the word *sex* in them would eliminate many respectable uses of that term.

Hence, we take a more nuanced view of offensive vocabulary. We use tf/idf scores to identify English vocabulary that is typical for websites that have the substring *porn* in their domain name. This yields words that are very frequently used on such web sites compared to full crawl for a language pair. We start with a list of 100 terms for each language pair, merge that list and curate it to remove, for instance, terms that refer to ethnicities (e.g., *Asian*). This list comprises 141 words.

Using this words list, we proceed to filter out websites. We compute the average tf/idf score across all the words for each website, and if it is above a certain threshold (we use 0.02), we eliminate all content from that website.

## 4 Corpora

### 4.1 Corpus Statistics

We apply the processing pipeline to 9 languages. Table 4 gives detailed statistics. The pipeline succeeded to process between 5,854 (Burmese) and 32,765 (Vietnamese) website crawls. A small proportion (about 10%) of the crawls are repeat crawls, i.e., they crawled the same website again at a later time, typically after several months or even years.

The next step is document alignment, resulting in 492,723 (Lao) to 7,758,116 (Korean) document pairs. Then comes sentence alignment, creating a raw corpus of 7,513,409 (Lao) to 128,828,741 (Korean) sentence pairs.

This corpus is filtered and deduplicated. We report how many good sentence pairs are retained when applying filtering to corpora from each crawl — which also includes deduplication: ranging from 605,959 (Khmer) to 11,014,387 (Korean). Then, deduplication is done again on the corpus combined across all crawls, reducing these numbers further to 420,824 (Khmer) to 8,298,299 (Korean). These numbers are based on a filtering threshold of 1.05. For five of the languages we saw better results with a filtering threshold of 1.00, so we report these numbers as well. For Khmer, this retains 1,507,135 sentence pairs.

Working back from the filtered data, we can check how many document pairs had sentence pairs that survived quality filtering. For instance, this is the case for 4,035,376 of the 7,758,116 Korean–

| Language | Crawls | | | Documents | | | Sentences | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | all | good | detox | all | good | detox | all | good | dedup | detox |
| Hindi | 13,605 | 10,900 | 10,348 | 4,033,751 | 2,453,234 | 2,361,953 | 52,919,986 | 5,989,651 | 4,823,444 | 4,712,564 |
| Nepali | 6,095 | 4,556 | 4,508 | 694,238 | 431,808 | 429,615 | 8,312,728 | 1,305,921 | 1,090,690 | 1,085,057 |
| ≥1.00 | | 5,136 | 5,074 | | 480,792 | 478,219 | | 2,706,360 | 2,254,055 | 2,243,954 |
| Burmese | 5,854 | 4,145 | 4,106 | 790,360 | 13,662 | 13,613 | 9,769,167 | 343,788 | 341,897 | 715,512 |
| ≥1.00 | | 4,817 | 4,760 | | 466,653 | 463,907 | | 2,002,212 | 1,674,072 | 1,666,530 |
| Thai | 14,012 | 11,131 | 10,556 | 3,349,364 | 1,409,191 | 1,357,692 | 61,466,936 | 1,470,556 | 1,190,997 | 1,176,111 |
| ≥1.00 | | 12,549 | 11,877 | | 2,232,342 | 2,152,042 | | 2,761,013 | 2,218,153 | 2,175,890 |
| Lao | 4,177 | 3,938 | 3,890 | 492,723 | 353,048 | 351,047 | 7,513,409 | 1,158,534 | 936,986 | 931,456 |
| ≥1.00 | | 4,019 | 3,971 | | 454,348 | 451,824 | | 2,391,972 | 2,004,028 | 1,994,053 |
| Khmer | 6,025 | 4,453 | 4,411 | 890,264 | 306,030 | 304,014 | 10,981,209 | 605,959 | 420,824 | 418,991 |
| ≥1.00 | | 5,102 | 5,048 | | 546,357 | 543,412 | | 1,884,419 | 1,507,135 | 1,501,304 |
| Vietnamese | 32,765 | 19,035 | 18,267 | 6,951,765 | 2,845,099 | 2,768,498 | 80,256,711 | 8,735,317 | 6,473,708 | 6,291,407 |
| Indonesian | 20,031 | 13,143 | 12,557 | 5,443,448 | 2,302,037 | 2,239,685 | 77,507,912 | 10,304,822 | 7,260,778 | 7,133,323 |
| Korean | 24,500 | 20,423 | 19,154 | 7,758,116 | 4,035,376 | 3,759,849 | 128,828,741 | 11,014,387 | 8,298,299 | 7,709,312 |

Table 4: Detailed statistics on the crawled datasets, in terms of number of crawls of websites, number of aligned document pairs, and sentence pairs. The numbers below *good* specify counts for these categories that have valid sentence pairs after LASER filtering with threshold 1.05 (extra rows for languages where we applied the threshold 1.00) and deduplication. For crawls and documents this number is inflated because the same good sentence pair may be in multiple documents and crawls. The deduplicated sentence pair count refers to a final global deduplication step. The table also reports these statistics after removing crawls due to toxic content.

| Language | Ours | CC-Matrix | OPUS |
|---|---|---|---|
| Hindi | 4.6m | 15.1m | 22.6m |
| Nepali | 2.2m | 19.6m | 1.9m |
| Burmese | 1.6m | 10.0m | 0.6m |
| Thai | 1.8m | – | 15.2m |
| Lao | 1.9m | 4.2m | 4.2m |
| Khmer | 1.3m | 5.9m | 0.6m |
| Vietnamese | 6.2m | 49.9m | 18.8m |
| Indonesian | 7.1m | 56.8m | 9.8m |
| Korean | 7.7m | 19.4m | 19.7m |

(a) Number of Segments

| Language | Ours | CC-Matrix | OPUS |
|---|---|---|---|
| Hindi | 74m | 196m | 296m |
| Nepali | 32m | 176m | 12m |
| Burmese | 28m | 102m | 8m |
| Thai | 22m | – | 152m |
| Lao | 27m | 40m | 40m |
| Khmer | 23m | 66m | 6m |
| Vietnamese | 93m | 780m | 211m |
| Indonesian | 109m | 624m | 88m |
| Korean | 114m | 205m | 151m |

(b) Number of English Words

Table 5: Size of parallel corpora, in millions, after length (≤80 words) and length ratio (≤9) filtering, compared to existing parallel data in OPUS (without CC-Matrix) and CC-Matrix.

English document pairs. Applying the same calculation for web crawls, 20,423 of the 24,500 Korean web crawls yielded at least one sentence pair in the final filtered corpus. Note that the number of crawls and documents after filtering is inflated because the same good sentence pair may be in multiple documents and crawls.

Finally, we remove toxic content from the corpus. This reduces only a small percentage of the data. The biggest reduction is for Korean–English, about 7%, from 8,298,299 to 7,709,312 sentence pairs.

## 4.2 Comparison to OPUS and CC-Matrix

We compare the size of the obtained corpora to pre-existing data sets in Table 5. We combined all corpora available in OPUS, the popular platform for parallel data. We separated out CC-Matrix (which is also available on OPUS) since it is the

method that is most similar to our approach and it is also typically the largest corpus on OPUS. CC-Matrix collected parallel sentences by matching sets of sentences from CommonCrawl solely based on the similarity of their LASER embeddings.

The table shows the number of segments and number of English words for each language. We count the number of English words because it is a consistent measure across all languages and counting words for languages like Thai is problematic due to the lack of word spacing. The numbers are computed after another filtering step typically done for translation: we remove sentences longer than 80 words and sentence pairs where one sentence has more than 9 times as many words as the other.

Note that the sizes of the obtained corpora are smaller than CC-Matrix and only for Nepali,

Figure 3: BLEU scores on neural machine translation systems build with our corpora, compared to existing corpora. We obtain better parallel corpora than anything previously existing for Nepali, Burmese, Lao, Khmer, Indonesian, and Korean, by a difference of +13.0, +8.0, +10.2, +3.1, +1.1 BLEU, respectively, compared to the better of CC-Matrix or OPUS (without CC-Matrix). Our Thai corpus matches OPUS, For Hindi and Vietnamese, existing corpora are better. CC-Matrix does not contain Thai.

Burmese, and Khmer bigger than what already exists in OPUS (excluding CC-Matrix). We obtain a larger Indonesian corpus than what exists in OPUS in terms of number of words but not in number of segments. Our smallest corpus is Khmer–English (1.3 million segment pairs, 23 million words), the largest corpus is Korean–English (8.2 million segment pairs, 118 million words). Note that CC-Matrix does not contain Thai.

## 5 Evaluation

Since our main motivation is to create parallel corpora for training machine translation systems, we evaluate them by training a system on each corpus and measuring each system's translation quality with spmBLEU (scarebleu -tok flores200) on Flores-200 (NLLB Team et al., 2022). We chose this test set and metric since they cover all our languages. Flores-200 comprises professional translations of English content drawn from Wikinews, Wikijunior, and Wikivoyage. We also computed scores with chrF++ which closely mirrors the spmBLEU results in terms of system ranking, so we do not report them here for sake of clarity.

Machine translation systems were trained using Marian (Junczys-Dowmunt et al., 2018) using the setup as for our distilled translation models (see

Section 3.2).

Results are shown in Figure 3. By our measure, we obtain better parallel corpora than anything previously existing for Nepali, Burmese, Lao, Khmer, Indonesian, and Korean, by a difference of +13.0, +8.0, +10.2, +3.1, +1.1 BLEU, respectively, compared to the better of CC-Matrix or OPUS. Our Thai–English corpus is as good as what is currently in OPUS ($\pm 0$). Only for Hindi and Vietnamese our data fares worse (–0.9 and –3.9 BLEU, respectively). We tried to investigate this discrepancy but did not gain any substantial insights.

It is worth noting that although our corpora are much smaller than CC-Matrix (by a factor of 2–8), we generally achieve better translation quality with them, indicating that the data is cleaner. These findings, however, allow only limited conclusions on the performance of the underlying methods (our hierarchical mining approach vs. the global mining approach of CC-Matrix) since they were executed on different, albeit quite similar, datasets (targeted crawling vs. pre-existing CommonCrawl).

A clean apples-to-apples comparison of the two approaches would be very difficult to carry given the scale of the data and the different data sources used. Nevertheless, we believe that the two large-scale efforts for these methods (CC-Matrix and

| Language | Ours | CC-M | OPUS | OPUS+CC-M | Ours+OPUS | Ours+OPUS +CC-M | NLLB Distilled |
|---|---|---|---|---|---|---|---|
| Hindi | 31.8 | 32.7 | 32.1 | 35.2 | 34.3 | 35.1 | 36.2 |
| Nepali | 26.0 | 7.8 | 13.0 | 21.5 | 25.2 | 25.8 | 33.6 |
| Burmese | 16.8 | 7.8 | 8.8 | 11.1 | 18.4 | 16.7 | 22.7 |
| Thai | 19.3 | – | 19.3 | – | 20.4 | – | 25.3 |
| Lao | 23.3 | 11.6 | 13.1 | 12.5 | 24.4 | 23.1 | 28.4 |
| Khmer | 13.7 | 8.9 | 10.6 | 17.0 | 19.8 | 21.3 | 26.3 |
| Vietnamese | 30.1 | 34.0 | 31.2 | 34.7 | 32.5 | 34.2 | 32.4 |
| Indonesian | 37.6 | 31.5 | 26.5 | 32.8 | 37.7 | 32.9 | 41.3 |
| Korean | 23.4 | 22.3 | 19.5 | 22.9 | 22.2 | 23.7 | 26.0 |

Table 6: Combining corpora: When combining our corpus with CC-Matrix and OPUS, we typically see improvements. The corpora are simply concatenated. The table reports spmBLEU scores on Flores-200 devtest for the models trained on the data.

ours) give strong evidence to the advantage of our approach.

## 6 Analysis

### 6.1 Combining Corpora

The three corpora we compare — OPUS, CC-Matrix, and ours — are obtained in quite different ways. Hence, we would expect that combining these corpora would lead to even better translation results.

Table 6 shows spmBLEU scores on Flores-200 devtest for the combinations OPUS+CC-Matrix, Ours+OPUS, and Ours+OPUS+CC-Matrix. For 3 languages (Khmer, Thai, and Korean) and almost Hindi–English, we do achieve the best results this way, while for Vietnamese the addition of our data slightly hurts (–0.5 BLEU) and for 3 languages (Burmese, Lao, Indonesian) the addition of the CC-Matrix corpus leads to worse results (–1.7, –1.3, and –4.8, respectively).

Note that we simply concatenated the corpora, and the CC-Matrix corpus has bigger impact on the results due to its typically larger size. There are many other ways to combine and weigh corpora which should be explored in future work by any researcher using this data.

### 6.2 Comparison with NLLB Distilled Data

Table 6 also contrasts the quality of the systems trained on the various combinations of corpora with systems built on data distilled with the NLLB model (these are the same numbers as in Table 2). Notably, the distilled data yields better quality systems for all languages except for Vietnamese. This observation is mirrored by Finkelstein et al.

| Language | Ours | | Statistical | |
|---|---|---|---|---|
| | BLEU | Words | BLEU | Words |
| Nepali | 26.0 | 32m | 23.8 | 31m |
| Burmese | 16.8 | 28m | 11.5 | 13m |
| Khmer | 13.7 | 23m | 9.4 | 9m |
| Vietnamese | 30.1 | 94m | 31.1 | 123m |
| Korean | 23.4 | 118m | 21.8 | 88m |

Table 7: Comparison of our neural methods with the statistical Paracrawl methods for document and sentence alignment.

(2024)'s finding that a distilled data set synthesized from a PaLM-2 Bison LLM model outperforms WMT training data.

However, it would be wrong to conclude that there is no need for crawled data and we should instead build our systems with synthetic data. Models such as NLLB rest on a vast collection of diverse data sources for training to achieve high quality, so crawled data is required to get started.

Nevertheless, this finding illustrate the complex data selection choices when it comes to building the best possible system for a given language pair and domain. We expect that future work will explore how to best combine and sequence the diverse set of data resources in more detail.

### 6.3 Comparison with Statistical Methods

Our pipeline makes two changes to the Paracrawl pipeline: use of a neural machine translation model for document alignment and sentence alignment based on neural sentence embeddings. Paracrawl uses a Moses-based statistical machine translation model and the lexicon-based Hunalign sentence aligner.

By running both the original pipeline and the pipeline with these changes, we can directly compare if the changes lead to an improved corpus. Results are shown in Table 7. We carried out this comparison only for 5 of the 9 languages due to the computation cost involved. Nevertheless, we covered both lower-resourced and higher-resourced languages. Except for Vietnamese (–1.0 BLEU), the neural methods lead to better results by a difference of +1.6 BLEU (Korean) to +5.3 BLEU (Burmese).

Since Vietnamese is an outlier here again (our new parallel corpus is also worse than CC-Matrix), we checked the execution of our pipeline for that language but could not find any obvious errors.

### 6.4 Computational Cost

We processed a total number of 127,064 web crawls. The size of the crawls has a very skewed distribution, with relatively few large crawls and a long tail of crawls that have only few web pages in the targeted languages. So, we can only make rough estimates about the processing cost.

Having said that, our document aligner takes about half an hour on average, of which half is spent on translation, summing up to about 2600 CPU days.

The sentence aligner takes about 6 minutes on average, the biggest computational cost being embedding of sentences with LASER, summing up to about 500 GPU days.

There is also signifcant time spent on extracting text from the web pages — we do not have reliable numbers on this. Note that this involves processing web crawls for which we ultimately do not find any content in the targeted languages and that are not included in our statistics here.

Sentence pair filtering takes tens of hours, training a neural model on a dataset takes a handful of days at most. Both these steps require a GPU.

### 7 Open Source Release

The corpora are available at the offical Paracrawl website `http://www.paracrawl.eu/`. Rachel Wicks created a document-aligned version of the corpus which is available at `https://huggingface.co/datasets/jhu-clsp/paradocs` using the approach outlined by Wicks et al. (2024).

### 8 Limitations

The motivating goal for this work was to create high-quality parallel corpora for important languages that have previously not received much attention. The languages were also chosen due to their large difference to English, often even using non-Latin writing systems.

Given the vast computational cost involved, we only have limited results on the comparison of methods. For instance, a more fine-grained demonstration of the effectiveness of the document aligner and sentence aligner in isolation would be useful. We do show that both in combination lead to better outcomes.

There are many more experiments that could be done with the data, such as more closely tracking how the quality of the machine translation model impacts the effectiveness of the document aligner. Another big area for follow-up research is how to best combine and filter different corpora for a language pair.

We are aware that much of the crawled data may stem from machine translation (Thompson et al., 2024). However, we argue that data quality is a better guide than the origin of the translations. Hence, we take a holistic filtering approach. See also work by Kreutzer et al. (2022) and Ranathunga et al. (2024) on the discussion of quality of web-crawled corpora.

Finally, the only measure of translation quality that we offer is the translation quality of a machine translation system trained on a dataset. While this is ultimately what is most important for the consumer of this data, it also ignores many other aspects of data quality, such as toxic content or bias. We added a toxicity filtering step but did not evaluate it, partly due to the vagaries of this task.

### 9 Risks

Our corpora may include harmful and violent content. It may also contain content that is copyrighted. We claim that our use of web-crawled data follows fair-use exceptions but we will remove data if any specific requests are made, thus slightly altering the composition of the data.

### 10 Conclusions

We deployed neural methods to the Paracrawl processing pipeline, demonstrated their superiority against the previous statistical methods and the

global mining approach, added a novel toxicity filtering method, and created high-quality parallel corpora for South and East Asian languages. We show that for 7 of the 9 languages our data leads to improvements in translation quality when building neural machine translation systems, for some languages dramatically.

We also spend significant effort on distilling NLLB models, reducing the computational cost by roughly doubling translation speeds, while using only a single CPU core vs. a full GPU — or a thousand-fold speed increase when calculated in terms of compute cores.

We release[4] all our corpora and models open source, with a liberal license for commercial and research use.

# References

Haluk Açarçiçek, Talha Çolakoğlu, Pınar Ece Aktan Hatipoğlu, Chong Hsuan Huang, and Wei Peng. 2020. Filtering noisy parallel corpus using transformers with proxy task learning. In *Proceedings of the Fifth Conference on Machine Translation*, pages 940–946, Online. Association for Computational Linguistics.

Mikko Aulamo, Nikolay Bogoychev, Shaoxiong Ji, Graeme Nail, Gema Ramírez-Sánchez, Jörg Tiedemann, Jelmer van der Linde, and Jaume Zaragoza. 2023. HPLT: High performance language technologies. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 517–518, Tampere, Finland. European Association for Machine Translation.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

Peter F. Brown, Jennifer C. Lai, and Robert L. Mercer. 1991. Aligning sentences in parallel corpora. In *Proceedings of the 29th Annual Meeting on Association for Computational Linguistics*, ACL '91, pages 169–176, Stroudsburg, PA, USA. Association for Computational Linguistics.

Christian Buck and Philipp Koehn. 2016. Quick and reliable document alignment via TF/IDF-weighted

cosine distance. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 672–678, Berlin, Germany. Association for Computational Linguistics.

Vishrav Chaudhary, Yuqing Tang, Francisco Guzmán, Holger Schwenk, and Philipp Koehn. 2019. Low-resource corpus filtering using multilingual sentence embeddings. In *Proceedings of the Fourth Conference on Machine Translation (WMT)*.

Aswarth Abhilash Dara and Yiu-Chang Lin. 2016. Yoda system for wmt16 shared task: Bilingual document alignment. In *Proceedings of the First Conference on Machine Translation*.

Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2020. CCAligned: A massive collection of cross-lingual web-document pairs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5960–5969, Online. Association for Computational Linguistics.

Miquel Esplà-Gomis, Víctor M. Sánchez-Cartagena, Jaume Zaragoza-Bernabeu, and Felipe Sánchez-Martínez. 2020. Bicleaner at WMT 2020: Universitat d'alacant-prompsit's submission to the parallel corpus filtering shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 952–958, Online. Association for Computational Linguistics.

Mara Finkelstein, David Vilar, and Markus Freitag. 2024. Introducing the NewsPaLM MBR and QE dataset: LLM-generated high-quality parallel data outperforms traditional web-crawled data. In *Proceedings of the Ninth Conference on Machine Translation (WMT)*.

William A Gale and Kenneth W Church. 1993. A program for aligning sentences in bilingual corpora. *Computational linguistics*, 19(1):75–102.

Mandy Guo, Yinfei Yang, Keith Stevens, Daniel Cer, Heming Ge, Yun-hsuan Sung, Brian Strope, and Ray Kurzweil. 2019. Hierarchical document encoder for parallel corpus mining. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 64–72, Florence, Italy. Association for Computational Linguistics.

Kevin Heffernan, Onur Çelebi, and Holger Schwenk. 2022. Bitext mining using distilled sentence representations for low-resource languages.

Marcin Junczys-Dowmunt. 2018. Dual conditional cross-entropy filtering of noisy parallel corpora. In *Proceedings of the Third Conference on Machine Translation*, Belgium, Brussels. Association for Computational Linguistics.

Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast

neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.

Huda Khayrallah and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83, Melbourne, Australia. Association for Computational Linguistics.

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of the Tenth Machine Translation Summit (MT Summit X)*, Phuket, Thailand.

Philipp Koehn, Vishrav Chaudhary, Ahmed El-Kishky, Naman Goyal, Peng-Jen Chen, and Francisco Guzmán. 2020. Findings of the WMT 2020 shared task on parallel corpus filtering and alignment. In *Proceedings of the Fifth Conference on Machine Translation*, pages 726–742, Online. Association for Computational Linguistics.

Philipp Koehn, Francisco Guzmán, Vishrav Chaudhary, and Juan Pino. 2019. Findings of the WMT 2019 shared task on parallel corpus filtering for low-resource conditions. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 54–72, Florence, Italy. Association for Computational Linguistics.

Philipp Koehn, Huda Khayrallah, Kenneth Heafield, and Mikel L. Forcada. 2018. Findings of the WMT 2018 shared task on parallel corpus filtering. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 726–739, Belgium, Brussels. Association for Computational Linguistics.

Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iroro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.

Murathan Kurfalı and Robert Östling. 2019. Noisy parallel corpus filtering through projected word embeddings. In *Proceedings of the Fourth Conference on Machine Translation (WMT)*.

Thanh Le, Hoa Trong Vu, Jonathan Oberländer, and Ondřej Bojar. 2016. Using term position similarity and language modeling for bilingual document alignment. In *Proceedings of the First Conference on Machine Translation*, pages 710–716, Berlin, Germany. Association for Computational Linguistics.

Robert C Moore. 2002. Fast and accurate sentence alignment of bilingual corpora. In *Conference of the Association for Machine Translation in the Americas*, pages 135–144. Springer.

Makoto Morishita, Katsuki Chousa, Jun Suzuki, and Masaaki Nagata. 2022. JParaCrawl v3.0: A large-scale English-Japanese parallel corpus. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6704–6710, Marseille, France. European Language Resources Association.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.

Wannaphong Phatthiyaphaibun, Korakot Chaovavanich, Charin Polpanumas, Arthit Suriyawongkul, Lalita Lowphansirikul, Pattarawat Chormai, Peerat Limkonchotiwat, Thanathip Suntorntip, and Can Udomcharoenchaikit. 2023. PyThaiNLP: Thai natural language processing in Python. In *Proceedings of the 3rd Workshop for Natural Language Processing Open Source Software (NLP-OSS 2023)*, pages 25–36, Singapore. Association for Computational Linguistics.

Surangika Ranathunga, Nisansa De Silva, Velayuthan Menan, Aloka Fernando, and Charitha Rathnayake. 2024. Quality does matter: A detailed look at the quality and utility of web-mined parallel corpora. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 860–880, St. Julian's, Malta. Association for Computational Linguistics.

Spencer Rarrick, Chris Quirk, and Will Lewis. 2011. MT detection in web-scraped parallel corpora. In *Proceedings of the 13th Machine Translation Summit (MT Summit XIII)*, pages 422–430. International Association for Machine Translation.

Philip Resnik. 1999. Mining the web for bilingual text. In *Proceedings of the 37th Annual Meeting of the Association of Computational Linguistics (ACL)*.

Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021. CCMatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500, Online. Association for Computational Linguistics.

Rico Sennrich and Martin Volk. 2010. MT-based sentence alignment for OCR-generated parallel texts. In *Proceedings of the Ninth Conference of the Association for Machine Translation in the Americas*.

Lei Shi, Cheng Niu, Ming Zhou, and Jianfeng Gao. 2006. A dom tree alignment model for mining parallel data from the web. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, pages 489–496. Association for Computational Linguistics.

Shashank Siripragada, Jerin Philip, Vinay P. Namboodiri, and C V Jawahar. 2020. A multilingual parallel corpora collection effort for Indian languages. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3743–3751, Marseille, France. European Language Resources Association.

Steve Sloto, Brian Thompson, Huda Khayrallah, Tobias Domhan, Thamme Gowda, and Philipp Koehn. 2023. Findings of the wmt 2023 shared task on parallel data curation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 95–102, Singapore. Association for Computational Linguistics.

Felipe Soares and Marta R. Costa-jussà. 2019. Unsupervised corpus filtering and mining. In *Proceedings of the Fourth Conference on Machine Translation (WMT)*.

Weiting Tan, Kevin Heffernan, Holger Schwenk, and Philipp Koehn. 2023. Multilingual representation distillation with contrastive learning. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1477–1490, Dubrovnik, Croatia. Association for Computational Linguistics.

Brian Thompson, Mehak Dhaliwal, Peter Frisch, Tobias Domhan, and Marcello Federico. 2024. A shocking amount of the web is machine translated: Insights from multi-way parallelism. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 1763–1775, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Brian Thompson and Philipp Koehn. 2019. Vecalign: Improved sentence alignment in linear time and space. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1342–1348, Hong Kong, China. Association for Computational Linguistics.

Jörg Tiedemann. 2009. News from OPUS - a collection of multilingual parallel corpora with tools and interfaces. In N. Nicolov, K. Bontcheva, G. Angelova, and R. Mitkov, editors, *Recent Advances in Natural Language Processing*, volume V, pages 237–248. John Benjamins.

Jakob Uszkoreit, Jay Ponte, Ashok Popat, and Moshe Dubiner. 2010. Large scale parallel document mining for machine translation. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 1101–1109, Beijing, China. Coling 2010 Organizing Committee.

Dániel Varga, Péter Halaácsy, András Kornai, Voktor Nagy, László Németh, and Viktor Trón. 2005. Parallel corpora for medium density languages. In *Proceedings of the RANLP 2005 Conference*, pages 590–596.

Rachel Wicks, Matt Post, and Philipp Koehn. 2024. Recovering document annotations for sentence-level bitext. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 9876–9890, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Runxin Xu, Zhuo Zhi, Jun Cao, Mingxuan Wang, and Lei Li. 2020. Volctrans parallel corpus filtering system for WMT 2020. In *Proceedings of the Fifth Conference on Machine Translation*, pages 985–990, Online. Association for Computational Linguistics.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Yuming Zhai, Lufei Liu, Xinyi Zhong, Gbariel Illouz, and Anne Vilnat. 2020. Building an English-Chinese parallel corpus annotated with sub-sentential translation techniques. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4024–4033, Marseille, France. European Language Resources Association.

Michał Ziemski, Marcin Junczys-Dowmunt, and Bruno Pouliquen. 2016. The United Nations parallel corpus v1.0. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3530–3534, Portorož, Slovenia. European Language Resources Association (ELRA).

## A  Marian Configuration

The following configuration is used both for the distilled translation models that are used by the document aligner as well as for evaluating different corpora. We guided alignment training, with alignments generated by fast-align.

**Model Configuration**
```
dec-cell: ssru
dec-cell-base-depth: 2
dec-cell-high-depth: 1
dec-depth: 2
dim-emb: 256
enc-cell: gru
enc-cell-depth: 1
enc-depth: 6
enc-type: bidirectional
tied-embeddings-all: true
transformer-decoder-autoreg: rnn
transformer-dim-ffn: 1536
transformer-ffn-activation: relu
transformer-ffn-depth: 2
transformer-guided-alignment-layer: last
transformer-heads: 8
transformer-no-projection: false
transformer-postprocess: dan
transformer-postprocess-emb: d
transformer-preprocess: ""
transformer-tied-layers:
    []
transformer-train-position-embeddings:
false
type: transformer
```

**Decoder Configuration**
```
models
- model.intgemm.alphas.bin
shortlist:
- lex.s2t.gz
- false
beam-size: 1
normalize: 1.0
word-penalty: 0
mini-batch: 64
maxi-batch: 1000
maxi-batch-sort: src
workspace: 2000
max-length-factor: 2.5
gemm-precision: int8shiftAlphaAll
```

**Training Parameters**
```
–dim-vocabs 32000 32000
```

```
–max-length 200
–exponential-smoothing
–cost-type ce-mean-words
–mini-batch-fit -w 3000
–mini-batch 300
–maxi-batch 500
–sync-sgd –optimizer-delay 2
–learn-rate 0.0003 –lr-report
–lr-warmup 16000
–lr-decay-inv-sqrt 32000
–optimizer-params 0.9 0.98 1e-09
–clip-norm 0
–valid-freq 5000 –save-freq 5000
–disp-freq 1000
–valid-metrics bleu-detok ce-mean-words
–valid-mini-batch 64 –beam-size 1
–normalize 1
–early-stopping 100
```

**Decoding Parameters**
```
–beam-size 1 –mini-batch 32
–maxi-batch 100 –maxi-batch-sort src -w
128
–skip-cost –cpu-threads 1
```

# Plug, Play, and Fuse: Zero-Shot Joint Decoding via Word-Level Re-ranking Across Diverse Vocabularies

**Sai Koneru**[1], **Matthias Huck**[2], **Miriam Exel**[2], and **Jan Niehues**[1]

[1] Karlsruhe Institute of Technology

[2] SAP SE, Dietmar-Hopp-Allee 16, 69190 Walldorf, Germany

{sai.koneru, jan.niehues}@kit.edu

{matthias.huck, miriam.exel}@sap.com

## Abstract

Recent advancements in NLP have resulted in models with specialized strengths, such as processing multimodal inputs or excelling in specific domains. However, real-world tasks, like multimodal translation, often require a combination of these strengths, such as handling both translation and image processing. While individual translation and vision models are powerful, they typically lack the ability to perform both tasks in a single system. Combining these models poses challenges, particularly due to differences in their vocabularies, which limit the effectiveness of traditional ensemble methods to post-generation techniques like N-best list re-ranking. In this work, we propose a novel zero-shot ensembling strategy that allows for the integration of different models during the decoding phase without the need for additional training. Our approach re-ranks beams during decoding by combining scores at the word level, using heuristics to predict when a word is completed. We demonstrate the effectiveness of this method in machine translation scenarios, showing that it enables the generation of translations that are both speech- and image-aware while also improving overall translation quality[1].

## 1 Introduction

A broad spectrum of Large Language Models (LLMs) are being developed at an increasing pace, with efforts focused alone or together on adapting them to specific domains (Roziere et al., 2023; Bolton et al., 2024; Colombo et al., 2024), enhancing their ability to process multiple modalities (Liu et al., 2023; Tang et al., 2023; Li et al., 2024; Beyer et al., 2024), or training general-purpose LLMs using high-quality data, advanced architectures, and

larger numbers of parameters (Touvron et al., 2023; Dubey et al., 2024; Jiang et al., 2023a; Mesnard et al., 2024). As a result, numerous models are now publicly available, each with its own unique strengths and weaknesses.

Many use cases, such as image-aware translation in movie subtitling, require combining these strengths because visual cues can be essential for disambiguating the text and ensuring accurate translations.. Currently, LLMs, such as Tower (Alves et al., 2024), Alma-R (Xu et al., 2024a), and Madlad-400 (Kudugunta et al., 2024), excel at translation tasks (Kocmi et al., 2024), while models like PaliGemma (Beyer et al., 2024) and LLava (Li et al., 2024) are leading in vision-related tasks. To effectively address image-aware translation, it is essential to harness the strengths of both translation and vision models.

One way to address such a task is to train a multimodal LLM to enhance its translation capabilities without compromising its vision abilities or vice versa. However, this approach requires additional training and task-specific data. Another approach is to leverage ensembling the two models via shallow fusion (Gulcehre et al., 2015) or re-ranking the N-best list (Hasan et al., 2007). The disadvantage of shallow fusion is that it assumes both models share the same vocabulary, which is often not the case with current open-source models.

Additionally, re-ranking the N-best list is insufficient because it doesn't allow models to influence each other during decoding. For example, in Figure 1, translating from English to gender-marked language French using audio and transcript shows this limitation. The Speech Translation (ST) model correctly uses the speaker's voice to translate "fell" into the right gender form but misidentifies the name "Ples." On the other hand, the Machine Trans-

---

[1] Code can be found at: https://ai4lt.anthropomatik.kit.edu/english/projects_kontextmt.php

Figure 1: The source sentence to be translated is ambiguous because the translation of the word **"fell"** can be either masculine (**"tombé"**) or feminine (**"tombée"**), depending on the speaker's gender. Seamless-Large V2 (Barrault et al., 2023) utilizes audio cues to correctly determine the gender form but struggles to accurately translate the name **"Mrs Ples"** using audio alone. In contrast, the text translation model Madlad-400-10b-mt (Kudugunta et al., 2024) relies on the gold transcript to correctly translate the name but fails to resolve the gender ambiguity. By combining both models using our approach, the translation correctly captures both the gender form and the named entity.

lation (MT) model correctly translates the name but can't use the speaker's voice for gender disambiguation. Thus, re-ranking falls short, as the correct forms may not even be in the N-best list due to low probability with missing cues.

Furthermore, re-ranking during the decoding process is impractical because the hypotheses are partial and may not align with the tokenization of the ranker model, leading to incorrect probability estimates (Section 2.1). Thus, resolving vocabulary mismatches by mapping the vocabulary of one model to another (Minixhofer et al., 2024; Xu et al., 2024b) is necessary to allow the merging of probabilities during decoding. However, this approach requires significant additional training steps and can lead to deviations from the original model. Therefore, developing a plug-and-play approach that seamlessly combines different models without requiring additional training or task-specific data is highly advantageous.

This work aims to enable the ranker model to influence the decoding process (online) without any constraints compared to conventional offline N-best list re-ranking. We address this by ensuring that the ranker model only influences the scores for completed words and not for the last word if it is unfinished. Additionally, we propose using the ranker model to determine whether the last word is finished rather than relying on look-ahead

approaches to maintain efficiency.

Our main contributions are summarized below:

1. **Online Re-Ranking Algorithm**: We introduce a novel re-ranking algorithm that operates at the word level during decoding at sub-word level, allowing for more accurate tokenization and better integration of information from different models

2. **Plug-and-Play Approach**: Our method does not require additional training or task-specific data, making it a flexible and practical solution for integrating multiple models with different strengths.

3. **Context-aware Translations**: We demonstrate through experiments including targeted multimodal test sets, which require information from both modalities, that our approach effectively combines the strengths of different models and improves translation quality (Illustrated in Figure 1).

## 2 Methodology

Given that many models are trained on different tasks, architectures, modalities, and data types, combining these models to leverage inputs from multiple modalities and facilitate knowledge sharing is highly beneficial. Moreover, it is ideal if the

1468

ensembling approaches satisfy the following constraints: 1) It should not rely on shared vocabularies for flexibility in choosing models and maximizing potential combinations. 2) Effective knowledge sharing should occur during decoding to better navigate the search space exploiting this knowledge at each step. 3) Avoid requiring additional training, parameters, or major dependence on task-specific data for maximum applicability and not cause deviations from the pre-trained model.

This section presents our algorithm for ensembling models with different vocabularies that satisfy the aforementioned constraints. First, we explain why re-ranking partial hypotheses can lead to incorrect probability estimates if the word is incomplete. Next, we introduce and justify a heuristic-based approach that predicts whether a hypothesis is at the end of a word, allowing for accurate re-ranking of completed words in partial hypotheses. Finally, we formally describe the complete algorithm, detailing how we merge probabilities from different models and how this process can be integrated with decoding strategies.

## 2.1 Challenges of Re-Ranking Partial Hypotheses

Current Neural Machine Translation (NMT) and LLM-based models can utilize various tokenization methods, such as byte-pair encoding (BPE) (Sennrich et al., 2016) or SentencePiece (Kudo and Richardson, 2018). These methods often result in distinct vocabularies due to variations in the data and tokenizer training processes. Despite these differences, techniques like re-ranking can still enable estimating the probability of sentences generated from another model. This is achieved by detokenizing the hypothesis from generator model and re-tokenizing it using the ranker model's vocabulary. This process enables the ranker model to produce accurate probability estimates based on its own tokenization scheme.

Now, consider the case of re-ranking while the hypotheses are still being decoded. Assume we have model $\mathcal{M}_G$ (the generator) and model $\mathcal{M}_R$ (the ranker), each using different tokenizers assign all the tokens in the sentence "Decoding is awesome" with a probability of p for a particular input. However, $\mathcal{M}_G$ tokenizes the sentence with subword tokens as "Dec od ing _is _awe some," while

$\mathcal{M}_R$ would tokenize it as "Dec od ing _is _awes ome."

If we attempt to re-rank during the decoding process, $\mathcal{M}_R$ will provide correct probability estimates up until "_is" is generated. However, when the generator predicts "_awe," $\mathcal{M}_R$ would incorrectly estimate the probability because it expects "_awes" instead. Even though both models aim to generate the same sentence, this tokenization mismatch leads to incorrect probability estimates during the decoding process, making online re-ranking challenging.

## 2.2 End-of-Word Prediction in Decoding for Accurate Re-Ranking

While the partially generated hypothesis cannot be accurately ranked at every time step, consider the cases when each word is finished. At that time, we can re-rank the complete hypothesis as the last word is fully generated and the ranker model can tokenize the completed word as it would have done naturally, thereby providing accurate probability estimates. If we know that the last word is incomplete, we can use this information to wait and only rank the previously completed words. Knowing the end of the word enables more precise re-ranking during decoding, even with models that use different tokenization schemes.

Nonetheless, a significant challenge remains: how do we determine when the last word is completed? If the tokenizer places spaces at the right of characters, we could check the predicted token to see if it includes a space, signaling the end of a word. However, this approach is not universal, as many tokenizers do not follow this pattern, and we aim to develop a tokenizer-agnostic solution.

One alternative is to perform a look-ahead step to check if the word has been completed, but this method is also sub-optimal, as it would require decoding twice for each step in the generation process, significantly increasing computational complexity and reducing efficiency. We need a more efficient and generalizable method to determine when a word has been completed during decoding.

To address these challenges, we propose using the ranker model to predict the next token and determine if the word has been completed. This approach offers two key advantages.

Firstly, if the ranker model predicts a space as

---

**Algorithm 1** Computing merged score of candidate with generator and ranker models.

1: **procedure** MERGESCORE
2:     **Input:** Generator tokens $g_1, g_2, g_3, \ldots, g_n$, Reranker tokens $r_1, r_2, r_3, \ldots, r_m$, Generator Model $\mathcal{M}_G$, Ranker model $\mathcal{M}_R$, Generator Input $\mathcal{I}_G$, Ranker Input $\mathcal{I}_R$, Re-ranking weight $\alpha$,
3:     **Output:** $merged\_score$
4:     $next\_tok \leftarrow \arg\max \log \mathcal{P}(y|r_1, \ldots, r_m; \mathcal{I}_R; \mathcal{M}_R)$
5:     **if** $next\_tok[0] ==$ "_" or $next\_tok ==$ "<eos>" **then**
6:         $full_G \leftarrow \frac{1}{n} \sum \log \mathcal{P}(g_1, g_2, \ldots, g_n | \mathcal{I}_G; \mathcal{M}_G)$         ▷ Generator Score for all words
7:         $full_R \leftarrow \frac{1}{m} \sum \log \mathcal{P}(r_1, r_2, \ldots, r_m | \mathcal{I}_R; \mathcal{M}_R)$         ▷ Ranker Score for all words
8:         $merged\_score \leftarrow (\alpha) \times full_G + (1 - \alpha) \times full_R$
9:     **else**
10:        $[g_1, \ldots, g_j], [g_{j+1}, \ldots, g_n] \leftarrow$ split_candidate$(g_1, \ldots, g_n)$     ▷ Last word from j+1 token
11:        $[r_1, \ldots, r_k], [r_{k+1}, \ldots, r_m] \leftarrow$ split_candidate$(r_1, \ldots, r_m)$     ▷ Last word from k+1 token
12:        $prev_G \leftarrow \frac{1}{j} \sum \log \mathcal{P}(g_1, g_2, \ldots, g_j | \mathcal{I}_G; \mathcal{M}_G)$     ▷ Generator Score for previous words
13:        $prev_R \leftarrow \frac{1}{k} \sum \log \mathcal{P}(r_1, r_2, \ldots, r_k | \mathcal{I}_R; \mathcal{M}_R)$     ▷ Ranker Score for previous words
14:        $prev_{GR} \leftarrow (\alpha) \times prev_G + (1 - \alpha) \times prev_R$
15:        $last_G \leftarrow \sum \log \mathcal{P}(g_{j+1}, \ldots, g_n | \mathcal{I}_G; \mathcal{M}_G)$
16:        $merged\_score \leftarrow \frac{1}{n}[prev_{GR} \times j + last_G]$     ▷ Re-normalized merged score
17:     **end if**
18: **end procedure**

---

the next top character, it indicates that the current last word has been completed. The hypothesis will be tokenized correctly, given that it is the prediction from the ranker model itself. Secondly, this prediction can be done together with the re-ranking process by simply also predicting the next token given the previous tokens of the current hypothesis to the ranker model.

This method is more efficient than the look-ahead approach, requiring only one pass of the generator and the ranker model. In contrast, the look-ahead method would require two passes of the generator and one pass of the ranker model. Using the ranker model in this way, we can ensure proper tokenization and accurate probability estimates during the decoding process (online) without additional computational overhead.

## 2.3 Integrating Online Re-Ranking with Search

This section formalizes achieving online re-ranking at a word level using beam search as an example of a decoding strategy. Note that the approach can also be applied to other strategies, with slight modifications when necessary.

A set of candidate sequences is typically maintained during the search, with the number of candidates equal to the configured beam size $b$. At each time step, for each of the $b$ candidate sequences, the model computes likelihood scores for all possible token extensions based on the vocabulary size $V$. This results in a total of $b \times V$ possible extensions. From these $b \times V$ extensions, the top $b$ sequences with the highest scores are selected to form the new set of candidate sequences. This process is repeated iteratively, updating the candidate sequences at each step until enough beams are generated that include end-of-sentence tokens or until a predefined length limit is reached.

To enable re-ranking during the decoding process, we need to adjust the scores of the possible extensions using the ranker model. Directly calculating the likelihood of all extensions would be computationally impractical. Therefore, we introduce a new parameter $topk$, which selects the top $topk$ extensions for each beam during re-ranking.

Hence, at each time step, the generator model calculates the likelihood scores for all $V$ possible extensions for each of the $b$ candidate sequences, resulting in $b \times V$ extensions. Instead of re-ranking all $b \times V$ extensions, the top $topk$ extensions with the highest likelihood scores are selected for each beam. Thus, only $b \times topk$ extensions are considered during re-ranking. For the selected $b \times topk$ extensions, the ranker model estimates their scores

and combines them with the original generator scores. For the remaining $b \times (V - topk)$ extensions, the scores are set to $-\infty$ (logically equivalent to discarding them) since they would not be selected in the top beams.

This method significantly reduces computational complexity while allowing effective re-ranking of the most promising candidate extensions, improving the decoding process.

At every decoding step, the problem can be reformulated as determining the merged score of the top candidates according to both models.

When calculating the merged score during decoding, it's essential to exclude the ranker model's probability if the last word in the current beam is incomplete. This prevents incomplete words from skewing the final score. For beams with incomplete final words, we combine the joint scores of the preceding words with the generator's score for the last word, ensuring proper normalization to address scale differences between finished and unfinished beams.

After computing the merged scores, we select the top extensions and repeat the process until all beams reach the end-of-sentence token. This method ensures that the final translation is based on fully formed words, optimizing the ranker model's effectiveness and maintaining consistent scoring across all candidates.

### 2.3.1 Unified Scoring with Generator and Ranker

The algorithm to compute the merged score is formally defined in Algorithm 1 and explained below.

Let us consider two models: the Generator $\mathcal{M}_G$ and the Ranker $\mathcal{M}_R$. Let $\mathcal{C}$ denote the current candidate for re-ranking and inputs $\mathcal{I}_G$ and $\mathcal{I}_R$ for $\mathcal{M}_G$ and $\mathcal{M}_R$ respectively.

Let the full candidate $\mathcal{C}$ consist of tokens $g_1, g_2, g_3, \ldots, g_n$ and $r_1, r_2, r_3, \ldots, r_m$ according to $\mathcal{M}_G$ and $\mathcal{M}_R$, respectively. Note that $n$ and $m$ denote the length of the sequence, and they may differ due to different tokenization.

The key idea is to rank and merge scores for completed words. We use the ranker model to predict the next token and determine if the last word is finished (**Line 4**).

**If the last word is finished**: We can calculate the probability of the full sequence in this case, sim-

ilar to the case of N-best list re-ranking. First, we calculate the likelihood of the candidate by averaging the log probabilities for both the generator and the ranker (**Line 6-7**). Then, we merge the scores from both models to determine the final score for the candidate sequence using a hyper-parameter $\alpha$ for weighting (**Line 8**). This combined score considers the estimates from both models, allowing for contributions from both models.

**If the last word is incomplete**: We cannot rank the last word due to potential incorrect tokenization. However, we can still estimate the tokens preceding the last word using the ranker model and merge their probabilities. First, we split the candidate into previous and last words based on the ranker and generator (**Lines 10-11**). We compute the merged score for the previous words using the weighting parameter $\alpha$ (**Lines 12-14**). For the last word, we rely solely on the generator's scores. To address length normalization issues when combining scores from both models, we re-normalize the merged score for the previous words by multiplying it by the length of the previous word tokens $j$ from the generator, adding the last word's score, and normalizing by the total length $n$ (**Lines 15-16**).

This integration process ensures that the re-rankers are utilized at the appropriate decoding stages, thereby enhancing the overall quality of the generated sequences by combining the strengths of both models.

## 3 Test Suites

The major advantage of combining models with different vocabularies zero-shot is that it leverages the strengths of available pre-trained models to generate more accurate and robust output. This is particularly relevant in multimodal scenarios, where unimodal systems excel in their respective modalities but are weaker or incapable of processing other modalities. Furthermore, it can also enhance quality compared to N-best list re-ranking when used as an ensembling technique as it waits until the complete sequence is generated. Hence, to validate our approach, we consider three MT scenarios as a test bed where quality can be improved by combining different sources and evaluating with targeted test sets that require information from both models. An overview of test suites is provided in Table 1.

## 3.1 Unimodal MT

We evaluate the use case of ensembling different LLM models to enhance translation quality. This is particularly relevant given the rapid development of various translation LLMs, where combining different systems can improve quality and robustness. We use the WMT 2022 *English → German* test set (Kocmi et al., 2022) to validate our approach and focus solely on assessing translation quality.

| Test Set | Language Pair | # Examples | Phenomena |
|---|---|---|---|
| MuST-SHE | *En → Fr* | 315 (1108) | Gender Disambiguation Translation |
| CoMMuTE | *En → De* | 300 | Word Disambiguation Translation |
| WMT22 | *En → De* | 2037 | Translation |

Table 1: Overview of test suites. For MuST-SHE, 315 examples are utterances where information is available in audio. However, we use the full test set with other types of bias when reporting translation quality.

## 3.2 Multimodal MT

Translating from English to gender-marked languages is challenging when the source text lacks clear gender cues. To evaluate bias in current NMT systems, Bentivogli et al. (2020) developed the *MuST-SHE* test suite, which includes examples with varying forms of gender bias. This suite features cases where gender information is conveyed through audio cues, such as the speaker's voice.

While End-to-End ST systems can handle such cases, they often fall short compared to advanced translation LLMs (Agarwal et al., 2023). Therefore, we use MuST-SHE for *English → French* to investigate if combining ST and translation LLMs can improve translation quality and address gender ambiguity.

Similarly, images can assist in disambiguating text and enhancing translation quality. However, translation LLMs typically do not process images, and vision LLMs alone are inadequate for translation tasks. We combine these models to leverage their strengths for better image-aware translations.

Existing vision translation test sets often lack ambiguity, making image inputs unnecessary (Vijayan et al., 2024). To address this, Futeral et al. (2023) introduced CoMMuTE, which features ambiguous source sentences with two images and their translations. We use CoMMuTE for *English*

*→ German* translation in a generative framework to evaluate if images can enhance translations without compromising overall quality.

## 4 Results

This section presents the experiments conducted using our ensembling approach across various test suites. Since each test suite has a distinct experimental setup, we will address them individually. First, we will specify the models and evaluation metrics applied in each scenario. Then, we will present the results and highlight our main findings.

### 4.1 Ensembling for Improving Translations

**Models:** We aim to combine two models that excel in translation but possess different strengths. For this purpose, we chose the *Madlad-10B*[2], an encoder-decoder architecture trained on extensive parallel data, and *ALMA-13B-R*[3], a decoder model trained using contrastive preference optimization and selecting high quality data (Xu et al., 2024b).

**Metrics:** As the models that we would like to ensemble are high quality, we report with several neural metrics to reliably validate the improvements. For reference-based we report with COMET (Rei et al., 2022a) and BLUERT (Sellam et al., 2020; Pu et al., 2021) whereas for reference-free we report with COMET-KIWI (Rei et al., 2022b), COMET-KIWI-XXL (Rei et al., 2023) and XCOMET-XXL (Guerreiro et al., 2023) metrics.

**Hyper-parameters:** We set the re-ranking weight $\alpha$ to $0.5$ given that both models have high quality and should be weighted equally. Furthermore, we set the $topk$ to 5 and the number of beams for the generator as 5.

To validate our combined model and online re-ranking approach, we compare it against several baselines. First, we check if the ensemble outperforms each individual model. Next, we evaluate if our method surpasses offline re-ranking techniques, indicating a more effective ranker influence and improved search space exploration during decoding.

We evaluate our approach using N-best list reranking, with *Madlad* as the generator and *Alma* as the ranker. We generate an N-best list of 25 hypotheses with $\alpha$ set to 0.5 to facilitate a fair

---

[2]https://huggingface.co/google/madlad400-10b-mt

[3]https://huggingface.co/haoranxu/ALMA-13B-R

| Generator | Ranker | Online | COMET 22 | COMET KIWI 22 QE | COMET KIWI XXL QE | XCOMET-XXL | BLEURT |
|---|---|---|---|---|---|---|---|
| | | | *No re-ranking* | | | | |
| GPT-4 | × | N/A | 87.29 | 83.48 | 84.91 | 97.56 | _ |
| Madlad-10B | × | N/A | 86.60 | 83.14 | 82.65 | 96.77 | 76.79 |
| Alma-13B-R | × | N/A | 86.40 | 83.28 | 84.25 | 97.48 | 77.20 |
| | | | *Offline re-ranking* | | | | |
| Madlad-10B | Alma-13B-R | × | 87.27 | 83.68 | 84.11 | 97.12 | 77.66 |
| Madlad-10B, Alma-13B-R | Madlad-10B, Alma-13B-R | × | 87.54 | **83.95** | 84.97 | 97.39 | 78.20 |
| | | | *Online re-ranking (ours)* | | | | |
| Madlad-10B | Alma-13B-R | ✓ | **87.69** | 83.94 | **85.20** | **97.68** | **78.36** |

Table 2: Performance of models on the WMT 22 *English → German* test set. Scores are highlighted in **bold** if it is the best in all configurations. Results for GPT-4 and Alma-13B-R are reported from (Xu et al., 2024b)

.

| Model | COMET 22 | COMET KIWI 22 QE | COMET KIWI XXL QE | XCOMET-XXL | BLEURT |
|---|---|---|---|---|---|
| GPT-4 | 87.29 | 83.48 | 84.91 | 97.56 | _ |
| Madlad | 86.60 | 83.14 | 82.65 | 96.77 | 76.79 |
| (Madlad) 5-best + QE | 87.33 | 83.83 | 86.45 | 97.25 | 77.78 |
| (Madlad + Alma Online re-rank) 5-best + QE | **87.66** | **84.12** | **87.86** | **97.91** | **78.31** |

Table 3: Performance of models on the WMT 22 *English → German* test set with Quality Estimation based re-ranking via selecting from 5-best list using comet-kiwi-xxl. Scores are highlighted in **bold** if it is the best in all configurations.

.

comparison between offline and online re-ranking methods. Additionally, we test a scenario where the N-best lists from both models are concatenated and jointly re-ranked on 50 hypotheses. We reports the results for the baselines and our approach in Table 2.

**Ensembling enables to reach state-of-the-art quality:** Both *Madlad* and *Alma* produce high-quality translations, though they still lag behind GPT-4 across all metrics. However, after applying offline re-ranking, their performance improves consistently, becoming competitive with GPT-4. When using our online re-ranking approach, the ensemble outperforms GPT-4 across all metrics and shows our proposed approach can improve the translation quality by a substantial margin.

**Online re-ranking outperforms offline joint re-ranking:** When *Madlad* serves as the generator and *Alma* as the ranker in our approach, the results are superior to those achieved with joint re-ranking, where both models are used simultaneously. Our approach enhances knowledge sharing and collaboration during the decoding process, leading to better translation quality.

### 4.1.1 Quality of N-best list

The primary motivation behind our approach was to influence the decoding process in real-time, rather than waiting until the end. If this is effective, we expect the N-best list to improve with online re-ranking. Additionally, using quality estimation should enhance the selection of the best hypothesis from the N-best list. To validate this, we utilize COMET-KIWI-XXL for selecting the best candidate from the top 5 beams of *Madlad*, comparing scenarios with and without online re-ranking and report the scores in Table 3.

We observe that integrating quality estimation significantly enhances Madlad's performance across all metrics. Using COMET-KIWI-XXL to select the best candidate from the top 5 beams improves score from $82.65 \rightarrow 86.45$. This improvement is also evident in the BLUERT score, increasing from $76.79 \rightarrow 77.78$. Additionally, comparing the top 5 beams with our approach, we find that the quality is superior, demonstrating that the early influence of *ALMA* in decoding. Furthermore, this allows to integration of multiple NMT models to generate the N-best list together and later combined with quality estimation for maximum performance.

### 4.2 Speech-Aware Translations

**Models:** To tackle gender ambiguity in text translation using speaker voice information, we combine a robust text translation model with a speech-based model that excels at disambiguating gender, even if it is not as strong in translation. We use the *Madlad* model (Kudugunta et al., 2024) for high-quality text translation with gold transcript and the *Seamless*[4] model for speech translation. Our approach employs *Madlad* as the *generator* and *Seamless* as the *ranker*, allowing us to leverage the speech model's ability to correct gendered forms in the translation.

However, we observed that the *Seamless* model exhibited a bias toward the masculine gender and struggled to effectively resolve gender ambiguities using speech. To mitigate this, we conducted additional fine-tuning using LoRA (Hu et al., 2021) on a balanced speaker dataset derived from MuST-C (TED talks) with gender annotations (Di Gangi et al., 2019; Gaido et al., 2020) (Training details in Appendix A.1). We remove talks that are present in MuST-SHE for no overlap. This "debiasing" process improved the model's ability to disambiguate gender based on speech. Consequently, we use the *Madlad* and adapted *Seamless* models to generate high-quality, speech-aware translations.

**Metrics:** To evaluate the effectiveness of our approach in disambiguating gender and improving translation quality, we use several key metrics. For gender disambiguation, we follow the methodology of Bentivogli et al. (2020) and report two metrics: accuracy (correct gender form is present) and coverage (either gender form is present).

For overall translation quality, we report BLEU (Papineni et al., 2002), ChrF2 (Popović, 2016) calculated using SacreBLEU (Post, 2018), and COMET (Rei et al., 2022a) (*wmt22-comet-da*) for brevity.

Additionally, we report *Sensitivity*, which measures the difference between the scores of correctly and incorrectly gendered references, as suggested by Bentivogli et al. (2020).

**Hyper-parameters:** For decoding with *Madlad*, we use beam search with 5 beams. Our proposed algorithm involves two key parameters: $\alpha$ and $topk$.

We set $topk$ to 5, resulting in a total of 25 candidates being ranked by *Seamless* at each step.

We optimized $\alpha$ through grid search on the MuST-C development set (Appendix A.3) via offline re-ranking and setting it to 0.8 based on these results. We also create an N-best list of 25 hypotheses with $\alpha$ at 0.8 for offline comparison and perform joint re-ranking on the combined 50 N-best lists. Results are summarized in in Table 4.

***Madlad* and *Seamless* complement each other:** *Madlad* excels in overall translation quality (83.5) compared to *Seamless* (79.31). While *Seamless* initially favors masculine terms, fine-tuning on balanced data improves overall quality to 80.48, significantly reducing masculine bias (90.44 to 65.89) and increasing feminine representation (25.92 to 50.18). Thus, the adapted *Seamless* demonstrates improved gender disambiguation, though *Madlad* remains superior in overall translation. Hence, combining the models can be highly beneficial.

**Online re-ranking improves overall translation quality:** After re-ranking with N-best list, we see that the translation quality is improved when *Madlad* as a generator and *Seamless Bal* as a ranker model (83.50 → 83.66). In the opposite scenario where *Seamless Bal* uses *Madlad* as a ranker model, the quality also improves (80.48 → 81.31) but is lower than *Madlad* alone. However, during online re-ranking, we see that we achieved the best performance of 83.78. This suggests that our approach facilitates knowledge sharing between the models during decoding, leading to significant quality enhancements.

**Balance between translation quality and gender disambiguation through online re-ranking** We observe that the highest accuracies for feminine terms (1F) are achieved when *Seamless Bal* is employed as a generator. Nevertheless, the overall translation quality in these instances is considerably lower compared to scenarios where *Madlad* is the generator. By using *Madlad* as a generator, we attain a higher average 1F score of 60.32 compared to offline re-ranking without compromising overall translation quality and better distribution across gender. Moreover, we achieved the highest sensitivity score of 1.1 across all configurations. This shows that our approach can consistently perform better than traditional N-best list re-ranking.

While the scores for the disambiguation are not

| Generator | Ranker | Online | 1F (Acc %) | 1F (Term Cov %) | 1M (Acc %) | 1M (Term Cov %) | Avg (Acc %) | COMET Correct | △ |
|---|---|---|---|---|---|---|---|---|---|
| | | | | *No re-ranking* | | | | | |
| Madlad | × | N/A | 25.92 | 68.39* | 90.44* | 63.65* | 58.18 | 83.52 | 0.90 |
| Seamless | × | N/A | 20.28 | 63.20 | 88.30 | 62.43 | 54.29 | 79.31 | 0.73 |
| Seamless Bal | × | N/A | 50.18 * | 62.73 | 65.89 | 59.02 | 58.03 | 80.48 | 0.83 |
| | | | | *Offline re-ranking* | | | | | |
| Madlad | Seamless Bal | × | 28.81 | 67.92 | **89.59** | 63.41 | 59.20 | 83.66 | 0.96 |
| Seamless Bal | Madlad | × | **40.90** | 65.09 | 77.99 | 60.97 | 59.44 | 81.31 | 0.90 |
| Madlad, Seamless Bal | Madlad, Seamless Bal | × | 29.83 | 67.92 | **89.59** | 63.41 | 59.71 | 83.64 | 0.96 |
| | | | | *Online re-ranking (ours)* | | | | | |
| Madlad | Seamless Bal | ✓ | 33.78 | **68.16** | 86.86 | **63.65\*** | **60.32\*** | **83.78\*** | **1.1\*** |

Table 4: Performance of models on the MuST-SHE test set for speech-aware translations. *Seamless Bal* indicates the adapted model trained on balanced gender data. △ denotes the *sensitivity*, i.e., the difference in scores between correct and incorrect references. Scores are highlighted in **bold** if online re-ranking improves over offline re-ranking and ∗ if it is the best in all configurations.

| Generator | Ranker | Online | BLEU Correct | △ | Chrf2 Correct | △ | COMET Correct | △ |
|---|---|---|---|---|---|---|---|---|
| Madlad-10B | × | N/A | 45.9 | 0.4 | 62.3 | 1.3 | 82.90 | 0.06 |
| PaliGemma-3B MT | × | N/A | 27.6 | **5.7** | 51.0 | **7.3** | 79.58 | **8.25** |
| Madlad-10B | PaliGemma-3B MT | × | 46.1 | 1.9 | **62.6** | 1.7 | **83.45** | 1.17 |
| Madlad-10B | PaliGemma-3B MT | ✓ | **46.2** | 1.8 | **62.6** | 1.9 | 83.25 | **1.34** |

Table 5: Performance of models on the CoMMuTE *English → German* test set for image-aware translations. △ indicates the sensitivity i.e difference between correct and incorrect references. Scores are highlighted in **bold** if it is the best in all configurations.

.

high, we would like to highlight that we focused on combining the strengths of the models. However, one can use targeted systems such as Gaido et al. (2020) to further improve the performance for the desired tasks.

## 4.3 Image-Aware Translations

**Models:** To integrate image information for disambiguating source text, a robust multimodal machine translation (MT) system is essential. Initially, we experimented with the off-the-shelf instruction-tuned Llava model[5] (Li et al., 2024). While Llava provided reasonable results, its performance was sub-par for our needs. Consequently, we chose to fine-tune the *PaliGemma* model[6] (Beyer et al., 2024), which was originally trained to generate captions in multiple languages. We fine-tuned *PaliGemma* using the Multi30k image captions

dataset (Elliott et al., 2016), adapting it with Q-LoRA (Appendix A.2) for enhanced image-aware translations (*PaliGemma-3B MT*).

**Metrics:** For evaluating this task, we use BLEU, ChrF2, and COMET scores, as we do not have specific annotations for words in the target sentences. To assess the impact of contextual information provided by the images, we also report the sensitivity metric △, to estimate how much the image context influences the translation quality.

**Hyper-parameters:** Vision LLMs require more memory because the image is encoded into a long sequence of tokens. Consequently, we were limited to using a beam size of 3 with a top-k of 3. Additionally, tuning the parameter $\alpha$ was challenging due to the lack of a dedicated ambiguous test set; using a standard test set would result in no weight being given to the vision model. Therefore, we report the oracle $\alpha$ of 0.9, which represents the best-performing weight on the test set, determined through a grid search with offline re-ranking. We report the scores in Table 5.

**PaliGemma is highly sensitive to image context:** We observe that the sensitivity $\triangle$ of our fine-tuned *PaliGemma* model for MT is notably high across all metrics (e.g., 5.7 BLEU), demonstrating that the model is effectively using the image information to influence its translations. This suggests that *PaliGemma* does not disregard the visual context during translation. However, despite this sensitivity, *PaliGemma's* overall translation quality significantly lags behind that of *Madlad*, as indicated by the lower COMET score (difference of 3.32). This disparity highlights the potential benefit of combining the strengths of both models to achieve more accurate and image-aware translations.

**No clear winner between offline and online re-ranking:** Comparing offline and online re-ranking, we find that re-ranking with *PaliGemma* enhances translations, evidenced by a sensitivity $\triangle$ increase of up to 1.28 COMET. There's also a slight improvement in overall translation quality after re-ranking. However, the difference between the two approaches is modest, especially given the small test set size of 300 examples.

We hypothesize two main factors behind the results. First, *Madlad* assigns very low probabilities to translations of ambiguous words it isn't biased toward, while *PaliGemma* avoids extremely high probabilities. As a result, merging probabilities tends to favor the incorrect translation with the highest overall score. Second, the test sentences are short, averaging 4-5 words, so the N-best list includes diverse variations, making offline re-ranking similar to the online approach. However, we believe our online re-ranking method could benefit longer sentences and stronger vision translation models.

## 5   Related Work

**Fusion for MT:** Integrating additional language models into MT systems via shallow or deep fusion, or through re-ranking, to improve translation quality is a well-studied area (Chen et al., 2006; Hasan et al., 2007; Gulcehre et al., 2015; Li and Jurafsky, 2016; Gulcehre et al., 2017; Herold et al., 2023). Stahlberg et al. (2018) explored advanced fusion method where an NMT model is trained from scratch while keeping a pre-trained language model fixed, allowing the model to learn only what

is missing. There has also been growing interest in combining NMT with document-level language models (Stahlberg et al., 2019; Petrick et al., 2023; Hoang et al., 2024). Unlike previous works that utilize static weights for merging probabilities, Jean and Cho (2020) propose dynamic coefficients, which are crucial for effectively combining models with different strengths.

**Ensembling:** System combination, which involves merging multiple hypotheses to generate a better version, is one approach to leveraging the strengths of different models (Bangalore et al., 2001; Matusov et al., 2006; Heafield and Lavie, 2010; Freitag et al., 2014). Another approach is to merge model parameters (Junczys-Dowmunt et al., 2016) or distill knowledge from the models (Freitag et al., 2017). With the increasing diversity of LLMs, recent research has explored methods to combine them through vocabulary merging (Xu et al., 2024b), generating new outputs based on hypotheses (Jiang et al., 2023b), or dynamically selecting different models at each step (Shen et al., 2024).

Our work differs from these approaches as it neither relies on vocabulary matching nor requires additional training data.

## 6   Conclusion

We proposed a novel ensembling strategy that operates at the word level during the decoding process to enhance knowledge sharing. Our approach demonstrated significant benefits across multiple scenarios. It proved effective for ensembling translation systems, and even when combined with quality estimation models, it achieved state-of-the-art translation quality. Additionally, experiments on targeted multimodal test sets revealed that our method facilitates better knowledge sharing compared to traditional re-ranking techniques.

For future work, we propose to explore unsupervised dynamic selection, enabling models to generate outputs only when they are better equipped for the task. We believe this approach could address the current limitations and lead to more significant improvements in image-aware translation.

## 7   Limitations

The major limitation of this work is that we operate at word-level which is not compatible for several

languages that are character based. Hence, it is not trivial to merge models for generating such languages. Further analysis is necessary on character-level tokenization to accurately re-rank during the decoding steps.

Another drawback is that, although re-ranking enhances translation quality, it incurs a latency cost. Unlike offline re-ranking, our approach employs the ranker model at each time step, resulting in significantly slower performance.

Finally, we focused mainly on ensembling the two models using static weights. However, since the models have different strengths, it is crucial to determine when to rely on one model or ensemble both. This dynamic approach would better exploit each model's strengths while avoiding the integration of their weaknesses.

## References

Milind Agarwal, Sweta Agrawal, Antonios Anastasopoulos, Luisa Bentivogli, Ondřej Bojar, Claudia Borg, Marine Carpuat, Roldano Cattoni, Mauro Cettolo, Mingda Chen, et al. 2023. Findings of the iwslt 2023 evaluation campaign. In *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*, pages 1–61.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Bangalore Bangalore, German Bordel, and Giuseppe Riccardi. 2001. Computing consensus translation from multiple machine translation systems. In *IEEE Workshop on Automatic Speech Recognition and Understanding, 2001. ASRU'01.*, pages 351–354. IEEE.

Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, et al. 2023. Seamless: Multilingual expressive and streaming speech translation. *arXiv preprint arXiv:2312.05187*.

Luisa Bentivogli, Beatrice Savoldi, Matteo Negri, Mattia A. Di Gangi, Roldano Cattoni, and Marco Turchi. 2020. Gender in danger? evaluating speech translation technology on the MuST-SHE corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6923–6933, Online. Association for Computational Linguistics.

Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. 2024. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*.

Elliot Bolton, Abhinav Venigalla, Michihiro Yasunaga, David Hall, Betty Xiong, Tony Lee, Roxana Daneshjou, Jonathan Frankle, Percy Liang, Michael Carbin, et al. 2024. Biomedlm: A 2.7 b parameter language model trained on biomedical text. *arXiv preprint arXiv:2403.18421*.

Boxing Chen, Roldano Cattoni, Nicola Bertoldi, Mauro Cettolo, and Marcello Federico. 2006. The itc-irst smt system for iwslt 2006. In *Proceedings of the Third International Workshop on Spoken Language Translation: Evaluation Campaign*.

Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, et al. 2024. Saullm-7b: A pioneering large language model for law. *arXiv preprint arXiv:2403.03883*.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.

Mattia A Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019. Must-c: a multilingual speech translation corpus. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2012–2017. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30k: Multilingual english-german image descriptions. In *Proceedings of the 5th Workshop on Vision and Language*, pages 70–74. Association for Computational Linguistics.

Markus Freitag, Yaser Al-Onaizan, and Baskaran Sankaran. 2017. Ensemble distillation for neural machine translation. *arXiv preprint arXiv:1702.01802*.

Markus Freitag, Matthias Huck, and Hermann Ney. 2014. Jane: Open source machine translation system combination. In *Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 29–32.

Matthieu Futeral, Cordelia Schmid, Ivan Laptev, Benoît Sagot, and Rachel Bawden. 2023. Tackling ambiguity with images: Improved multimodal machine translation and contrastive evaluation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5394–5413, Toronto, Canada. Association for Computational Linguistics.

Marco Gaido, Beatrice Savoldi, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2020. Breeding gender-aware direct speech translation systems. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3951–3964.

Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.

Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Loic Barrault, Huei-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2015. On using monolingual corpora in neural machine translation. *arXiv preprint arXiv:1503.03535*.

Caglar Gulcehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, and Yoshua Bengio. 2017. On integrating a language model into neural machine translation. *Computer Speech & Language*, 45:137–148.

Saša Hasan, Richard Zens, and Hermann Ney. 2007. Are very large n-best lists useful for smt? In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Companion Volume, Short Papers*, pages 57–60.

Kenneth Heafield and Alon Lavie. 2010. Combining machine translation output with open source. *The carnegie mellon multi-engine machine translation scheme. The Prague Bulletin of Mathematical Linguistics*, 93(1):27–36.

Christian Herold, Yingbo Gao, Mohammad Zeineldeen, and Hermann Ney. 2023. Improving language model integration for neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7114–7123.

Hieu Hoang, Huda Khayrallah, and Marcin Junczys-Dowmunt. 2024. On-the-fly fusion of large language models and machine translation. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 520–532.

Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

Sébastien Jean and Kyunghyun Cho. 2020. Log-linear reformulation of the noisy channel model for document-level neural machine translation. In *Proceedings of the Fourth Workshop on Structured Prediction for NLP*, pages 95–101.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023a. Mistral 7b. *Preprint*, arXiv:2310.06825.

Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023b. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14165–14178.

Marcin Junczys-Dowmunt, Tomasz Dwojak, and Rico Sennrich. 2016. The amu-uedin submission to the wmt16 news translation task: Attention-based nmt models as feature functions in phrase-based smt. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 319–325.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2024. Preliminary wmt24 ranking of general mt systems and llms. *arXiv preprint arXiv:2407.19884*.

Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, et al. 2022. Findings of the 2022 conference on machine translation (wmt22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71.

Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2024. Madlad-400: A multilingual and document-level large audited dataset. *Advances in Neural Information Processing Systems*, 36.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2024. Llava-med: Training a large language-and-vision assistant

for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36.

Jiwei Li and Dan Jurafsky. 2016. Mutual information and diverse decoding improve neural machine translation. *arXiv preprint arXiv:1601.00372*.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. Improved baselines with visual instruction tuning. *Preprint*, arXiv:2310.03744.

Evgeny Matusov, Nicola Ueffing, and Hermann Ney. 2006. Computing consensus translation for multiple machine translation systems using enhanced hypothesis alignment. In *11th Conference of the European Chapter of the Association for Computational Linguistics*, pages 33–40.

GemmaTeam Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open models based on gemini research and technology. *Preprint*, arXiv:2403.08295.

Benjamin Minixhofer, Edoardo Maria Ponti, and Ivan Vulić. 2024. Zero-shot tokenizer transfer. *arXiv preprint arXiv:2405.07883*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the*

40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Frithjof Petrick, Christian Herold, Pavel Petrushkov, Shahram Khadivi, and Hermann Ney. 2023. Document-level language models for machine translation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 375–391, Singapore. Association for Computational Linguistics.

Maja Popović. 2016. chrf deconstructed: beta parameters and n-gram weights. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 499–504.

Matt Post. 2018. A call for clarity in reporting bleu scores. *arXiv preprint arXiv:1804.08771*.

Amy Pu, Hyung Won Chung, Ankur P Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. Learning compact metrics for mt. In *Proceedings of EMNLP*.

Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022a. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.

Ricardo Rei, Nuno M Guerreiro, Daan van Stigt, Marcos Treviso, Luísa Coheur, José GC de Souza, André FT Martins, et al. 2023. Scaling up cometkiwi: Unbabel-ist 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848.

Ricardo Rei, Marcos Treviso, Nuno M Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José GC de Souza, Taisiya Glushkova, Duarte Alves, Luísa Coheur, et al. 2022b. Cometkiwi: Ist-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.

Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. In *Proceedings of ACL*.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725.

Shannon Zejiang Shen, Hunter Lang, Bailin Wang, Yoon Kim, and David Sontag. 2024. Learning to decode collaboratively with multiple language models. *arXiv preprint arXiv:2403.03870*.

Felix Stahlberg, James Cross, and Veselin Stoyanov. 2018. Simple fusion: Return of the language model. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 204–211.

Felix Stahlberg, Danielle Saunders, Adrià de Gispert, and Bill Byrne. 2019. Cued@ wmt19: Ewc&lms. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 364–373.

Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, MA Zejun, and Chao Zhang. 2023. Salmonn: Towards generic hearing abilities for large language models. In *The Twelfth International Conference on Learning Representations*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Vipin Vijayan, Braeden Bowen, Scott Grigsby, Timothy Anderson, and Jeremy Gwinnup. 2024. The case for evaluating multimodal translation models on text datasets. *arXiv preprint arXiv:2403.03014*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024a. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. In *Forty-first International Conference on Machine Learning*.

Yangyifan Xu, Jinliang Lu, and Jiajun Zhang. 2024b. Bridging the gap between different vocabularies for llm ensemble. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7133–7145.

# A    Appendix

## A.1    Adapting Seamless

We use the gender annotations from Gaido et al. (2020) to select talks with feminine speaker pronouns and an equal amount of randomly sampled masculine talks that are in the training set. We use the huggingface transformer's library (Wolf et al., 2019) for fine-tuning Seamless. We use LoRA (Hu et al., 2021) to fine-tune Seamless on this data. We set the *rank* to 16, *lora_alpha* to 64 and *lora_dropout* to 0.1. We apply adapters on the following modules: *q_proj, v_proj, linear_q, linear_v*. We set *batch_size* to 16, *gradient_accumulation_steps* to 8 and train with *fp16* for 20 epochs validating at every 200 steps. The *learning_rate* is set to $1e^{-5}$. The other parameters are set to default in the transformers library.

## A.2    Adapting PaliGemma

We also fine-tune the PaliGemma model with the huggingface transformer's library (Wolf et al., 2019) but use Q-LoRA (Dettmers et al., 2023) with *4-bit* quantization as the vision models require more VRAM. We set the *rank* to 8, *lora_alpha* and *lora_dropout* to default. We apply adapters on the following modules: *q_proj, k_proj, v_proj, gate_proj, up_proj, down_proj*. We set *batch_size* to 2, *gradient_accumulation_steps* to 6 and train with *bf16* for 5 epochs validating at every 200 steps. The *learning_rate* is set to $2e^{-5}$ with *AdamW* optimizer. The other parameters are set to default in the transformers library.

## A.3    Hyper-parameter Tuning for Speech-Aware Translations

To find the re-ranking weight $\alpha$, we generate the 25-best list of *Madlad* and *Seamless* on the MuST-C development set. Then, we calculate the scores of the models on these hypothesis and perform a grid search to find the optimal weight. Here, $\alpha = 1$ means that the score is only from *Madlad* and $\alpha = 0.5$ means equal contribution. The grid search is plotted in Figure 2.

We see that $\alpha$ as 0.8 is always achieving higher scores. Furthermore, we see that using *Seamless* as generator (Figure 2b) leads to poor translation quality and $\alpha$ as 1. However, in the case of *Madlad* as a generator (Figure 2a), we see that $\alpha$ as 1 is not optimal showing that re-ranking with *Seamless* is indeed beneficial. Finally in the case of both models as generator (Figure 2c), we again see that $\alpha$ as 1 achieves highest quality showing that *Seamless* is not beneficial.

(a) Re-ranking with Madlad N-best list



(b) Re-ranking with Seamless N-best list



(c) Re-ranking with Joint N-best list

Figure 2: Grid Search on $\alpha$ with *Madlad* and *Seamless Bal* on the MuST-C development set with N-best lists from different generators and rankers.

# Author Index

1486