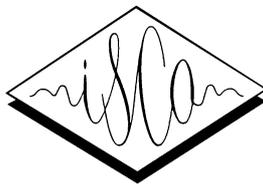
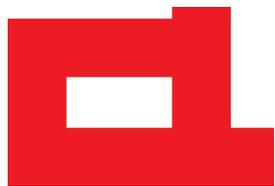


SIGDIAL 2024



**25th Annual Meeting of  
the Special Interest Group on  
Discourse and Dialogue**



**Proceedings of the Conference**

September 18 - 20, 2024  
Kyoto University, Kyoto, Japan

**In cooperation with:**

Association for Computational Linguistics (ACL)

International Speech Communication Association (ISCA)

Association for the Advancement of Artificial Intelligence (AAAI)

**We thank our sponsors:**

**Gold**



 SB Intuitions

 CyberAgent®

**PKSHA**  
TECHNOLOGY



**EQUMENOPOLIS**

**Silver**

The Google logo in its multi-colored font.

**Bronze**

 Fairy Devices

©2024 The Association for Computational Linguistics

Association for Computational Linguistics (ACL)  
209 N. Eighth Street  
Stroudsburg, PA 18360  
USA  
Tel: +1-570-476-8006  
Fax: +1-570-476-0860  
[acl@aclweb.org](mailto:acl@aclweb.org)

ISBN 979-8-89176-161-2

## Preface

We are glad to address the first few words for the proceedings of SIGDIAL 2024, the 25th Annual Meeting of the Special Interest Group on Discourse and Dialogue. The SIGDIAL conference is a premier venue for research publication and exchanges in discourse and dialogue. This year, the conference is organized as a fully in-person event on September 18-20, 2024, at Kyoto University, Kyoto, Japan.

The SIGDIAL 2024 program features two keynote talks, six sessions of oral presentations, including the special session on “GEMINI – Graph-based knowlEdge for Modelling Intelligent Natural Interaction”, two sessions of poster presentations and demonstrations, and a panel discussion. Two satellite workshops are held on the preceding days: the 2024 Young Researchers’ Roundtable on Spoken Dialog Systems (YRRSDS 2024) and the Workshop on Spoken Dialogue Systems for Cybernetic Avatars (SDS4CA).

SIGDIAL received a record of 156 submissions (excluding desk rejects and withdrawals) this year, comprising 106 long papers, 44 short papers, and 6 demo descriptions. Only a small minority of submissions was received via ARR (3 submissions). We had 15 Senior Program Committee (SPC) members, leading the discussion process and also writing meta-reviews. Each submission was assigned to an SPC member and received at least three reviews. Decisions carefully considered the original reviews, meta-reviews, and discussions among reviewers facilitated by the SPCs. We are immensely grateful to the members of the Program Committee and Senior Program Committee for their efforts in providing excellent, thoughtful reviews of the large number of submissions. Their contributions have been essential to selecting the accepted papers and providing a high-quality technical program for the conference. We have aimed to develop a broad, varied program spanning the many positively-rated papers identified by the review process. We therefore accepted 65 papers in total: 44 long papers (41.5%), 17 short papers (38.6%), and 4 demo descriptions, for an overall acceptance rate of 41.7%. The topics to be presented demonstrate the current breadth of research in discourse and dialogue.

We thank the two keynote speakers: Koji Inoue (Kyoto University, Kyoto) and Zhou Yu (Columbia University, New York, and Articulate.AI), for their inspiring talks. We also thank the organizers of the special session: “GEMINI – Graph-based knowlEdge for Modelling Intelligent Natural Interaction”. We are grateful for the mentoring chairs of Shikib Mehri and David Howcroft.

SIGDIAL 2024 is made possible by the dedication and hard work of our community. First, we express our gratitude to the SIGDIAL officers for their continuous support and advice. We also appreciate the volunteer works of all members of the organizing committee: Sponsorships Chair, Ramesh Manuvinakurike; Publication Chair, Kazunori Komatani; Publicity Chair, Ryuichiro Higashinaka, and Local Arrangements Chair, Koji Inoue. We gratefully acknowledge the support of our sponsors: Apple (Gold), SB Intuitions (Gold), CyberAgent (Gold), PKSHA (Gold), Equmenopolis (Gold), Google (Silver), and Fairy Devices (Bronze).

Finally, it is our great pleasure to welcome you to Kyoto and Kyoto University. Kyoto was the old capital of Japan from the 8th to the 19th century; thus, it has many national treasures and cultural heritages and yet incubates many innovative companies. Kyoto University is the second oldest (founded in 1897) national university in Japan, which encourages liberal and original research, resulting in the largest number of Nobel laureates in Asian institutes. You will meet robot interviewers at the conference. The banquet will be hosted in Heian Jingu Shrine, which models the old imperial palace. We hope that you will have an enjoyable and productive experience and leave with fond memories of SIGDIAL 2024. With our best wishes for a successful conference.

Tatsuya Kawahara, General Chair

Vera Demberg and Stefan Ultes, Program Co-Chairs

# Organizing Committee

## General Chair:

Tatsuya Kawahara, Kyoto University, Japan

## Program Chairs:

Vera Demberg, Saarland University, Germany

Stefan Ultes, University of Bamberg, Germany

## Local Arrangements Chair:

Koji Inoue, Kyoto University, Japan

## Sponsorship Chair:

Ramesh Manuvinakurike, Intel Labs, United States

## Mentoring Chairs:

Shikib Mehri, Contextual AI, United States

David Howcroft, Edinburgh Napier University, United Kingdom

## Publication Chair:

Kazunori Komatani, Osaka University, Japan

## Publicity Chair:

Ryuichiro Higashinaka, Nagoya University, Japan

## Social Media Chair:

Angus Addlesee, Heriot-Watt University, United Kingdom

## Local Arrangement Team:

Divesh Lala, Kyoto University, Japan

Keiko Ochi, Kyoto University, Japan

Mikey Elmers, Kyoto University, Japan

Mayumi Abe, Kyoto University, Japan

## SIGdial Officers:

President: Dilek Hakkani-Tur, University of Illinois Urbana-Champaign, United States

Vice President: Milica Gasic, University of Dusseldorf, Germany

Secretary: Alexandros Papangelis, Amazon, United States

Treasurer: Casey Kennington, Boise State University, United States

President Emeritus: Gabriel Skantze, KTH, Sweden

## Senior Program Committee:

Heriberto Cuayahuitl, University of Lincoln, United Kingdom

Kristiina Jokinen, AIRC, AIST, Japan

Casey Kennington, Boise State University, United States

Matthias Kraus, Augsburg University, Germany

Staffan Larsson, University of Gothenburg, Sweden

Pierre Lison, Norwegian Computing Centre, Norway

Wolfgang Maier, Mercedes-Benz AG, Germany

Mikio Nakano, C4A Research Institute, Inc., Japan  
Attapol Rutherford, Chulalongkorn University, Thailand  
Tatjana Scheffler, Ruhr University Bochum, Germany  
Gabriel Skantze, KTH Speech Music and Hearing, Sweden  
Svetlana Stoyanchev, Toshiba Europe, United Kingdom  
Maria Ines Torres, Universidad del Pais Vasco UPV/EHU, Spain  
Koichiro Yoshino, RIKEN, Nara Institute of Science and Technology, Japan  
Frances Yung, Saarland University Germany

**Program Committee:**

Angus Addlesee, Heriot-Watt University, United Kingdom  
Annalena Aicher, University of Augsburg, Germany  
Jan Alexandersson, DFKI GmbH, Germany  
Ron Artstein, USC Institute for Creative Technologies, United States  
Masayuki Asahara, National Institute for Japanese Language and Linguistics, Japan  
Timo Baumann, Ostbayerische Technische Hochschule Regensburg, Germany  
Frederic Bechet, Aix Marseille Universite - LIS/CNRS, France  
Alexander Berman, CLASP, University of Gothenburg, Sweden  
Nitu Bharati, TheOpenUniversity, United Kingdom  
Nathaniel Blanchard, Colorado State University, United States  
Nate Blaylock, Canary Speech, United States  
Jerome Boudy, Telecom SudParis, France  
Kevin Bowden, University of California Santa Cruz, United States  
Johan Boye, KTH, Sweden  
Kristy Boyer, University of Florida, United States  
Trung Bui, Adobe Research, United States  
Hendrik Buschmeier, Bielefeld University, Germany  
Zoraida Callejas, University of Granada, Spain  
Giuseppe Carenini, university of british columbia, Canada  
Justine Cassell, Carnegie Mellon University and Inria Paris, United States  
Alberto Cetoli, Private, United Kingdom  
Senthil Chandramohan, Staples, United States  
Akshay Chaturvedi, Institut de Recherche en Informatique de Toulouse, France  
Lin Chen, Meta Platform Inc., United States  
Zhiyu Chen, The University of Texas at Dallas, United States  
Alberto Chierici, New York University Abu Dhabi, United Arab Emirates  
Tahiya Chowdhury, Davis Institute for Artificial Intelligence, Colby College, United States  
Chenhui Chu, Kyoto University, Japan  
Luisa Coheur, INESC-ID/Instituto Superior Tecnico, Portugal  
Paul Crook, Meta, United States  
Souvik Das, University at Buffalo, United States  
Gael de Chalendar, CEA LIST, France  
Nina Dethlefs, University of Hull, United Kingdom  
David DeVault, Anticipant Speech, Inc., United States  
Bosheng Ding, Nanyang Technological University, Singapore  
Rama Sanand Doddipatla, Toshiba Cambridge Research Laboratory, United Kingdom  
Ondrej Dusek, Charles University, Czech Republic  
Mikey Elmers, Kyoto University, Japan  
Alex Fabbri, Salesforce AI Research, United States  
Shutong Feng, Heinrich-Heine-Universitat Dusseldorf, Germany  
Elisa Ferracane, Abridge AI, Inc., United States

Yingxue Fu, University of St Andrews, United Kingdom  
Shinya Fujie, Chiba Institute of Technology, Japan  
Kotaro Funakoshi, Tokyo Institute of Technology, Japan  
Angel Garcia Contreras, RIKEN, Japan  
Kallirroi Georgila, University of Southern California Institute for Creative Technologies, United States  
Felix Gervits, US Army Research Laboratory, United States  
Emer Gilmartin, Trinity College Dublin, Ireland  
Jonathan Ginzburg, Universite Paris Cite, France  
Akhilesh Deepak Gotmare, Salesforce Research, Singapore  
Venkata Subrahmanyam Govindarajan, University of Texas at Austin, United States  
David Griol, University of Granada, Spain  
Yulia Grishina, Amazon, Germany  
David Gros, University of California - Davis, United States  
Prakhar Gupta, Carnegie Mellon University, United States  
joakim gustafson, KTH, Sweden  
Nizar Habash, New York University Abu Dhabi, United Arab Emirates  
Sherzod Hakimov, University of Potsdam, Germany  
Dilek Hakkani-Tur, UIUC, United States  
Jie He, University of Edinburgh, United Kingdom  
Michael Heck, Heinrich Heine University, Germany  
Behnam Hedayatnia, Amazon, United States  
Ryuichiro Higashinaka, Nagoya University/NTT, Japan  
Jeff Higginbotham, University at Buffalo, United States  
Julia Hirschberg, Columbia University in the City of New York, United States  
Julian Hough, Swansea University, United Kingdom  
Christine Howes, University of Gothenburg, Sweden  
Jessica Huynh, Carnegie Mellon University, United States  
Michimasa Inaba, The University of Electro-Communications, Japan  
Mert Inan, Northeastern University, United States  
Koji Inoue, Kyoto University, Japan  
Bahar Irfan, KTH Royal Institute of Technology, Sweden  
Aditya Joshi, UNSW, Australia  
Georgi Karadzhov, University of Cambridge, United Kingdom  
Alexey Karpov, St. Petersburg Federal Research Center of the Russian Academy of Sciences (SPC RAS), Russian Federation  
Tatsuya Kawahara, Kyoto University, Japan  
Seiya Kawano, RIKEN, Japan  
Simon Keizer, Toshiba Europe Ltd, United Kingdom  
Hongjin Kim, Konkuk University, Republic of Korea  
Kazunori Komatani, Osaka University, Japan  
Vasily Konovalov, ,  
Divesh Lala, Kyoto University, Japan  
Kornel Laskowski, Carnegie Mellon University, United States  
Fabrice Lefevre, Avignon Univ., France  
Oliver Lemon, Heriot-Watt University, United Kingdom  
Chuyuan Li, The University of British Columbia, Canada  
Mohan Li, Toshiba Europe Ltd, United Kingdom  
Sheng Li, NICT, Advanced Speech Technology Laboratory, Japan  
Xiangci Li, Amazon, The University of Texas at Dallas, United States  
Yu Li, Columbia University, United States

Hsien-chin Lin, Heinrich Heine University, Germany  
Yang Janet Liu, Georgetown University, United States  
Zhengyuan Liu, Institute for Infocomm Research, A\*STAR, Singapore  
Eduardo Lleida Solano, University of Zaragoza, Spain  
Wanqiu Long, The University of Edinburgh, United Kingdom  
Bogdan Ludusan, Bielefeld University, Germany  
Cristina Luna Jimenez, Chair for Human-Centered Artificial Intelligence - Uni Augsburg, Germany  
Brielen Madureira, University of Potsdam, Germany  
Vladislav Maraev, University of Gothenburg, Sweden  
Alessandro Mazzei, Universita degli Studi di Torino, Italy  
Michael McTear, Ulster University, United Kingdom  
Teruhisa Misu, Honda Research Institute USA, United States  
Anhad Mohanany, Google, United States  
Tasnim Mohiuddin, Nanyang Technological University, Singapore  
Satoshi Nakamura, Nara Institute of Science and Technology, Japan  
Anna Nedoluzhko, Charles University in Prague, Czech Republic  
Vincent Ng, University of Texas at Dallas, United States  
Noriki Nishida, RIKEN Center for Advanced Intelligence Project, Japan  
Douglas O'Shaughnessy, INRS-EMT (Univ. of Quebec), Canada  
Win Pa Pa, University of Computer Studies, Yangon, Myanmar  
Suraj Pandey, The Open University, United Kingdom  
Patrick Paroubek, University Paris-Saclay - CNRS - LISN, France  
Paul Piwek, The Open University, United Kingdom  
Ondrej Platek, Charles University in Prague, Faculty of Mathematics and Physics, Institute of Formal and Applied Linguistics, Czech Republic  
Massimo Poesio, Queen Mary University of London and University of Utrecht, United Kingdom  
Laurent Prevot, Aix Marseille Universite & CNRS, France  
Dongqi Pu, Saarland University, Germany  
Stephen Pulman, Apple Inc., United Kingdom  
Matthew Purver, Queen Mary University of London, United Kingdom  
James Pustejovsky, Brandeis University, United States  
Kun Qian, Columbia University, United States  
Chengwei Qin, Nanyang Technological University, Singapore  
Vikram Ramanarayanan, University of California, San Francisco, United States  
Mathieu Ravaut, Nanyang Technological University, Singapore  
Ehud Reiter, University of Aberdeen, United Kingdom  
Lina M. Rojas Barahona, Orange Innovation Research, France  
Antonio Roque, Spark Interventions, United States  
Rimvydas Rubavicius, The University of Edinburgh, United Kingdom  
Ramon Ruiz-Dolz, University of Dundee, United Kingdom  
Benjamin Ruppik, Heinrich Heine University Dusseldorf, Germany  
Saurav Sahay, Intel Labs, United States  
Sakriani Sakti, Nara Institute of Science and Technology / Japan Advanced Institute of Science and Technology (JAIST/NAIST), Japan  
Hugues Sansen, Telecom SudParis, France  
David Schlangen, University of Potsdam, Germany  
Phillip Schneider, Technical University of Munich, Germany  
Samira Shaikh, University of North Carolina at Charlotte, United States  
Anthony Sicilia, University of Pittsburgh, United States  
Georg Stemmer, Intel Corp., Germany  
Matthew Stone, Rutgers University, United States

Carl Strathearn, Edinburgh Napier University, United Kingdom  
Hiroaki Sugiyama, NTT Communication Science Labs., Japan  
Alessandro Suglia, Heriot-Watt University, United Kingdom  
Ekaterina Svikhnushina, EPFL, Switzerland  
Antonio Teixeira, DETI/IEETA, University of Aveiro, Portugal  
Alberto Testoni, University of Amsterdam, Netherlands  
Kate Thompson, IRIT, France  
Maria Ines Torres, Universidad del Pais Vasco UPV/EHU, Spain  
Isabel Trancoso, INESC-ID / IST Univ. Lisbon, Portugal  
David Traum, University of Southern California Institute for Creative Technologies, United States  
Yuma Tsuta, The University of Tokyo, Japan  
Jingxuan Tu, Brandeis University, United States  
Gokhan Tur, University of Illinois Urbana Champaign, United States  
Markku Turunen, Tampere University, Finland  
Carel van Niekerk, Heinrich Heine University, Germany  
David Vandyke, Apple, United Kingdom  
Renato Vukovic, Heinrich Heine University Dusseldorf, Germany  
Nicolas Wagner, Otto-Friedrich-University of Bamberg, Germany  
Marilyn Walker, University of California Santa Cruz, United States  
Nicholas Walker, University of Bamberg, Germany  
Yi-Chia Wang, Facebook AI, United States  
Leo Wanner, ICREA and Pompeu Fabra University, Spain  
Nigel Ward, University of Texas at El Paso, United States  
Taro Watanabe, Nara Institute of Science and Technology, Japan  
Bonnie Webber, University of Edinburgh, United Kingdom  
Michael White, The Ohio State University, United States  
Graham Wilcock, CDM Interact and University of Helsinki, Finland  
Chien-Sheng Wu, Salesforce, United States  
Qingyang Wu, Columbia University, United States  
Yanchao Yu, Edinburgh Napier University, United Kingdom  
Amir Zeldes, Georgetown University, United States  
Chao Zhang, Tsinghua University, China  
Tiancheng Zhao, Binjiang Institute of Zhejiang University, China  
Mingyang Zhou, Post Doctoral Research Scientist at Columbia University, United States  
Sarka Zikanova, Charles University, Prague, Czech Republic  
Ingrid Zukerman, Monash University, Australia

**Secondary Reviewers:**

Aye Mya Hlaing, Pier Balestrucci, Hay Mar Soe Naing, Renato Vukovic, Virgile Socal,  
Alena Velichko, Ahmed Njifenjou, Yahui Fu, Moritz Wolf, Danila Mamontov, Mina Ameli,  
Tanja Simeonovski

**Mentors:**

Nandita Shankar Naik, Contextual AI, United States  
Shuhaib Mehri, University of Illinois Urbana-Champaign (UIUC), United States  
Jessica Huynh, Carnegie Mellon University, United States  
Behnam Hedayatnia, Amazon, United States  
Tejas Srinivasan, University of Southern California, United States  
Simon Mile, ADAPT Research Centre, Dublin City University, Ireland

**Invited Speakers:**

Koji Inoue, Kyoto University, Japan

Zhou Yu, Columbia University, United States

## Table of Contents

<i>Dialogue Discourse Parsing as Generation: A Sequence-to-Sequence LLM-based Approach</i> Chuyuan Li, Yuwei Yin and Giuseppe Carenini .....	1
<i>Rhetorical Strategies in the UN Security Council: Rhetorical Structure Theory and Conflicts</i> Karolina Zaczynska and Manfred Stede .....	15
<i>Elaborative Simplification for German-Language Texts</i> Freya Hewett, Hadi Asghari and Manfred Stede .....	29
<i>Examining Gender and Power on Wikipedia through Face and Politeness</i> Adil Soubki, Shyne E. Choi and Owen Rambow .....	40
<i>ReALM: Reference Resolution as Language Modeling</i> Joel Ruben Antony Moniz, Soundarya Krishnan, Melis Ozyildirim, Prathamesh Saraf, Halim Cagri Ates, Yuan Zhang and Hong Yu .....	51
<i>Dialog Flow Induction for Constrainable LLM-Based Chatbots</i> Stuti Agrawal, Pranav Pillai, Nishi Uppuluri, Revanth Gangi Reddy, Sha Li, Gokhan Tur, Dilek Hakkani-Tur and Heng Ji .....	66
<i>Knowledge-Grounded Dialogue Act Transfer using Prompt-Based Learning for Controllable Open-Domain NLG</i> Alain Vazquez Risco, Angela Maria Ramirez, Neha Pullabhotla, Nan Qiang, Haoran Zhang, Marilyn Walker and Maria Ines Torres .....	78
<i>Incremental Learning for Knowledge-Grounded Dialogue Systems in Industrial Scenarios</i> Izaskun Fernandez, Cristina Aceta, Cristina Fernandez, Maria Ines Torres, Aitor Etxalar, Ariane Mendez, Maia Agirre, Manuel Torralbo, Arantza Del Pozo, Joseba Agirre, Egoitz Artetxe and Iker Altuna .....	92
<i>Anticipating Follow-Up Questions in Exploratory Information Search</i> Graham Wilcock .....	103
<i>Bridging Information Gaps in Dialogues with Grounded Exchanges Using Knowledge Graphs</i> Phillip Schneider, Nektarios Machner, Kristiina Jokinen and Florian Matthes .....	110
<i>"Keep up the good work!": Using Constraints in Zero Shot Prompting to Generate Supportive Teacher Responses</i> E. Margaret Perkoff, Angela Maria Ramirez, Sean von Bayern, Marilyn Walker and James Martin .....	121
<i>HelloThere: A Corpus of Annotated Dialogues and Knowledge Bases of Time-Offset Avatars</i> Alberto Chierici and Nizar Habash .....	139
<i>It Couldn't Help but Overhear: On the Limits of Modelling Meta-Communicative Grounding Acts with Supervised Learning</i> Brielen Madureira and David Schlangen .....	149
<i>Data Augmentation Integrating Dialogue Flow and Style to Adapt Spoken Dialogue Systems to Low-Resource User Groups</i> Zhiyang Qi and Michimasa Inaba .....	159

<i>StyEmp: Stylizing Empathetic Response Generation via Multi-Grained Prefix Encoder and Personality Reinforcement</i>	
Yahui Fu, Chenhui Chu and Tatsuya Kawahara . . . . .	172
<i>Multi-Criteria Evaluation Framework of Selecting Response-worthy Chats in Live Streaming</i>	
Zhantao Lai and Kosuke Sato . . . . .	186
<i>Generating Unexpected yet Relevant User Dialog Acts</i>	
Lucie Galland, Catherine Pelachaud and Florian Pecune . . . . .	192
<i>Training LLMs to Recognize Hedges in Dialogues about Roadrunner Cartoons</i>	
Amie Paige, Adil Soubki, John Murzaku, Owen Rambow and Susan E. Brennan . . . . .	204
<i>On the Controllability of Large Language Models for Dialogue Interaction</i>	
Nicolas Wagner and Stefan Ultes . . . . .	216
<i>Divide and Conquer: Rethinking Ambiguous Candidate Identification in Multimodal Dialogues with Pseudo-Labeling</i>	
Bhathiya Hemanthage, Christian Dondrup, Hakan Bilen and Oliver Lemon . . . . .	222
<i>Self-Emotion Blended Dialogue Generation in Social Simulation Agents</i>	
Qiang Zhang, Jason Naradowsky and Yusuke Miyao . . . . .	228
<i>Enhancing Model Transparency: A Dialogue System Approach to XAI with Domain Knowledge</i>	
Isabel Feustel, Niklas Rach, Wolfgang Minker and Stefan Ultes . . . . .	248
<i>Affect Recognition in Conversations Using Large Language Models</i>	
Shutong Feng, Guangzhi Sun, Nurul Lubis, Wen Wu, Chao Zhang and Milica Gasic . . . . .	259
<i>Sentiment-Aware Dialogue Flow Discovery for Interpreting Communication Trends</i>	
Patrícia Sofia Pereira Ferreira, Isabel Carvalho, Ana Alves, Catarina Silva and Hugo Gonçalo Oliveira . . . . .	274
<i>Analyzing and Enhancing Clarification Strategies for Ambiguous References in Consumer Service Interactions</i>	
Changling Li, Yujian Gan, Zhenrong Yang, Youyang Chen, Xinxuan Qiu, Yanni Lin, Matthew Purver and Massimo Poesio . . . . .	289
<i>Coherence-based Dialogue Discourse Structure Extraction using Open-Source Large Language Models</i>	
Gaetano Cimino, Chuyuan Li, Giuseppe Carenini and Vincenzo Deufemia . . . . .	297
<i>Transforming Slot Schema Induction with Generative Dialogue State Inference</i>	
James D. Finch, Boxin Zhao and Jinho D. Choi . . . . .	317
<i>Using Respiration for Enhancing Human-Robot Dialogue</i>	
Takao Obi and Kotaro Funakoshi . . . . .	325
<i>Interactive Dialogue Interface for Personalized News Article Comprehension</i>	
Tomoya Higuchi and Michimasa Inaba . . . . .	329
<i>Enhancing Dialogue Speech Recognition with Robust Contextual Awareness via Noise Representation Learning</i>	
Wonjun Lee, San Kim and Gary Geunbae Lee . . . . .	333

<i>Local Topology Measures of Contextual Language Model Latent Spaces with Applications to Dialogue Term Extraction</i>	
Benjamin Matthias Ruppik, Michael Heck, Carel van Niekerk, Renato Vukovic, Hsien-chin Lin, Shutong Feng, Marcus Zibrowius and Milica Gasic . . . . .	344
<i>Adaptive Open-Set Active Learning with Distance-Based Out-of-Distribution Detection for Robust Task-Oriented Dialog System</i>	
Sai Keerthana Goruganthu, Roland R. Oruche and Prasad Calyam . . . . .	357
<i>Dialogue Ontology Relation Extraction via Constrained Chain-of-Thought Decoding</i>	
Renato Vukovic, David Arps, Carel van Niekerk, Benjamin Matthias Ruppik, Hsien-chin Lin, Michael Heck and Milica Gasic . . . . .	370
<i>InteLLA: Intelligent Language Learning Assistant for Assessing Language Proficiency through Interviews and Roleplays</i>	
Mao Saeki, Hiroaki Takatsu, Fuma Kurata, Shungo Suzuki, Masaki Eguchi, Ryuki Matsuura, Kotaro Takizawa, Sadahiro Yoshikawa and Yoichi Matsuyama . . . . .	385
<i>Curriculum-Driven Edubot: A Framework for Developing Language Learning Chatbots through Synthesizing Conversational Data</i>	
Yu Li, Shang Qu, Jili Shen, Shangchao Min and Zhou Yu . . . . .	400
<i>Going beyond Imagination! Enhancing Multi-modal Dialogue Agents with Synthetic Visual Descriptions</i>	
Haolan Zhan, Sameen Maruf, Ingrid Zukerman and Gholamreza Haffari . . . . .	420
<i>User Review Writing via Interview with Dialogue Systems</i>	
Yoshiki Tanaka and Michimasa Inaba . . . . .	428
<i>Conversational Feedback in Scripted versus Spontaneous Dialogues: A Comparative Analysis</i>	
Ildiko Pilan, Laurent Prévot, Hendrik Buschmeier and Pierre Lison . . . . .	440
<i>Exploring the Use of Natural Language Descriptions of Intents for Large Language Models in Zero-shot Intent Classification</i>	
Taesuk Hong, Youbin Ahn, Dongkyu Lee, Joongbo Shin, Seungpil Won, Janghoon Han, Stanley Jungkyu Choi and Jungyun Seo . . . . .	458
<i>Voice and Choice: Investigating the Role of Prosodic Variation in Request Compliance and Perceived Politeness Using Conversational TTS</i>	
Eva Szekely, Jeff Higginbotham and Francesco Possemato . . . . .	466
<i>A Dialogue Game for Eliciting Balanced Collaboration</i>	
Isidora Jeknic, David Schlangen and Alexander Koller . . . . .	477
<i>Improving Speech Recognition with Jargon Injection</i>	
Minh-Tien Nguyen, Dat Phuoc Nguyen, Tuan-Hai Luu, Xuan-Quang Nguyen, Tung-Duong Nguyen and Jeff Yang . . . . .	490
<i>Optimizing Code-Switching in Conversational Tutoring Systems: A Pedagogical Framework and Evaluation</i>	
Zhengyuan Liu, Stella Xin Yin and Nancy Chen . . . . .	500
<i>ECoh: Turn-level Coherence Evaluation for Multilingual Dialogues</i>	
John Mendonca, Isabel Trancoso and Alon Lavie . . . . .	516

<i>An Investigation into Explainable Audio Hate Speech Detection</i> Jinmyeong An, Wonjun Lee, Yejin Jeon, Jungseul Ok, Yunsu Kim and Gary Geunbae Lee . . . .	533
<i>Mhm... Yeah? Okay! Evaluating the Naturalness and Communicative Function of Synthesized Feedback Responses in Spoken Dialogue</i> Carol Figueroa, Marcel de Korte, Magalie Ochs and Gabriel Skantze . . . . .	544
<i>Generalizing across Languages and Domains for Discourse Relation Classification</i> Peter Bourgonje and Vera Demberg . . . . .	554
<i>BoK: Introducing Bag-of-Keywords Loss for Interpretable Dialogue Response Generation</i> Suvodip Dey and Maunendra Sankar Desarkar . . . . .	566
<i>Cross-lingual Transfer and Multilingual Learning for Detecting Harmful Behaviour in African Under-Resourced Language Dialogue</i> Tunde Oluwaseyi Ajayi, Mihael Arcan and Paul Buitelaar . . . . .	579
<i>A Few-shot Approach to Task-oriented Dialogue Enhanced with Chitchat</i> Armand Stricker and Patrick Paroubek . . . . .	590
<i>Exploration of Human Repair Initiation in Task-oriented Dialogue: A Linguistic Feature-based Approach</i> Anh Ngo, Dirk Heylen, Nicolas Rollet, Catherine Pelachaud and Chloé Clavel . . . . .	603
<i>Comparing Pre-Trained Embeddings and Domain-Independent Features for Regression-Based Evaluation of Task-Oriented Dialogue Systems</i> Kallirroi Georgila . . . . .	610
<i>Question Type Prediction in Natural Debate</i> Zlata Kikteva, Alexander Trautsch, Steffen Herbold and Annette Hautli-Janisz . . . . .	624
<i>MemeIntent: Benchmarking Intent Description Generation for Memes</i> Jeongsik Park, Khoi P. N. Nguyen, Terrence Li, Suyesh Shrestha, Megan Kim Vu, Jerry Yining Wang and Vincent Ng . . . . .	631
<i>Automating PTSD Diagnostics in Clinical Interviews: Leveraging Large Language Models for Trauma Assessments</i> Sichang Tu, Abigail Powers, Natalie Merrill, Negar Fani, Sierra Carter, Stephen Doogan and Jinho D. Choi . . . . .	644
<i>DialBB: A Dialogue System Development Framework as an Educational Material</i> Mikio Nakano and Kazunori Komatani . . . . .	664
<i>A Multimodal Dialogue System to Lead Consensus Building with Emotion-Displaying</i> Shinnosuke Nozue, Yuto Nakano, Shoji Moriya, Tomoki Ariyama, Kazuma Kokuta, Suchun Xie, Kai Sato, Shusaku Sone, Ryohei Kamei, Reina Akama, Yuichiroh Matsubayashi and Keisuke Sakaguchi . . . . .	669
<i>PersonaCLR: Evaluation Model for Persona Characteristics via Contrastive Learning of Linguistic Style Representation</i> Michimasa Inaba . . . . .	674
<i>DiagESC: Dialogue Synthesis for Integrating Depression Diagnosis into Emotional Support Conversation</i> Seungyeon Seo and Gary Geunbae Lee . . . . .	686

<i>Infusing Emotions into Task-oriented Dialogue Systems: Understanding, Management, and Generation</i> Shutong Feng, Hsien-chin Lin, Christian Geishauser, Nurul Lubis, Carel van Niekerk, Michael Heck, Benjamin Matthias Ruppik, Renato Vukovic and Milica Gasic .....	699
<i>Estimating the Emotional Valence of Interlocutors Using Heterogeneous Sensors in Human-Human Dialogue</i> Jingjing Jiang, Ao Guo and Ryuichiro Higashinaka .....	718
<i>The Gap in the Strategy of Recovering Task Failure between GPT-4V and Humans in a Visual Dialogue</i> Ryosuke Oshima, Seitaro Shinagawa and Shigeo Morishima .....	728
<i>MindDial: Enhancing Conversational Agents with Theory-of-Mind for Common Ground Alignment and Negotiation</i> Shuwen Qiu, Mingdian Liu, Hengli Li, Song-Chun Zhu and Zilong Zheng .....	746
<i>An Open Intent Discovery Evaluation Framework</i> Grant Anderson, Emma Hart, Dimitra Gkatzia and Ian Beaver .....	760
<i>Toximatics: Towards Understanding Toxicity in Real-Life Social Situations</i> Mayukh Das and Wolf-Tilo Balke .....	770



# Conference Program

**Wednesday September 18**

**09:30–10:00**    **Opening**

**10:00–11:40**    **Oral session 1: Discourse** (Chair: Peter Bourgonje)

*Dialogue Discourse Parsing as Generation: A Sequence-to-Sequence LLM-based Approach*

Chuyuan Li, Yuwei Yin and Giuseppe Carenini

*Rhetorical Strategies in the UN Security Council: Rhetorical Structure Theory and Conflicts*

Karolina Zaczynska and Manfred Stede

*Elaborative Simplification for German-Language Texts*

Freya Hewett, Hadi Asghari and Manfred Stede

*Examining Gender and Power on Wikipedia through Face and Politeness*

Adil Soubki, Shyne E. Choi and Owen Rambow

*ReALM: Reference Resolution as Language Modeling*

Joel Ruben Antony Moniz, Soundarya Krishnan, Melis Ozyildirim, Prathamesh Saraf, Halim Cagri Ates, Yuan Zhang and Hong Yu

**11:40–13:00**    **Lunch**

**Wednesday September 18 (continued)**

**13:00–14:40 Oral session 2: Special Session (GEMINI – Graph-based knowEdge for Modelling Intelligent Natural Interaction) (Chair: Kristiina Jokinen)**

*Dialog Flow Induction for Constrainable LLM-Based Chatbots*

Stuti Agrawal, Pranav Pillai, Nishi Uppuluri, Revanth Gangi Reddy, Sha Li, Gokhan Tur, Dilek Hakkani-Tur and Heng Ji

*Knowledge-Grounded Dialogue Act Transfer using Prompt-Based Learning for Controllable Open-Domain NLG*

Alain Vazquez Risco, Angela Maria Ramirez, Neha Pullabhotla, Nan Qiang, Haoran Zhang, Marilyn Walker and Maria Ines Torres

*Incremental Learning for Knowledge-Grounded Dialogue Systems in Industrial Scenarios*

Izaskun Fernandez, Cristina Aceta, Cristina Fernandez, Maria Ines Torres, Aitor Etxalar, Ariane Mendez, Maia Agirre, Manuel Torralbo, Arantza Del Pozo, Joseba Agirre, Egoitz Artetxe and Iker Altuna

*Anticipating Follow-Up Questions in Exploratory Information Search*

Graham Wilcock

*Bridging Information Gaps in Dialogues with Grounded Exchanges Using Knowledge Graphs*

Phillip Schneider, Nektarios Machner, Kristiina Jokinen and Florian Matthes

**14:40–15:00 Sponsors**

**15:00–15:30 Coffee break**

**15:30–17:00 Poster session 1**

*"Keep up the good work!": Using Constraints in Zero Shot Prompting to Generate Supportive Teacher Responses*

E. Margaret Perkoff, Angela Maria Ramirez, Sean von Bayern, Marilyn Walker and James Martin

*HelloThere: A Corpus of Annotated Dialogues and Knowledge Bases of Time-Offset Avatars*

Alberto Chierici and Nizar Habash

*It Couldn't Help but Overhear: On the Limits of Modelling Meta-Communicative Grounding Acts with Supervised Learning*

Brielen Madureira and David Schlangen

**Wednesday September 18 (continued)**

*Data Augmentation Integrating Dialogue Flow and Style to Adapt Spoken Dialogue Systems to Low-Resource User Groups*

Zhiyang Qi and Michimasa Inaba

*StyEmp: Stylizing Empathetic Response Generation via Multi-Grained Prefix Encoder and Personality Reinforcement*

Yahui Fu, Chenhui Chu and Tatsuya Kawahara

*Multi-Criteria Evaluation Framework of Selecting Response-worthy Chats in Live Streaming*

Zhantao Lai and Kosuke Sato

*Generating Unexpected yet Relevant User Dialog Acts*

Lucie Galland, Catherine Pelachaud and Florian Pecune

*Training LLMs to Recognize Hedges in Dialogues about Roadrunner Cartoons*

Amie Paige, Adil Soubki, John Murzaku, Owen Rambow and Susan E. Brennan

*On the Controllability of Large Language Models for Dialogue Interaction*

Nicolas Wagner and Stefan Ultes

*Divide and Conquer: Rethinking Ambiguous Candidate Identification in Multi-modal Dialogues with Pseudo-Labeling*

Bhathiya Hemanthage, Christian Dondrup, Hakan Bilen and Oliver Lemon

*Self-Emotion Blended Dialogue Generation in Social Simulation Agents*

Qiang Zhang, Jason Naradowsky and Yusuke Miyao

*Enhancing Model Transparency: A Dialogue System Approach to XAI with Domain Knowledge*

Isabel Feustel, Niklas Rach, Wolfgang Minker and Stefan Ultes

*Affect Recognition in Conversations Using Large Language Models*

Shutong Feng, Guangzhi Sun, Nurul Lubis, Wen Wu, Chao Zhang and Milica Gasic

*Sentiment-Aware Dialogue Flow Discovery for Interpreting Communication Trends*

Patrícia Sofia Pereira Ferreira, Isabel Carvalho, Ana Alves, Catarina Silva and Hugo Gonçalves Oliveira

*Analyzing and Enhancing Clarification Strategies for Ambiguous References in Consumer Service Interactions*

Changling Li, Yujian Gan, Zhenrong Yang, Youyang Chen, Xinxuan Qiu, Yanni Lin, Matthew Purver and Massimo Poesio

**Wednesday September 18 (continued)**

*Coherence-based Dialogue Discourse Structure Extraction using Open-Source Large Language Models*

Gaetano Cimino, Chuyuan Li, Giuseppe Carenini and Vincenzo Deufemia

*Transforming Slot Schema Induction with Generative Dialogue State Inference*

James D. Finch, Boxin Zhao and Jinho D. Choi

*Using Respiration for Enhancing Human-Robot Dialogue*

Takao Obi and Kotaro Funakoshi

*Interactive Dialogue Interface for Personalized News Article Comprehension*

Tomoya Higuchi and Michimasa Inaba

*Scoring Coreference Chains with Split-Antecedent Anaphors*

Silviu Paun, Juntao Yu, Nafise Sadat Moosavi and Massimo Poesio  
(paper from Dialogue & Discourse journal)

**17:00–18:00** **Keynote 1: Koji Inoue** (Chair: Vera Demberg)

*Yeah, Well, Haha: Generating Non-linguistic Behaviors For Human-like Conversational Robots*

**18:00–20:00** **Reception @Foyer**

## Thursday, September 19

**09:00–10:20 Oral session 3: ML for Dialogue** (Chair: Yun-Nung Chen)

*Enhancing Dialogue Speech Recognition with Robust Contextual Awareness via Noise Representation Learning*

Wonjun Lee, San Kim and Gary Geunbae Lee

*Local Topology Measures of Contextual Language Model Latent Spaces with Applications to Dialogue Term Extraction*

Benjamin Matthias Ruppik, Michael Heck, Carel van Niekerk, Renato Vukovic, Hsien-chin Lin, Shutong Feng, Marcus Zibrowius and Milica Gasic

*Adaptive Open-Set Active Learning with Distance-Based Out-of-Distribution Detection for Robust Task-Oriented Dialog System*

Sai Keerthana Goruganthu, Roland R. Oruche and Prasad Calyam

*Dialogue Ontology Relation Extraction via Constrained Chain-of-Thought Decoding*

Renato Vukovic, David Arps, Carel van Niekerk, Benjamin Matthias Ruppik, Hsien-chin Lin, Michael Heck and Milica Gasic

**10:20–10:40 Coffee break**

**10:40–11:40 Keynote 2: Zhou Yu** (Chair: Stefan Ultes)

*AI Agents beyond ChatGPT*

**11:40–13:00 Lunch**

**Thursday, September 19 (continued)**

**13:00–14:20 Oral session 4: Education / Tutoring** (Chair: Kallirroi Georgila)

*InteLLA: Intelligent Language Learning Assistant for Assessing Language Proficiency through Interviews and Roleplays*

Mao Saeki, Hiroaki Takatsu, Fuma Kurata, Shungo Suzuki, Masaki Eguchi, Ryuki Matsuura, Kotaro Takizawa, Sadahiro Yoshikawa and Yoichi Matsuyama

*Curriculum-Driven Edubot: A Framework for Developing Language Learning Chatbots through Synthesizing Conversational Data*

Yu Li, Shang Qu, Jili Shen, Shangchao Min and Zhou Yu

*Going beyond Imagination! Enhancing Multi-modal Dialogue Agents with Synthetic Visual Descriptions*

Haolan Zhan, Sameen Maruf, Ingrid Zukerman and Gholamreza Haffari

*User Review Writing via Interview with Dialogue Systems*

Yoshiki Tanaka and Michimasa Inaba

**14:20–14:30 Introduction to Heian Jingu**

**14:30–15:00 Coffee break**

Thursday, September 19 (continued)

15:00–16:30 Poster session 2

*Conversational Feedback in Scripted versus Spontaneous Dialogues: A Comparative Analysis*

Ildiko Pilan, Laurent Prévot, Hendrik Buschmeier and Pierre Lison

*Exploring the Use of Natural Language Descriptions of Intents for Large Language Models in Zero-shot Intent Classification*

Taesuk Hong, Youbin Ahn, Dongkyu Lee, Joongbo Shin, Seungpil Won, Janghoon Han, Stanley Jungkyu Choi and Jungyun Seo

*Voice and Choice: Investigating the Role of Prosodic Variation in Request Compliance and Perceived Politeness Using Conversational TTS*

Eva Szekely, Jeff Higginbotham and Francesco Possemato

*A Dialogue Game for Eliciting Balanced Collaboration*

Isidora Jeknic, David Schlangen and Alexander Koller

*Improving Speech Recognition with Jargon Injection*

Minh-Tien Nguyen, Dat Phuoc Nguyen, Tuan-Hai Luu, Xuan-Quang Nguyen, Tung-Duong Nguyen and Jeff Yang

*Optimizing Code-Switching in Conversational Tutoring Systems: A Pedagogical Framework and Evaluation*

Zhengyuan Liu, Stella Xin Yin and Nancy Chen

*ECoh: Turn-level Coherence Evaluation for Multilingual Dialogues*

John Mendonca, Isabel Trancoso and Alon Lavie

*An Investigation into Explainable Audio Hate Speech Detection*

Jinmyeong An, Wonjun Lee, Yejin Jeon, Jungseul Ok, Yunsu Kim and Gary Geunbae Lee

*Mhm... Yeah? Okay! Evaluating the Naturalness and Communicative Function of Synthesized Feedback Responses in Spoken Dialogue*

Carol Figueroa, Marcel de Korte, Magalie Ochs and Gabriel Skantze

*Generalizing across Languages and Domains for Discourse Relation Classification*

Peter Bourgonje and Vera Demberg

*BoK: Introducing Bag-of-Keywords Loss for Interpretable Dialogue Response Generation*

Suvodip Dey and Maunendra Sankar Desarkar

**Thursday, September 19 (continued)**

*Cross-lingual Transfer and Multilingual Learning for Detecting Harmful Behaviour in African Under-Resourced Language Dialogue*

Tunde Oluwaseyi Ajayi, Mihael Arcan and Paul Buitelaar

*A Few-shot Approach to Task-oriented Dialogue Enhanced with Chitchat*

Armand Stricker and Patrick Paroubek

*Exploration of Human Repair Initiation in Task-oriented Dialogue: A Linguistic Feature-based Approach*

Anh Ngo, Dirk Heylen, Nicolas Rollet, Catherine Pelachaud and Chloé Clavel

*Comparing Pre-Trained Embeddings and Domain-Independent Features for Regression-Based Evaluation of Task-Oriented Dialogue Systems*

Kallirroi Georgila

*Question Type Prediction in Natural Debate*

Zlata Kikteva, Alexander Trautsch, Steffen Herbold and Annette Hautli-Janisz

*MemeIntent: Benchmarking Intent Description Generation for Memes*

Jeongsik Park, Khoi P. N. Nguyen, Terrence Li, Suyesh Shrestha, Megan Kim Vu, Jerry Yining Wang and Vincent Ng

*Automating PTSD Diagnostics in Clinical Interviews: Leveraging Large Language Models for Trauma Assessments*

Sichang Tu, Abigail Powers, Natalie Merrill, Negar Fani, Sierra Carter, Stephen Doogan and Jinho D. Choi

*DialBB: A Dialogue System Development Framework as an Educational Material*

Mikio Nakano and Kazunori Komatani

*A Multimodal Dialogue System to Lead Consensus Building with Emotion-Displaying*

Shinnosuke Nozue, Yuto Nakano, Shoji Moriya, Tomoki Ariyama, Kazuma Kokuta, Suchun Xie, Kai Sato, Shusaku Sone, Ryohei Kamei, Reina Akama, Yuichiroh Mat-subayashi and Keisuke Sakaguchi

**16:30–18:00** Excursion to Heian Jingu Shrine

**18:00–20:30** Banquet @Heian Jingu Kaikan

## Friday, September 20

**09:00–10:20 Oral session 5: Persona / Emotions** (Chair: Mikio Nakano)

*PersonaCLR: Evaluation Model for Persona Characteristics via Contrastive Learning of Linguistic Style Representation*

Michimasa Inaba

*DiagESC: Dialogue Synthesis for Integrating Depression Diagnosis into Emotional Support Conversation*

Seungyeon Seo and Gary Geunbae Lee

*Infusing Emotions into Task-oriented Dialogue Systems: Understanding, Management, and Generation*

Shutong Feng, Hsien-chin Lin, Christian Geishauser, Nurul Lubis, Carel van Niekerk, Michael Heck, Benjamin Matthias Ruppik, Renato Vukovic and Milica Gasic

*Estimating the Emotional Valence of Interlocutors Using Heterogeneous Sensors in Human-Human Dialogue*

Jingjing Jiang, Ao Guo and Ryuichiro Higashinaka

**10:20–10:40 Coffee break**

**10:40–11:40 Panel: Future of Dialogue Research**

**11:40–13:00 Lunch**

**Friday, September 20 (continued)**

**13:00–14:20 Oral session 6: Interaction** (Chair: Koichiro Yoshino)

*The Gap in the Strategy of Recovering Task Failure between GPT-4V and Humans in a Visual Dialogue*

Ryosuke Oshima, Seitaro Shinagawa and Shigeo Morishima

*MindDial: Enhancing Conversational Agents with Theory-of-Mind for Common Ground Alignment and Negotiation*

Shuwen Qiu, Mingdian Liu, Hengli Li, Song-Chun Zhu and Zilong Zheng

*An Open Intent Discovery Evaluation Framework*

Grant Anderson, Emma Hart, Dimitra Gkatzia and Ian Beaver

*Toximatics: Towards Understanding Toxicity in Real-Life Social Situations*

Mayukh Das and Wolf-Tilo Balke

**14:20–14:50 Coffee break**

**14:50–15:50 Business meeting**

**15:50–16:20 Award & Closing**

## Keynote Abstracts

### **Keynote 1 - Yeah, Well, Haha: Generating Non-linguistic Behaviors For Human-like Conversational Robots**

Koji Inoue

*Kyoto University*

#### **Abstract**

The rise of multimodal large language models (MLLMs) has notably enhanced the capabilities of spoken dialogue systems and conversational robots, paving the way for practical applications. Yet, these models still struggle with specific non-linguistic behaviors crucial for the fluidity and dynamism of human conversations. This keynote will delve into these essential behaviors—such as backchanneling, laughter, and turn-taking—tracing the evolution of research from early studies to the latest Transformer-based models. The discussion will also address the persistent challenges in spoken dialogue research, aiming to advance human-like conversational robots in the era of evolving MLLMs.

#### **Biography**

Koji Inoue received his Ph.D. from the Graduate School of Informatics, Kyoto University, Japan, in 2018. He is currently an assistant professor at Kyoto University. In 2023, he was a visiting researcher at KTH Royal Institute of Technology in Sweden. His research team has developed a spoken dialogue system for the android ERICA. He was awarded the NETEXPLO Innovation 2022 Award.

## **Keynote 2 - AI Agents Beyond ChatGPT**

Zhou Yu

*Columbia University and Articulate.AI*

### **Abstract**

ChatGPT has significantly raised public expectations for conversational agents, with many now anticipating these agents to handle a wide range of tasks. However, deploying one single larger model with generalized capabilities is often impractical, in terms of accuracy, cost, and security, particularly in industry settings. Solving specific tasks requires a systematic combination of different models to form workflows. In this talk, we will explore various approaches to developing smaller, open-source models that can power AI agents to perform specialized tasks more effectively, using diverse fine-tuning techniques. In addition, we will talk about how AI Agent frameworks such as reflection could be applied in smaller model settings.

### **Biography**

Zhou(Jo) Yu is an Associate Professor at Columbia University's Computer Science Department. She obtained her Ph.D. from Carnegie Mellon University. Dr. Yu has received several best paper awards in top NLP conferences and has won Forbes 30 under 30 in 2018. Dr. Yu has developed various dialog system applications that have had a real impact, including winning the Amazon Alexa Prize. Dr. Yu co-founded Articulate.ai INC, democratizing AI Agent building with GenAI developer tools.

# Dialogue Discourse Parsing as Generation: a Sequence-to-Sequence LLM-based Approach

Chuyuan Li, Yuwei Yin, Giuseppe Carenini

Department of Computer Science  
University of British Columbia  
Vancouver, BC, Canada, V6T 1Z4  
chuyuan.li@ubc.ca; {yuweiyin, carenini}@cs.ubc.ca

## Abstract

Discourse analysis studies the sentence organization within a document, aiming to reveal its underlying structural information. Existing works on dialogue discourse parsing mostly use encoder-only models and sophisticated decoding strategies to extract structures. Despite recent advances in Large Language Models (LLMs), applying directly these models on discourse parsing is challenging. To fully leverage the rich semantic and discourse knowledge in LLMs, we propose to transform discourse parsing into a generation task using a text-to-text paradigm. Our approach is intuitive and requires no modification of the LLM architecture. Experimental results on STAC and Molwani datasets show that a sequence-to-sequence model such as T0 can perform reasonably well. Notably, our improved transition-based sequence-to-sequence system achieves new state-of-the-art performance on Molwani. Furthermore, our systems can generate richer discourse structures such as graphs, whereas previous methods are mostly limited to trees.<sup>1</sup>

## 1 Introduction

Discourse parsing is a Natural Language Processing task that aims to retrieve a structure from a document. The discursive structure contains clause-like text spans (known as Elementary Discourse Units) and are linked by semantic-pragmatic relations such as *Elaboration* and *Acknowledgment*. It plays a crucial role in natural language understanding and has demonstrated its usefulness in various downstream applications such as summarization (Feng et al., 2021) and dialogue comprehension (He et al., 2021; Ma et al., 2023).

Existing works on Dialogue Discourse Parsing (DDP) suggest that task-specific models are necessary to achieve state-of-the-art (SOTA) performance (Chi and Rudnicky, 2022; Li et al., 2023a).

<sup>1</sup>Code is available at <https://github.com/chuyuanli/Seq2Seq-DDP>.

They are based on complex architectures constructed on top of encoder-only pre-trained language models (PLMs) such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019). These models present a few limitations. First, they require task-specific architectures which oftentimes involve heavy engineering of utterance embeddings and specialized decoding strategies. Second, the predicted structures are typically limited to *trees*, neglecting other rich representations such as directed acyclic graphs (Asher et al., 2016). Third, they do not leverage rich latent knowledge in more recent Large decoder-only and encoder-decoder Language Models (LLMs) (Brown et al., 2020; Sanh et al., 2022; Chowdhery et al., 2023; Touvron et al., 2023).

Such LLMs have shown remarkable abilities in a wide range of applications, from text understanding and generation to coding to reasoning (Bang et al., 2023; Bubeck et al., 2023), resulting in a shift in focus from relatively small encoder-only PLMs to large-scale encoder-decoder and decoder-only LLMs. LLMs see a great amount of data: T0 model (Sanh et al., 2022), for instance, is pretrained on the C4 corpus (Habernal et al., 2016) containing 356 billion tokens; they are pretrained on a mixture of downstream tasks such as multi-document question answering (Yang et al., 2018) and natural language inference (Bowman et al., 2015). Since many of these tasks require an understanding of the inter-sentence structure, we hypothesize that LLMs have good contextual representation for sentence-level reasoning (e.g., discourse analysis).

However, in our preliminary experiments, we found that directly prompting LLMs does not perform well on the DDP task, confirming similar observations by Chan et al. (2023) who applied zero-shot prompting and in-context learning methods but found poor performance with GPT-3.5.

In this paper, we ask the question: *how to effectively transform the discourse parsing task into a*

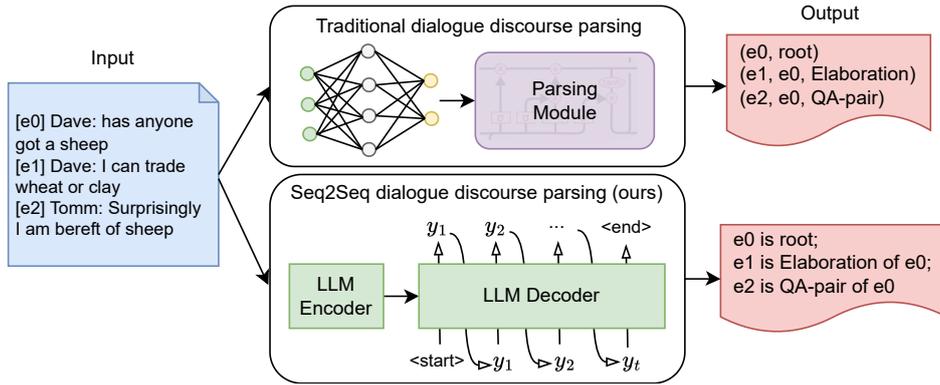


Figure 1: Traditional dialogue discourse parsing and our Seq2Seq dialogue discourse parsing systems.  $e_i$  denotes the discourse units and “QA-pair” represents the question-answer pair.

generation task?

To this end, we propose to tackle this problem within a text-to-text paradigm. We first formalize the parsing task as a Seq2Seq process and present a simple approach where a model takes a sequence of raw texts as input and produces a sequence of structures as output. We call this system **Seq2Seq-DDP**. The adopted model, such as T0, has a standard encoder-decoder architecture and is fine-tuned on parsing task. There is a great deal of flexibility in sequence representations, especially for the target sequence where *tree-like* and *graph-like* structures need to be expressed linearly. We design different schemes: one is close to natural language and another one is inspired by the *translation between augmented natural language* (TANL) formats (Paolini et al., 2021). This approach is straightforward, but it is constrained by weak supervision with lengthy inputs, which can lead to hallucinated or missing predictions for some utterances.

To tackle these issues, we propose to improve our system with transition-based algorithms which are widely used in dependency parsing (Nivre, 2003, 2008). A transition-based model receives the states of parsed sentences and the target sentence and predicts an action corresponding to the target sentence. A recent work on coreference resolution implemented such a system and achieved SOTA performance (Bohnet et al., 2023). Our enhanced system, **Seq2Seq-DDP+Transition**, processes one sentence at each step and predicts an action that establishes links and relations towards that sentence. We also adapt the sequence representations accordingly. Compared to the previous approach using full text input and output, the new system is more controllable with partial inputs and outputs.

We evaluate both systems on the STAC and Molwani datasets. The Seq2Seq-DDP model delivers promising results, matching the performance of SOTA models on Molwani. The transition-based system provides significant improvements across both datasets, setting new SOTA on Molwani. Through a series of analyses, we identify several key factors in converting a parsing task into a generation task, including the amount of supervision and the design of the representation scheme.

To summarize: (1) we propose a Seq2Seq-DDP method, along with an improved Seq2Seq-DDP+Transition variant, to transform discourse parsing into an LLM-based generation task, where our sophisticated sequence representations deliver promising performance gains; (2) we conduct extensive experiments and comprehensive analyses, which reveal insightful ideas on what makes a successful generative model for discourse parsing.

## 2 Related Work

**Discourse Parsing** Discourse parsing is a hard task, with low performance especially for multi-party dialogues which involve intricate relations between speakers, such as STAC (Asher et al., 2016) and Molwani (Li et al., 2020). Early approaches to discourse parsing used varied decoding strategies, such as Maximum Spanning Tree (Muller et al., 2012; Afantenos et al., 2012; Li et al., 2014) or Integer Linear Programming (Perret et al., 2016). Researchers soon applied neural models such as Gated Recurrent Units (Shi and Huang, 2019) and Graph Neural Networks (Wang et al., 2021b) to build contextualized embeddings, compared to hand-crafted features from the previous work. More recent works attempted to enhance the parsing task by utilizing Pre-trained Language

Models (PLMs) as backbone (Liu and Chen, 2021; Chi and Rudnicky, 2022), injecting external information such as speaker interactions (Yu et al., 2022; Li et al., 2023b), or joint learning with auxiliary tasks (Yang et al., 2021; He et al., 2021). Due to the small number of annotated examples, some also investigated semi-supervised approaches such as data programming (Badene et al., 2019), bootstrapping (Nishida and Matsumoto, 2022), and signals from the attention matrices in PLMs (Li et al., 2023a). However, much of this line of work dealt only with structure extraction while ignoring relations.

With LLMs on the scene, Chan et al. (2023) evaluated the performance of GPT-3.5 on discourse parsing using zero-shot and few-shot in-context-learning, but only to find that the model performs abysmally. Recently, Maekawa et al. (2024) employed decoder-only LLMs for Rhetorical Structure Theory (RST) discourse parsing in monologues, where conventional top-down and bottom-up strategies are transformed into prompts. On dialogues, only Wang et al. (2023) have investigated discourse parsing with a fine-tuned T5 model. However, their design of output sequences were overly simplified and we observed poor results with a similar abridged scheme in our experiments. In comparison, we explore the effectiveness of using Seq2Seq LLMs for this task with more sophisticated representations, such as an output closer to natural language.

### Structure Prediction with Generative Models

Loosely related to our work are papers about other structure prediction tasks which also apply generative modeling. For instance, on coreference resolution, Urbizu et al. (2020) conducted a proof-of-concept study where they literally translated the coreference annotation into a target sequence. Zhang et al. (2023) fine-tuned the T0 model with more sophisticated sequence representations that outperformed traditional coreference models. Bohnet et al. (2023) developed a transition-based Seq2Seq system based on mT5, which works on the same principle as our second approach. Paolini et al. (2021) proposed a unified framework that translates a series of structure tasks into *augmented natural languages* using T5. Their work aimed at creating a general and transferable model to solve many tasks. Generative models have also been used for semantic parsing (Rongali et al., 2020), syntactic parsing (He and Choi, 2023), and constituency parsing (Bai et al., 2023). Although

large generative models have been successfully applied to various structure prediction tasks, the DDP task, which requires inter-sentence reasoning in dialogues, remains under-explored.

## 3 A Formal Description of Discourse Parsing and Seq2Seq Modeling

### 3.1 Discourse Parsing

Given a document  $\mathcal{D} = \{e_0, e_1, \dots, e_n\}$  where  $e_i$  are clause-like text spans known as Elementary Discourse Units (EDU) and  $e_0$  is a dummy *root* node, the general goal of discourse parsing is to create a graph  $\mathcal{G}$  composed of  $(V, E, \ell)$  where  $V$  is a set of nodes or EDUs including  $\{e_0, e_1, \dots, e_n\}$ ,  $E_i \subset V \times V$  a set of edges pointing towards the node  $e_i$  with  $i \in [1, n]$ , and  $\ell$  a function  $\ell : (e_k, e_i) \mapsto r$  that maps an EDU pair with a rhetorical relation type  $r \in \mathbb{R}$ , with  $0 \leq k < i \leq n$ .

$$E_i = \{(e_k, e_i), e_i \in V, e_k \in V\} \quad (1)$$

Every  $E_i$  contains at least one pair of EDUs *pointing to* the node  $e_i$ . Here, we emphasize the uni-direction of edges given that in a dialogue, there are no “backwards” edges such that an EDU  $e_k$  by speaker a is rhetorically and anaphorically dependent upon a further EDU  $e_i$  of speaker b. This is known as *Turn Constraint* in the Segmented Discourse Representation Theory (SDRT) (Asher and Lascarides, 2003; Afantenos et al., 2015). The combination of all  $E_i$  is the set of all potential EDU pairs in document  $\mathcal{D}$ .

$$E = \cup_{i=1}^n E_i, \quad E_{\leq i} = \cup_{j=1}^i E_j \quad (2)$$

The equation 2 defines what we called *discourse structure prediction* where a “naked” graph can be extracted from  $\mathcal{D}$ . For *full parsing*, each edge must be assigned a relation with the function  $\ell$ . We can expand the pairs in  $E$  to triples in  $F$ :

$$F_i = \{(e_k, e_i, r_{ki}), e_i \in V, e_k \in V, r_{ki} \in \mathbb{R}\} \quad (3)$$

$$F = \cup_{i=1}^n F_i, \quad F_{\leq i} = \cup_{j=1}^i F_j \quad (4)$$

In a nutshell, discourse parsing takes a document  $\mathcal{D}$  as input and predicts the triples  $F$  as output. Assuming we have a training set of  $N$  examples,  $(\mathcal{D}_i, F_i)_{i=1}^N$  consists of  $N$  pairs of triples.

### 3.2 Seq2Seq Modeling

Let  $\mathcal{V}$  denote the vocabulary. Given a training pair  $(x, y)$  where  $x \in \mathcal{V}^{T'}$  is the source sequence of length  $T' \in \mathbb{N}$ ,  $y \in \mathcal{V}^T$  is the target sequence of length  $T \in \mathbb{N}$ , a Seq2Seq model computes the conditional probability  $p(y|x; \theta)$  autoregressively:

$$p(y|x; \theta) = \prod_{t=1}^T p(y_t|y_1, \dots, y_{t-1}, x; \theta) \quad (5)$$

Model parameters  $\theta$  are learned by maximizing the sum of conditional probabilities of all examples in the training set:

$$\theta^* = \arg \max_{\theta} \sum_{X, Y} \log p(Y|X; \theta) \quad (6)$$

### 3.3 Discourse Parsing as Seq2Seq Generation

To conduct discourse parsing with a Seq2Seq model, we translate  $(\mathcal{D}, F)$  into a pair of sequences  $(x, y)$ . The transformation from  $\mathcal{D}$  to  $x$  is straightforward since  $\mathcal{D}$  contains already a sequence of raw text. Our goal is to find a way to express  $F$  as a sequence  $y \in \mathcal{V}^T$ , which is also known as the “**linearization**” process for structured objects. A minimal approach is to literally predict the triples  $(e_k, e_i, r_{ik})$  in  $F$  as a sequence of strings. However, such a succinct format performs unsatisfactorily with limited training examples (see analysis in Section 6). We design several representation formats to explore a better solution for structure learning.

Another crucial issue is how to calculate the **conditional probability**  $p(y|x)$ . We can either feed  $x$  all at once and predict  $y$  in an end-to-end style or employ a transition system (Nivre, 2008), where the Seq2Seq model takes a single EDU as input and predicts an action corresponding to a set of discourse links involving that EDU as its output. In practice, we implement two Seq2Seq systems: a full text-in text-out system (Section 4) and an improved transition-based system (Section 5).

## 4 Seq2Seq Modeling for DDP

### 4.1 Methodology

**End-to-End System** A Seq2Seq-DDP system takes as input a document with raw text sequences and generates *structure-and-relation-labeled* output for each discourse unit autoregressively. Different from a classic pipeline approach where structure and relation are predicted subsequently (Afantenos et al., 2015; Shi and Huang, 2019; Liu and Chen, 2021; Li et al., 2024), our method jointly predicts link attachment  $(e_k, e_i)$  and relation  $(e_k, e_i) \mapsto r_{ki}$ .

**Representation Scheme** We investigate two output schemes: a natural scheme and an augmented scheme. For **natural scheme**, we hypothesize that the closer the output is to natural language, the more advantage the Seq2Seq model can take from

its pre-training. In other structure prediction tasks such as syntactic dependency parsing (He and Choi, 2023), natural language in the outputs has demonstrated its effectiveness. We use the following as a running example (*pilot01, STAC corpus*):

$\mathcal{D}$ : [e<sub>0</sub>] Dave: has anyone got a sheep, [e<sub>1</sub>] Dave: I can trade wheat or clay. [e<sub>2</sub>] Tomm: Surprisingly I am bereft of sheep.  
 $F$ : {(e<sub>0</sub>, e<sub>1</sub>, Elaboration), (e<sub>0</sub>, e<sub>2</sub>, QA-pair)}

We describe the triples in  $F$  with the template “ $e_i$  is  $r_{ki}$  of  $e_k$ ”:  $e_i$  and  $e_k$  are EDU markers;  $r_{ki}$  is a relation. In the input, we also append these markers as prefixes for each speech turn. The output joins all sequences with a semicolon. It reads:

$y_{nat}$ : [e<sub>0</sub>] is root; [e<sub>1</sub>] is Elaboration of [e<sub>0</sub>]; [e<sub>2</sub>] is Question-Answer-pair of [e<sub>0</sub>].

In cases where one node has multiple incoming edges, the template extends its tail to “ $e_i$  is  $r_{ki}$  of  $e_k$   $r_{mi}$  of  $e_m$   $r_{ni}$  of  $e_n$ ”, where  $e_m$  and  $e_n$  (resp.  $r_{mi}$  and  $r_{ni}$ ) are other linked nodes (resp. relations) to  $e_i$ . The advantage of this format is that each EDU uses exactly one sentence for structure description so that the length  $T$  of prediction  $y$  is fixed ( $T = T'$ ).

Inspired by the pioneering work on TANL (Paolini et al., 2021), we design an **augmented scheme**  $y_{aug}$  that replicates the input sentences and augments them with link and relation information:

$y_{aug}$ : [ Dave: has anyone got a sheep, | e<sub>0</sub> | root = e<sub>0</sub> ] [ Dave: I can trade wheat or clay | e<sub>1</sub> | Elaboration = e<sub>0</sub> ] [ Tomm: Surprisingly I am bereft of sheep. | e<sub>2</sub> | QA-pair = e<sub>0</sub> ]

Specifically, each EDU is enclosed by the special tokens [ ]. The pipe token | separates raw text, the EDU marker, and a list of relations in the format “ $r_{ki} = e_k$ ”. The EDU marker  $e_i$  is not preprend in the input. The model needs to use EDU markers to represent utterances and apply them on structure prediction. In other structure prediction tasks such as semantic role labeling (Paolini et al., 2021) and coreference resolution (Zhang et al., 2023; Bohnet et al., 2023), such a representation gives SOTA performance with Seq2Seq models.

**Decoding Structured Output** Once the model generates an output ( $y_{nat}$  or  $y_{aug}$ ), we decode the sentences to obtain  $F$  by following:

- Step1. Split the sequences with semicolons (resp. enclosed brackets) and remove all spe-

Scheme	STAC		Molweni		STAC		Molweni	
	Link	Full	Link	Full	Hallu	Miss	Hallu	Miss
Natural	65.6 $\pm$ 0.3	46.9 $\pm$ 1.8	81.4 $\pm$ 0.4	57.8 $\pm$ 0.1	3.1%	1.7%	0.4%	0
Augmented	66.7 $\pm$ 0.7	52.0 $\pm$ 0.1	82.4 $\pm$ 0.4	59.1 $\pm$ 1.0	0	0.2%	0	0

Table 1: Seq2Seq-DDP results on STAC and Molweni test sets (left) and error statistics (right). Scores are averaged micro-F<sub>1</sub>. “hallu” and “miss”: hallucinated and missed EDUs.

cial tokens (*is*, *of*, *|*, *=*) to extract triples in  $y_{nat}$  and quadruples in  $y_{aug}$ .

- Step2. Match the generated  $\hat{e}_i$  with the source  $e_i$  using heuristics. For  $y_{nat}$ , we match EDU markers; for  $y_{aug}$ , we match the input sentence and the cleaned output sentence at the token level using the Jaro distance (Jaro, 1989). We use 10 examples from the validation set in STAC and find that using the similarity value  $> 0.96$  can best cover the difference—most of times caused by more spacing between tokens—in generated and gold output. Once the  $\hat{e}_i$  and  $e_i$  is matched, we obtain the triples in  $(e_k, \hat{e}_i, r_{ki})$  which is the predicted structure for EDU  $e_i$ .
- Step3. Sanity check for *hallucinated* or *forgotten* EDUs in  $\hat{y}$ . The output sequence is designed in a way that its length matches the length of the input, so it is easy to spot erroneous generation. We introduce default rules for failure cases: remove the hallucination and add an adjacent attachment with a majority relation (i.e., *Question-answer-pair*) to the missed EDUs<sup>2</sup>.

We do not apply constrained decoding (Hokamp and Liu, 2017) as the output is well-aligned with the designed scheme and does not require extra vocabulary masking during generation.

## 4.2 Experimental Setup

We test our Seq2Seq-DDP system on two most commonly utilized datasets for dialogue discourse parsing: STAC (Asher et al., 2016) is composed of online multi-party conversations during the game *Settlers of Catan*. It contains 1,161 documents with in average 11 speech turns. We follow the subset split in Shi and Huang (2019) and set the maximum document length to 37, resulting in 911, 97, and 109 documents for training, validation, and testing, respectively. Molweni (Li et al., 2020) is a

<sup>2</sup>In reality, failure cases are few with a F<sub>1</sub>  $<$   $\pm$ 1%.

dataset derived from Ubuntu Chat Corpus (Lowe et al., 2015). It contains 10,000 documents with in average 8 utterances. We follow its original separation: 9,000 training, 500 validation, and 500 testing. Both corpora are annotated under the SDRT (Asher and Lascarides, 2003) and have the same relations ( $|\mathbb{R}| = 16$ ). We employ the traditional evaluation metrics, namely, the micro-averaged F<sub>1</sub> scores for link attachment ( $E$ ) and full structure ( $F$ ). All our experiments are conducted on T0 model (Sanh et al., 2022) with the 3B checkpoint, without any modification to the architecture. Most hyper-parameters in fine-tuning follow the suggestions in Raffel et al. (2020) (details in Appendix A).

## 4.3 Results and Analysis

The left part in Table 1 shows the parsing results on STAC and Molweni. Despite the simplicity of the Seq2Seq modeling, the fine-tuned T0 model can well perform dialogue discourse parsing, reaching 66–80 F<sub>1</sub> on the *naked* structure and 47–60 F<sub>1</sub> on the full structure. The outputs are well-aligned with the desired formats and only in rare cases do we observe erroneous generation (see below). Both *natural* and *augmented* formats produce satisfactory results on Molweni (link F<sub>1</sub>  $>$  81, full F<sub>1</sub>  $>$  57), whereas on STAC, we observe a more pronounced performance difference. The *natural* scheme is a succinct format that utilizes EDU markers in target sequences. This abridgment may cause ambiguity. In fact, the utterances in STAC are short (4.4 tokens/sentence) and similar texts can occur (e.g., the same answer from different speakers towards the same question). In comparison, *augmented* scheme replicates all tokens including speaker markers in the target sequence, helping to reduce ambiguity. Aligned with our observation, Paolini et al. (2021) also reported performance drops when using an abridged format for the entity and relation extraction task.

On the other hand, we observe a few problems originating from the Seq2Seq-DDP design, such as hallucinated or missed EDUs during generation,

as shown on the right part in Table 1. Since no explicit constraints are placed on the model’s output, there is potential for the model to produce invalid EDUs. However, this does not happen often: *natural* scheme generates 3% hallucinated and 1.7% missed EDUs on STAC (resp. 0.4% hallucinated and 0 missed on Molweni); while *augmented* scheme bypasses this issue completely. These erroneous outputs happen typically in longer documents when the number of speech turns exceeds thirty. In practice, we apply refinement rules in post-processing (included in Appendix B) to effectively eliminate this kind of generation.

## 5 Improve Seq2Seq-DDP Model with Transition-based Algorithm

An inherent drawback of the basic Seq2Seq-DDP system is the weak supervision in long sequences. The longer the document, the harder it is for the model to retrace previous predictions, as evidenced by the hallucinated or forgotten EDUs. Additionally, the act of consecutive output requires extra attention to some properties such as *counting*, which LLMs struggle with (Kojima et al., 2022). To provide more guidance during the generation and bypass the counting issue, we improve the Seq2Seq model with transition-based algorithms. The new Seq2Seq-DDP+Transition system takes a single EDU at each step and predicts an action corresponding to a set of links involving that EDU.

### 5.1 Methodology

**Transition-based System** The system we considered is closely related to the deterministic dependency parsing algorithm (Nivre, 2003, 2008). It starts with the dummy *root*  $e_0$  on the stack, all the EDUs in the buffer, and an empty set  $F$ . The parse ends once the buffer is empty and  $F$  contains triples of all EDUs (Equation 3). The transitions are composed of two actions: *link* action creates a right-arc from one EDU in the stack to the first EDU (i.e., target) in the buffer; *assign* action labels the arc. The target EDU in the buffer is then moved to the stack and a new round of transition will be conducted on the next EDU in the buffer.

**States.** A state  $c_i$  keeps track of which EDU is being processed through the index  $i$ , the established pairs  $E_{<i}$ , and associated relations  $F_{<i}$  up to  $i$ . We define the following states:

- $C$  is the set of all possible states.
- $c_s = (e_0, \epsilon, \epsilon)$  is the initial state, where two  $\epsilon$

are the empty sets  $E$  and  $F$ .

•  $C_t = \{c \in C \mid c = (e_n, E, F)\}$  is the set of the final states.

**Actions.** Given an intermediate state  $c_i = (e_i, E_{<i}, F_{<i})$ , we implement  $a_i$  which contains a link action  $\mathcal{L}(\cdot)$  and an assign action  $\mathcal{A}(\cdot)$ :

$$\mathcal{L}(e_i, F_{<i}) = \{e_k \rightarrow e_i, 0 \leq k < i\} \quad (7)$$

$$\mathcal{A}(e_i, E_i, F_{<i}) = \{(e_k \rightarrow e_i) \mapsto r_{ki}, r \in \mathbb{R}\} \quad (8)$$

The transition function  $\phi$  gives an updated state  $c_i$  accordingly:

$$\begin{aligned} & \phi(c_i, (e_k \rightarrow e_i), (e_k \rightarrow e_i) \mapsto r_{ki}) \\ &= (e_i, E_{<i} \oplus (e_k \rightarrow e_i), F_{<i} \oplus r_{ki}) \\ &= (e_i, E_i, F_i) \end{aligned} \quad (9)$$

Our transition system is a quadruple  $S = (C, c_s, T, C_t)$  where  $C$ ,  $c_s$ , and  $C_t$  are the states defined previously.  $T$  is the set of transitions, each of which is a function  $\phi : C \rightarrow C$ . The parsing path  $K$  is a sequence composed of actions and states:  $K = \{c_s, a_0, c_1, a_1, \dots, c_i, a_i, \dots, c_n\}$  where for  $i \in [1, n]$ ,  $c_{i+1} = \phi(c_i, a_i)$ , and where  $a_i = \mathcal{L}_i \cup \mathcal{A}_i$ ,  $c_n = C_t$ .

**Representation Scheme** Our goal is to encode the parsing path  $K$  into input and output strings. Specifically, each state-action pair  $(c_i, a_i)$  is mapped to an input-output pair  $(x_i, y_i)$ . Similar to Seq2Seq-DDP, we design output strings close to natural language. We illustrate two input-output pairs in the **natural scheme**, where the predicted action (underlined) is appended to the next state:

$x_1$ : [e<sub>0</sub>] [Dave: has anyone got a sheep,] is root;  
 $[e_1]$  [Dave: I can trade wheat or clay.] is  
 $y_{nat_1}$ : Elaboration of [e<sub>0</sub>]

$x_2$ : [e<sub>0</sub>] [Dave: has anyone got a sheep,] is root;  
 $[e_1]$  [Dave: I can trade wheat or clay.] is  
 Elaboration of [e<sub>0</sub>]; [e<sub>2</sub>] [Tomm: Surprisingly I am  
 bereft of sheep.] is  
 $y_{nat_2}$ : QA-pair of [e<sub>0</sub>].

We also implement a new format called **focused scheme** that utilizes special tokens **\*\*** to emphasize the target EDU ( $e_i$ ) and a pipe token **|** to separate the text with prediction, as depicted in Figure 2.

**Decoding and Sliding Window Strategy** Compared to the previous system, decoding the structured output from a transition-based model is easier: the generation is incremental with no mismatched or hallucinated EDUs. At each stage, we split  $\hat{y}$  on token *of* to obtain  $e_k$  and  $r_{ki}$ .

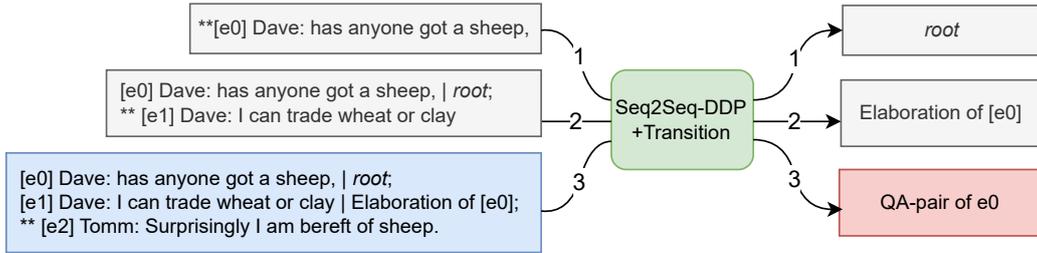


Figure 2: Seq2Seq-DDP+Transition system with *focused* scheme. It takes as input the previous state, the predicted action, and the next EDU; as output, actions for the current state. In blue: current input ( $c_i$ ); in red: current output ( $a_i$ ); in grey: parsed input ( $C_{<i}$ ).

The input grows longer as we continue adding the predicted structures. To comply with the maximum input length of pretrained models, we employ a sliding window strategy that reserves the closest EDUs for the next stage of prediction. Naturally, the closest EDUs are most relevant to the target EDU, so we frame a window with a set maximum length and slide it to the right at each stage. We set the window length to 18, as this is the longest link attachment in the validation set. The model is required to focus only on the target EDU  $e_i$  and its nearest preceding neighbors in the context  $c_i$ <sup>3</sup>.

## 5.2 Experiments and Analysis

We test our new system by fine-tuning T0-3B on STAC and Molweni datasets, results are shown in the first two rows in Table 2. Clearly, the transition-based system outperforms its Seq2Seq-DDP counterpart on all metrics: 5–8 and 1–3 points improvements on STAC and Molweni, respectively.

In the last four rows, we compare with the SOTA models (Shi and Huang, 2019; Liu and Chen, 2021; Chi and Rudnicky, 2022; Li et al., 2023c). Most of which use pre-trained language models such as RoBERTa to provide contextualized representations and task-specific techniques for decoding. Tellingly, our approach obtains new SOTA results on Molweni, surpassing the latest model proposed by Li et al. (2023c). We also achieve comparable results on STAC. Moreover, our approach is not limited to *tree-style* structures. Discourse-aware Seq2Seq models are capable of producing graphs (see Section 6). Although SOTA models use relatively small language models (110M - 340M parameters), it is important to point out that full comparability is challenging due to the numerous ways our approaches differ. First, the complexity of the parsing systems: SOTA models are built upon heavily

<sup>3</sup>In the oracle structures in test set, the longest edge distance is 13, so this approach does not affect any distant edges.

engineered architecture and require specific decoding strategies such as the Maximum Spanning Tree (MST). Our approach, on the other hand, directly leverages the standard encoder-decoder models and does not require any architecture modification. Second, scaling up encoder-only models does not always result in improvements in downstream applications. These models are also more difficult to deploy. Megatron-BERT (Shoeybi et al., 2019) with 1.3B and 3.9B parameters, for instance, are not publicly available. For generative models consisting of decoder networks, scaling tends instead to be closely associated with improved performance on many tasks (Ganguli et al., 2022).

Compared to Seq2Seq-DDP, the improved system does not suffer from EDU mismatch in the source and generation. However, the model sometimes predicts repetitive structures, such as “Acknowledgment of  $[e_2]$  Acknowledgment of  $[e_2]$ ”. In reality, failure cases are few: only 13 cases (1%) in all 1.2k triples in the development set. This occurs typically when the oracle output contains multiple incoming edges and the model tries to predict a graph structure.

## 6 Further Investigation

### 6.1 Masked Labels and Abridged Output

We investigate the influence of label semantics. The semantics of rhetorical relation types can be different in a pre-trained model. To prevent the model from understanding the relation through label semantics, we replace these words with special tokens, such as “rel1” and “rel2”, to the model vocabulary. This format is called  $y_{lmask}$ :

$y_{nat}$ : $[e_0]$ is <i>root</i> ; $[e_1]$ is Elaboration of $[e_0]$ ; $[e_2]$ is QA-pair of $[e_0]$ . <i>(baseline)</i>
$y_{lmask}$ : $[e_0]$ is <i>root</i> ; $[e_1]$ is rel4 of $[e_0]$ ; $[e_2]$ is rel0 of $[e_0]$ . <i>(label masked)</i>

System		STAC		Molweni	
		Link ( $\Delta$ )	Full ( $\Delta$ )	Link ( $\Delta$ )	Full ( $\Delta$ )
Natural (ours)	Seq2Seq-DDP+Transition	70.8 $\pm$ 0.9 ( $\uparrow$ 5.2)	55.1 $\pm$ 1.0 ( $\uparrow$ 8.2)	83.5 $\pm$ 0.2 ( $\uparrow$ 2.1)	60.3 $\pm$ 0.1 ( $\uparrow$ 2.5)
Focused (ours)	Seq2Seq-DDP+Transition	72.3 $\pm$ 0.6 ( $\uparrow$ 5.5)	56.6 $\pm$ 0.6 ( $\uparrow$ 4.6)	83.4 $\pm$ 0.6 ( $\uparrow$ 1.0)	60.0 $\pm$ 0.5 ( $\uparrow$ 0.9)
Shi and Huang (2019)	GRU+Pointer*	72.9 $\pm$ 0.4	54.2 $\pm$ 0.5	77.9 $\pm$ 0.4	54.1 $\pm$ 0.6
Liu and Chen (2021)	RoBERTa+Pointer	72.9 $\pm$ 1.5	57.0 $\pm$ 1.0	79.0 $\pm$ 0.4	55.4 $\pm$ 1.8
Chi and Rudnicky (2022)	RoBERTa+CLE $\dagger$	73.0 $\pm$ 0.5	58.1 $\pm$ 0.7	81.0 $\pm$ 0.7	58.6 $\pm$ 0.6
Li et al. (2023c)	BERT+BiAffine+Pointer	73.0	58.5	83.2	59.8

Table 2: Parsing results with our Seq2Seq-DDP+Transition models (top) and replicated SOTA models (bottom) on STAC and Molweni test sets. Scores are averaged micro-F<sub>1</sub>. Teal  $\uparrow$  shows performance gains compared to Seq2Seq-DDP systems. Pointer\*: pointer network (Vinyals et al., 2015). CLE $\dagger$ : Chu-Liu-Edmonds algorithm (Chu, 1965; Edmonds et al., 1967).

Additionally, to analyze the impact of sequence representations, we design abridged formats ( $y_{abr}$ ) for *natural* and *augmented* schema:

$y_{nat}$ : [ $e_0$ ] is root ; [ $e_1$ ] is Elaboration of [ $e_0$ ] ; [ $e_2$ ] is QA-pair of [ $e_0$ ].	(baseline)
$y_{abr}$ : [ $e_0$ ] root ; [ $e_1$ ] [ $e_0$ ] rel4 ; [ $e_2$ ] [ $e_0$ ] rel0.	(abridged)
$y_{aug}$ : [ Dave: has anyone got a sheep,   $e_0$   root = $e_0$ ] [ Dave: I can trade wheat or clay   $e_1$   Elaboration = $e_0$ ] [ Tomm: Surprisingly I am bereft of sheep.   $e_2$   QA-pair = $e_0$ ]	(baseline)
$y_{abr}$ : $e_0$   root = $e_0$ ; $e_1$   Elaboration = $e_0$ ; $e_2$   QA-pair = $e_0$ .	(abridged)

For the abridged version of *natural* representation, we transform the output into a triple  $(x, y, r)$  where  $x$  and  $y$  are respectively the dependent and head of an EDU pair;  $r$  is the masked relation type. It reads: EDU  $x$  is linked to EDU  $y$  with relation  $r$ . This is the expected output  $F$  from document  $\mathcal{D}$  (Equation 3), but such an extremely short linearization creates the most challenging representation: the model not only needs to learn the semantics of masked labels but also the implicit output pattern. For the abridged version of *augmented* representation, we do not repeat the input utterance and only keep EDU markers. The pipe (|) tag still denotes the start of the area of interest. Without the original text sequence, the abridged scheme requires extra reasoning to map the text with EDU markers.

We present the results of masked labels and abridged output in Table 3. On STAC, masking out the labels substantially hurt the performance with  $-2.5$  points in link prediction and  $-9.6$  in full. This demonstrates that label semantics are useful, especially for datasets containing smaller training examples. In terms of abridged output, both *natural abridged* and *augmented abridged* formulations underperform the baselines significantly ( $-12$  and  $-9.7$  points on full prediction). Interestingly, we do not observe a similar performance drop on Mol-

weni. Label-masked models obtain similar results as the *natural* baseline. The differences in link and full gains are not significant:  $p > 0.7$ ,  $p > 0.4$ . The most challenging abridged formulation also continues to perform well on Molweni. We think the amount of supervision is key. Molweni contains 9,000 documents in the training set whereas STAC only  $\approx 900$ . In terms of utterance length and token number, STAC is also very limited (see Table 5). These results are informative, indicating that a more “natural language”-like output generally brings more accurate predictions, especially when the amount of training data is low. On the other hand, sufficient supervision enables us to use the simpler paradigm of a text-to-text model successfully.

## 6.2 Pretrained LLMs and Model Sizes

We compare three LLMs in the T5 family: T5 (Raffel et al., 2020), Flan-T5 (Chung et al., 2022), and T0 (Sanh et al., 2022). In Table 4, we find that the model performance improves as the model size increases, which is in line with the observations in Zhang et al. (2023). In terms of different models in the T5 family, there is a notable difference between models with and without instruction finetuning such as FLAN (Wei et al., 2022). For models of the same size, the performance of the Flan-T5 and T0 is comparable (link 68.5 vs. 69.2; full 50.4 vs. 50.2), and both greatly exceed the performance of the original T5 model (+8 points in link attachment and +10 points in full prediction). Even the much smaller Flan-T5-base model (250M) outperforms T5-3B on link prediction by 2 points. This is not surprising: Chung et al. (2022) demonstrate that on some challenging BIG-Bench tasks (Srivastava et al., 2023), Flan-T5-11B outperforms the same size T5 by double-digit performances. This proves that instruction tuning can significantly enhance

Sequence representation	STAC		Molweni	
	Link (F <sub>1</sub> )	Full (F <sub>1</sub> )	Link (F <sub>1</sub> )	Full (F <sub>1</sub> )
Natural baseline	69.2 ± 0.5	50.2 ± 0.7	83.2 ± 1.4	58.6 ± 0.8
Label masked	↓ -2.5 ± 0.9	↓ -9.6 ± 0.4	↑ +0.3 ± 0.4	↑ +0.6 ± 0.5
Label masked + abridged	↓ -2.7 ± 0.2	↓ -12.4 ± 3.0	↑ +1.3 ± 1.0	↑ +0.6 ± 0.2
Augmented baseline	70.0 ± 0.8	54.2 ± 0.4	84.5 ± 0.4	59.0 ± 1.0
Abridged	↓ -2.6 ± 0.9	↓ -9.7 ± 0.4	~ ± 0.9	↑ +0.7 ± 1.1

Table 3: Sequence representation study on STAC and Molweni development sets. Red ↓, teal ↑, and ~ symbols refer to resp. lower, higher, and same scores compared to the baselines.

Pre-trained model	#Params	Link (F <sub>1</sub> )	Full (F <sub>1</sub> )
T5-large	738M	59.3 ± 0.6	36.4 ± 0.6
T5-3B	3B	60.7 ± 1.3	40.5 ± 0.9
Flan-T5-base	250M	63.0 ± 0.5	36.7 ± 0.1
Flan-T5-large	780M	67.2 ± 1.4	46.6 ± 1.8
Flan-T5-xl	3B	<u>68.5 ± 0.5</u>	<b>50.4 ± 0.1</b>
T0-3B	3B	<b>69.2 ± 0.5</b>	<u>50.2 ± 0.7</u>

Table 4: Study of different models in the T5 family on STAC development set (*natural* scheme). The best and second-best scores are **bolded** and underlined.

the model’s ability to learn complex language tasks, such as dialogue discourse parsing, thereby advancing it towards human-like language reasoning.

### 6.3 Richer Output Structures

We observe some distinctive features in the predicted structures such as directed acyclic graphs with Seq2Seq models. This is an exciting and big advantage over other SOTA models (Shi and Huang, 2019; Liu and Chen, 2021; Wang et al., 2021a; Chi and Rudnicky, 2022; Li et al., 2023a) that can only generate trees using MST algorithms in decoding (Eisner, 1996; Chu, 1965; Edmonds et al., 1967). Among all the proposed schemes, the *focused* scheme in Seq2Seq-DDP+Transition system achieves the highest performance in capturing multiple incoming edges, with a precision rate of 13% for graph structures. Other schemes such as *natural* and *augmented* also correctly predict around 10% graph structures. This is non-trivial: these structures are few and difficult to learn ( $\approx 5\%$  of nodes,  $< 7\%$  of links in STAC; none in Molweni) and demonstrate interesting and unique structures in dialogues.

### 6.4 Different Document Lengths

Since long documents can pose challenges for Seq2Seq models, we analyze the parsing performance under different document lengths, as shown

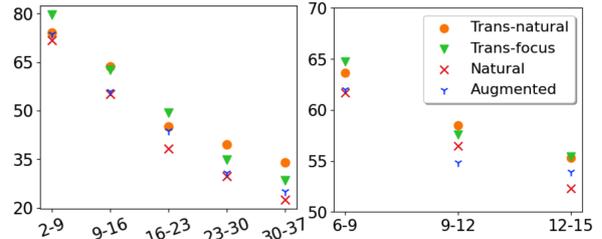


Figure 3: STAC (left) and Molweni (right) Full parsing performance under different Seq2Seq models and document lengths. x axis: #EDUs in a document. y axis: F1.

in Figure 3. On STAC, we split the length range into five even buckets between the shortest (2 EDUs) and longest (37 EDUs) document, resulting in 60, 25, 16, 4, and 4 data points per bucket. On Molweni, we split the documents into three buckets with 276, 154, and 70 data points in each group. Both the Seq2Seq-DDP and Seq2Seq-DDP+Transition systems exhibit a decline in performance with longer documents. However, our transition-based models (“Trans-\*) show a superior ability to handle long documents compared to their counterparts, as validated across both datasets.

## 7 Conclusion

We investigate an effective transformation approach for the DDP task by leveraging Seq2Seq LLMs. We adopt the pretrained encoder-decoder model T0 and fine-tune it to produce structured sequences. Without using any specific parsing module or modifying LLM architecture, our Seq2Seq-DDP system performs reasonably well on STAC and Molweni datasets. Excitingly, our Seq2Seq-DDP+Transition system yields comparable results with task-specific SOTA models, with richer discourse structures. Building on this work, we intend to explore various generative model architectures and sequence representations, and eventually extend our method to other discourse parsing tasks.

## Limitations

Longer documents tend to be more difficult to parse due to the growing number of possible discourse parse trees and the inherent drawbacks such as *counting* in LLMs. Our Transition-based systems mitigate this issue to some extent by using a sliding window strategy that focuses only on the closest EDUs.

In terms of decoding speed and performance, end2end systems demonstrate lower F<sub>1</sub> score but faster inference compared to transition-based systems. On the development set of STAC, the inference time for the end2end system is 2.5 seconds per document, whereas the transition-based system takes around 1.8 seconds per sequence, summing up to around 20 seconds for a complete document prediction.

## Acknowledgments

We thank the anonymous reviewers for their insightful comments and suggestions. We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC). Nous remercions le Conseil de recherches en sciences naturelles et en génie du Canada (CRSNG) de son soutien. The computing resources are provided by the Digital Research Alliance of Canada (alliance-can.ca).

## References

- Stergos Afantenos, Nicholas Asher, Farah Benamara, Anais Cadilhac, Cedric Dégremont, Pascal Denis, Markus Guhe, Simon Keizer, Alex Lascarides, Oliver Lemon, et al. 2012. [Modelling strategic conversation: model, annotation design and corpus](#). In *Proceedings of the 16th Workshop on the Semantics and Pragmatics of Dialogue (Seinedial)*, Paris.
- Stergos Afantenos, Eric Kow, Nicholas Asher, and Jérémy Perret. 2015. [Discourse parsing for multi-party chat dialogues](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 928–937, Lisbon, Portugal. Association for Computational Linguistics.
- Nicholas Asher, Julie Hunter, Mathieu Morey, Benamara Farah, and Stergos Afantenos. 2016. [Discourse structure and dialogue acts in multiparty dialogue: the STAC corpus](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 2721–2727, Portorož, Slovenia. European Language Resources Association (ELRA).
- Nicholas Asher and Alex Lascarides. 2003. *Logics of conversation*. Cambridge University Press.
- Sonia Badene, Kate Thompson, Jean-Pierre Lorré, and Nicholas Asher. 2019. [Data programming for learning discourse structure](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 640–645, Florence, Italy. Association for Computational Linguistics.
- Xuefeng Bai, Jialong Wu, Yulong Chen, Zhongqing Wang, and Yue Zhang. 2023. [Constituency parsing using llms](#). *arXiv preprint arXiv:2310.19462*.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. [A multitask, multilingual, multimodal evaluation of ChatGPT on reasoning, hallucination, and interactivity](#). In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 675–718, Nusa Dua, Bali. Association for Computational Linguistics.
- Bernd Bohnet, Chris Alberti, and Michael Collins. 2023. [Coreference resolution through a seq2seq transition-based system](#). *Transactions of the Association for Computational Linguistics*, 11:212–226.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. 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, et al. 2020. [Language models are few-shot learners](#). *Advances in neural information processing systems*, 33:1877–1901.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. [Sparks of artificial general intelligence: Early experiments with gpt-4](#). *arXiv preprint arXiv:2303.12712*.
- Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song. 2023. [Chatgpt evaluation on sentence level relations: A focus on temporal, causal, and discourse relations](#). *arXiv preprint arXiv:2304.14827*.
- Ta-Chung Chi and Alexander Rudnicky. 2022. [Structured dialogue discourse parsing](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 325–335.
- 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.
- Yoeng-Jin Chu. 1965. On the shortest arborescence of a directed graph. *Scientia Sinica*, 14:1396–1400.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- 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.
- Jack Edmonds et al. 1967. Optimum branchings. *Journal of Research of the national Bureau of Standards B*, 71(4):233–240.
- Jason M. Eisner. 1996. Three new probabilistic models for dependency parsing: An exploration. In *COLING 1996 Volume 1: The 16th International Conference on Computational Linguistics*.
- Xiachong Feng, Xiaocheng Feng, and Bing Qin. 2021. A survey on dialogue summarization: Recent advances and new frontiers. *arXiv preprint arXiv:2107.03175*.
- Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova Dassarma, Dawn Drain, Nelson Elhage, et al. 2022. Predictability and surprise in large generative models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1747–1764.
- Ivan Habernal, Omnia Zayed, and Iryna Gurevych. 2016. C4Corpus: Multilingual web-size corpus with free license. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 914–922, Portorož, Slovenia. European Language Resources Association (ELRA).
- Han He and Jinho D Choi. 2023. Unleashing the true potential of sequence-to-sequence models for sequence tagging and structure parsing. *Transactions of the Association for Computational Linguistics*, 11:582–599.
- Yuchen He, Zhuosheng Zhang, and Hai Zhao. 2021. Multi-tasking dialogue comprehension with discourse parsing. In *Proceedings of the 35th Pacific Asia Conference on Language, Information and Computation*, pages 551–561, Shanghai, China. Association for Computational Linguistics.
- Chris Hokamp and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1535–1546, Vancouver, Canada. Association for Computational Linguistics.
- Matthew A Jaro. 1989. Advances in record-linkage methodology as applied to matching the 1985 census of tampa, florida. *Journal of the American Statistical Association*, 84(406):414–420.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Chuyuan Li, Chloé Braud, Maxime Amblard, and Giuseppe Carenini. 2024. Discourse relation prediction and discourse parsing in dialogues with minimal supervision. In *Proceedings of the 5th Workshop on Computational Approaches to Discourse (CODI 2024)*, pages 161–176, St. Julians, Malta. Association for Computational Linguistics.
- Chuyuan Li, Patrick Huber, Wen Xiao, Maxime Amblard, Chloe Braud, and Giuseppe Carenini. 2023a. Discourse structure extraction from pre-trained and fine-tuned language models in dialogues. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 2562–2579, Dubrovnik, Croatia. Association for Computational Linguistics.
- Jiaqi Li, Ming Liu, Min-Yen Kan, Zihao Zheng, Zekun Wang, Wenqiang Lei, Ting Liu, and Bing Qin. 2020. Molweni: A challenge multiparty dialogues-based machine reading comprehension dataset with discourse structure. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2642–2652, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Jiaqi Li, Ming Liu, Yuxin Wang, Daxing Zhang, and Bing Qin. 2023b. A speaker-aware multiparty dialogue discourse parser with heterogeneous graph neural network. *Cognitive Systems Research*, 79:15–23.
- Jiwei Li, Rumeng Li, and Eduard Hovy. 2014. Recursive deep models for discourse parsing. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2061–2069.
- Wei Li, Luyao Zhu, Wei Shao, Zonglin Yang, and Erik Cambria. 2023c. Task-aware self-supervised framework for dialogue discourse parsing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14162–14173, Singapore. Association for Computational Linguistics.
- 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*.
- Zhengyuan Liu and Nancy Chen. 2021. [Improving multi-party dialogue discourse parsing via domain integration](#). In *Proceedings of the 2nd Workshop on Computational Approaches to Discourse*, pages 122–127, Punta Cana, Dominican Republic and Online. Association for Computational Linguistics.
- Ilya Loshchilov and Frank Hutter. 2018. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.
- Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. [The Ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems](#). In *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 285–294, Prague, Czech Republic. Association for Computational Linguistics.
- Xinbei Ma, Zhuosheng Zhang, and Hai Zhao. 2023. [Enhanced speaker-aware multi-party multi-turn dialogue comprehension](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Aru Maekawa, Tsutomu Hirao, Hidetaka Kamigaito, and Manabu Okumura. 2024. [Can we obtain significant success in RST discourse parsing by using large language models?](#) In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2803–2815, St. Julian’s, Malta. Association for Computational Linguistics.
- Philippe Muller, Stergos Afantenos, Pascal Denis, and Nicholas Asher. 2012. [Constrained decoding for text-level discourse parsing](#). In *Proceedings of COLING 2012*, pages 1883–1900, Mumbai, India. The COLING 2012 Organizing Committee.
- Noriki Nishida and Yuji Matsumoto. 2022. [Out-of-domain discourse dependency parsing via bootstrapping: An empirical analysis on its effectiveness and limitation](#). *Transactions of the Association for Computational Linguistics*, 10:127–144.
- Joakim Nivre. 2003. [An efficient algorithm for projective dependency parsing](#). In *Proceedings of the Eighth International Conference on Parsing Technologies*, pages 149–160, Nancy, France.
- Joakim Nivre. 2008. [Algorithms for deterministic incremental dependency parsing](#). *Computational Linguistics*, 34(4):513–553.
- Giovanni Paolini, Ben Athiwaratkun, Jason Krone, Jie Ma, Alessandro Achille, Rishita Anubhai, Cícero Nogueira dos Santos, Bing Xiang, and Stefano Soatto. 2021. [Structured prediction as translation between augmented natural languages](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Jérémy Perret, Stergos Afantenos, Nicholas Asher, and Mathieu Morey. 2016. [Integer linear programming for discourse parsing](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 99–109, San Diego, California. 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.
- Subendhu Rongali, Luca Soldaini, Emilio Monti, and Wael Hamza. 2020. [Don’t parse, generate! A sequence to sequence architecture for task-oriented semantic parsing](#). In *WWW ’20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*, pages 2962–2968. ACM / IW3C2.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. 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 V. 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 Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. [Multi-task prompted training enables zero-shot task generalization](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Zhouxing Shi and Minlie Huang. 2019. [A deep sequential model for discourse parsing on multi-party dialogues](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7007–7014.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. [Megatron-lm: Training multi-billion parameter language models using model parallelism](#). *arXiv preprint arXiv:1909.08053*.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adri Garriga-Alonso, et al. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *Transactions on machine learning research*.
- 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*.

- Gorka Urbizu, Ander Soraluze, and Olatz Arregi. 2020. [Sequence to sequence coreference resolution](#). In *Proceedings of the Third Workshop on Computational Models of Reference, Anaphora and Coreference*, pages 39–46, Barcelona, Spain (online). Association for Computational Linguistics.
- Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. [Pointer networks](#). *Advances in neural information processing systems*, 28.
- Ante Wang, Linfeng Song, Hui Jiang, Shaopeng Lai, Junfeng Yao, Min Zhang, and Jinsong Su. 2021a. [A structure self-aware model for discourse parsing on multi-party dialogues](#). In *Proceedings of the Thirtieth International Conference on International Joint Conferences on Artificial Intelligence*.
- Ante Wang, Linfeng Song, Lifeng Jin, Junfeng Yao, Haitao Mi, Chen Lin, Jinsong Su, and Dong Yu. 2023. [D 2 psg: Multi-party dialogue discourse parsing as sequence generation](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Jinfeng Wang, Longyin Zhang, and Fang Kong. 2021b. [Multi-level cohesion information modeling for better written and dialogue discourse parsing](#). In *Natural Language Processing and Chinese Computing - 10th CCF International Conference, NLPCC 2021, Qingdao, China, October 13-17, 2021, Proceedings, Part I*, volume 13028 of *Lecture Notes in Computer Science*, pages 40–52. Springer.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. [Finetuned language models are zero-shot learners](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Jingxuan Yang, Kerui Xu, Jun Xu, Si Li, Sheng Gao, Jun Guo, Nianwen Xue, and Ji-Rong Wen. 2021. [A joint model for dropped pronoun recovery and conversational discourse parsing in Chinese conversational speech](#). 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 1752–1763, Online. Association for Computational Linguistics.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Nan Yu, Guohong Fu, and Min Zhang. 2022. [Speaker-aware discourse parsing on multi-party dialogues](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 5372–5382, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Wenzheng Zhang, Sam Wiseman, and Karl Stratos. 2023. [Seq2seq is all you need for coreference resolution](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11493–11504, Singapore. Association for Computational Linguistics.

## A Experimental Setup

The data statistics are given in Table 5. All our experiments are conducted on T0 model (Sanh et al., 2022) with the 3B checkpoint: [https://huggingface.co/bigscience/T0\\_3B](https://huggingface.co/bigscience/T0_3B). The hyper-parameters for fine-tuning are kept as simple as possible. We do not apply parameter efficient fine-tuning techniques nor use lower precision during training. We apply a constant learning rate ( $5e-5$ ) using the AdamW optimizer (Loshchilov and Hutter, 2018). The mini-batch sizes are set to 4 for both *natural* and *augmented* schemes. The maximum input and output lengths are set to 512 and 1024. To fit in the positional embeddings of T0, we discard 36 and 6 documents in the STAC train and development sets, respectively. The actual training and development sets thus contain 911 and 97 documents, respectively. The test set is not affected. No document is discarded for Molweni. On Seq2Seq-DDP system, we train for a maximum of 20 epochs on STAC (resp. 10 epochs on Molweni) for 3B models, which takes about 5 hours (resp. 13 hours) on 1 A100 80G GPU. On Seq2Seq-DDP+Transition system, we train for a maximum of 10 epochs on STAC (resp. 5 epochs on Molweni), which takes around 12 hours (resp. 60 hours).

## B Seq2Seq-DDP System Examples of Erroneous Generation

Table 6 presents a few concrete examples of the error generations using Seq2Seq-DDP system. Specifically, we find that in *natural* scheme, 38 EDUs from 19 documents are [hallucinated](#); 22 EDUs are [not predicted \(missed\)](#) in 11 documents, all of which have an EDU index greater than 18. In the *augmented* scheme, there are no hallucinated EDUs; 3 EDUs are [not predicted](#). The typical error in this format comes from the wrong counting of large EDU index, as shown in [orange](#) in the Table.

Dataset	Train			Development			Test		
	#Doc	#Sent	#Token	#Doc	#Sent	#Token	#Doc	#Sent	#Token
STAC	911	10k	47k	97	1k	5k	109	1k	5k
Molweni	9000	79k	945k	500	4k	52k	500	4k	52k

Table 5: Data statistics for STAC and Molweni corpora. The numbers of documents, utterances, and tokens in training, development, and test sets are listed.

Error	$x$	$y$	$\hat{y}$
Hallu	$x_{nat}$ : ...[ $e_{13}$ ] Gaeilgeoir: I'll try again [ $e_{14}$ ] nareik15: cool :) [ $e_{15}$ ] Gaeilgeoir: I'm definitely giving wheat [ $e_{16}$ ] Gaeilgeoir: I have no sheep :)	$y_{nat}$ : ...[ $e_{14}$ ] is Acknowledgement of [ $e_{13}$ ] ; [ $e_{15}$ ] is Continuation of [ $e_{13}$ ] ; [ $e_{16}$ ] is Elaboration of [ $e_{15}$ ].	$\hat{y}_{nat}$ : [ $e_{14}$ ] is Acknowledgement of [ $e_{12}$ ] ; [ $e_{15}$ ] is <b>Result</b> of [ $e_{14}$ ] ; [ $e_{16}$ ] is <b>QA-pair</b> of [ $e_{15}$ ] ; [ $e_{17}$ ] is <b>Contrast</b> of [ $e_{16}$ ].
Miss	$x_{nat}$ : [ $e_0$ ] ztime: morning. [ $e_1$ ] Shawnus: hey [ $e_2$ ] Shawnus: good morning ...[ $e_{28}$ ] ztime: damn [ $e_{29}$ ] Shawnus: misplaced/ [ $e_{30}$ ] Shawnus: ? [ $e_{31}$ ] somdechn: Need to undo are you? [ $e_{32}$ ] ztime: no. [ $e_{33}$ ] ztime: you took the spot I was looking at. [ $e_{34}$ ] ztime: no it's fine. [ $e_{35}$ ] Shawnus: hahaha [ $e_{36}$ ] somdechn: Got to be mean here.	$y_{nat}$ : [ $e_0$ ] is root; [ $e_1$ ] is Acknowledgement of [ $e_0$ ] ; [ $e_2$ ] is Elaboration of [ $e_1$ ] ; ...[ $e_{29}$ ] is Clarification_question of [ $e_{28}$ ] ; [ $e_{30}$ ] is Correction of [ $e_{29}$ ] ; [ $e_{31}$ ] is Clarification_question of [ $e_{28}$ ] ; [ $e_{32}$ ] is QA-pair of [ $e_{29}$ ] ; [ $e_{33}$ ] is Explanation of [ $e_{32}$ ] ; [ $e_{34}$ ] is QA-pair of [ $e_{31}$ ] ; [ $e_{35}$ ] is Comment of [ $e_{32}$ ] ; [ $e_{36}$ ] is Comment of [ $e_{32}$ ].	$\hat{y}_{nat}$ : [ $e_0$ ] is root; [ $e_1$ ] is Acknowledgement of [ $e_0$ ] ; [ $e_2$ ] is <b>Continuation</b> of [ $e_1$ ] ; ...[ $e_{29}$ ] is <b>Comment</b> of [ $e_{28}$ ] ; [ $e_{30}$ ] is <b>Comment</b> of [ $e_{28}$ ] ; [ $e_{30}$ ] is <b>Comment</b> of [ $e_{28}$ ] ; [ $e_{30}$ ] is <b>Comment</b> of [ $e_{28}$ ] ; [ $e_{30}$ ] is <b>Comment</b> of [ $e_{28}$ ]
Count	$x_{aug}$ : [ ztime: morning ] [ Shawnus: hey ] [ Shawnus: good morning ] ... [ ztime: damn ] [ Shawnus: misplaced/ ] [ Shawnus: ? ] [ somdechn: Need to undo are you? ] [ ztime: no.. ] [ ztime: you took the spot I was looking at. ] [ ztime: no it's fine ] [ Shawnus: hahaha ] [ somdechn: Got to be mean here. ]	$y_{aug}$ : [ ztime: morning   $e_1$   root = $e_0$ ] [ Shawnus: hey   $e_1$   Acknowledgement = $e_0$ ] [ Shawnus: good morning   $e_2$   Elaboration = $e_1$ ] ... [ Shawnus: misplaced/   $e_{29}$   Clarification_question = $e_{28}$ ] [ Shawnus: ?   $e_{30}$   Correction = $e_{29}$ ] [ somdechn: Need to undo are you?   $e_{31}$   Clarification_question = $e_{28}$ ] [ ztime: no.   $e_{32}$   QA-pair = $e_{29}$ ] [ ztime: you took the spot I was looking at.   $e_{33}$   Explanation = $e_{32}$ ] [ ztime: no it's fine.   $e_{34}$   QA-pair = $e_{31}$ ] [ Shawnus: hahaha   $e_{35}$   Comment = $e_{32}$ ] [ somdechn: Got to be mean here.   $e_{36}$   Comment = $e_{32}$ ]	$\hat{y}_{aug}$ : [ ztime: morning   $e_1$   root = $e_0$ ] [ Shawnus: hey   $e_1$   Acknowledgement = $e_0$ ] [ Shawnus: good morning   $e_2$   <b>Continuation</b> = $e_1$ ] ... [ Shawnus: misplaced/   $e_{25}$   <b>QA-pair</b> = $e_{24}$ ] [ Shawnus: ?   $e_{25}$   <b>Continuation</b> = $e_{24}$ ] [ somdechn: Need to undo are you?   $e_{25}$   <b>Clarification_question</b> = $e_{24}$ ] [ ztime: no.   $e_{25}$   QA-pair = $e_{24}$ ] [ ztime: you took the spot I was looking at.   $e_{25}$   Explanation = $e_{24}$ ] [ ztime: no it's fine.   $e_{25}$   <b>Acknowledgement</b> = $e_{24}$ ] [ Shawnus: hahaha   $e_{25}$   Comment = $e_{24}$ ] [ Shawnus: hahaha   $e_{27}$   Comment = $e_{24}$ ] [ Shawnus: hahaha   $e_{27}$

Table 6: Error generation examples in STAC corpus.  $x$ ,  $y$ ,  $\hat{y}$  refer to resp. source input, target output, and generated output. ‘‘Hallu’’: hallucinated EDU in **teal**; ‘‘Miss’’: missing EDUs in **cyan**; ‘‘Count’’: wrong counting of EDU index in **orange**. False predictions are in **red**.

# Rhetorical Strategies in the UN Security Council: Rhetorical Structure Theory and Conflicts

Karolina Zaczynska and Manfred Stede

University of Potsdam, Applied Computational Linguistics

Potsdam, Germany

{lastname}@uni-potsdam.de

## Abstract

More and more corpora are being annotated with Rhetorical Structure Theory (RST) trees, often in a multi-layer scenario, as analyzing RST annotations in combination with other layers can lead to a deeper understanding of texts. To date, prior work on RST for the analysis of diplomatic language however, is scarce. We are interested in political speeches and investigate what rhetorical strategies diplomats use to communicate critique or deal with disputes. To this end, we present a new dataset with RST annotations of 82 diplomatic speeches aligned to existing Conflict annotations (UNSC-RST). We explore ways of using rhetorical trees to analyze an annotated multi-layer corpus, looking at both the relation distribution and the tree structure of speeches. In preliminary analyses we already see patterns that are characteristic for particular topics or countries.

## 1 Introduction

The United Nations Security Council (UNSC) meetings offer a unique longitudinal, cross-thematic resource on diplomatic interactions. Transcriptions of these meetings (Schönfeld et al., 2019) are a valuable corpus to study language use and communication style in an international relations context. In this paper, we study rhetorical style in diplomatic speech, by analyzing UNSC speeches from the perspective of Rhetorical Structure Theory (RST) (Mann and Thompson, 1988).

RST aims to capture the structure of a text by combining its elementary discourse units (EDUs) into one single, hierarchical tree structure. RST trees have proven to be useful in several downstream tasks, including characterizing genre distinctions (Sun et al., 2021; Liu and Zeldes, 2023), investigating text complexity (Hewett, 2023; Williams and Power, 2008) and fake news analysis (Rubin and Vashchilko, 2012; Popoola, 2017). However, little work has been done on RST in political and

diplomatic context, with a notable exception presented by Zeldes (2017). We address this gap by presenting a new corpus of 82 UNSC speeches annotated with RST trees. The resulting corpus (henceforth referred to as UNSC-RST) overlaps with our earlier work (Zaczynska et al., 2024), in which we annotated verbal Conflicts in UNSC speeches. In this paper, we present a multi-layer corpus of both RST trees and linguistic markers of Conflicts. We demonstrate how combining the two layers can reveal strategies in verbalizing disputes in a diplomatic setting. The main contributions of this paper are:

First, we present a new corpus with RST annotations for 82 diplomatic speeches from the UNSC. We adopt the RST annotation guidelines from earlier work (Carlson and Marcu, 2001; Zeldes, 2017; Stede et al., 2017), but make amendments tailored to the characteristics of diplomatic language. We include and discuss inter-annotator agreement, and publish our annotation guidelines.

Second, we combine our obtained RST annotations with earlier annotations of Conflict over the same texts, and use insights from argumentation analysis (Stede, 2016), to demonstrate how conclusions can be drawn on strategies to express Conflict. We compare the rhetorical style used by different countries (the five permanent members of the UNSC, plus Ukraine) and in different topics (debates concerning the situation in Ukraine, and the Women, Peace and Security agenda), and show, for example, that Conflicts are not as often supported by causal or justification relations as one might expect.

Our work provides an empirical basis for Political Science and International Relations researchers who are interested in understanding rhetorical styles used by representatives of different countries and in different contexts. The dataset, guidelines and code are available at: [https://github.com/linatal/rhetorical\\_UNSC](https://github.com/linatal/rhetorical_UNSC)

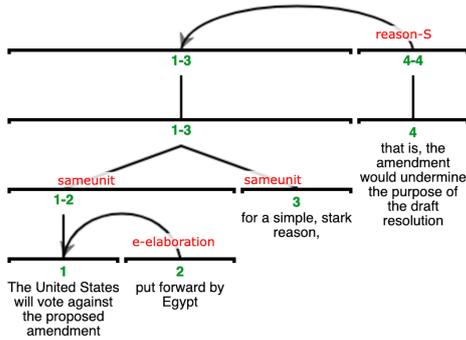


Figure 1: RST subtree from UNSC-RST (S/PV.7658, United States of America)

## 2 Background

This section first provides an overview of earlier work related to RST, and then describes the UNSC Conflicts Corpus that our work is based on.

### 2.1 RST Theory and Corpora

RST (Mann and Thompson, 1988) is a theory for analyzing the organization of texts and looks at discourse from an intention-driven perspective. It represents the structure of text in terms of coherence relations between text spans and captures the “plan” the author devised to influence their audience. Annotating texts with RST consists of two main steps: 1) segmenting the text into so-called Elementary Discourse Units (EDUs) and 2) organizing the EDUs into a single, hierarchical tree-structure. The result is a tree with hierarchically weighted EDUs, capturing the relative importance of each unit. Fig. 1<sup>1</sup> shows an RST tree with EDUs and discourse relations between EDUs. Most relations express a hierarchical relation between EDUs; they connect a less important EDU (called the *satellite*) to the more important one (the *nucleus*). In Fig. 1, EDU 4 is supporting the decision described in EDUs 1-3 by providing a REASON for the decision. Some relations, however, join equally-weighted EDUs, such as SAME-UNIT, which in the example connects two EDUs (1 and 3) that are interrupted by an E-ELABORATION (2).

Existing RST corpora such as the RST Discourse Treebank (RST-DT) (Carlson et al., 2001, 2002), the RST layer of the Georgetown University Multi-layer corpus (GUM) (Zeldes, 2017) and the RST layer of the Potsdam Commentary Corpus (henceforth: PCC-RST) (Stede et al., 2017) each come

<sup>1</sup>All RST examples are taken from UNSC-RST. We provide an official debate ID, beginning with S/PV and the country of the speaker, for each example.

with their own, slightly different versions of annotation guidelines. The guidelines of our UNSC-RST corpus are based on both the RST-DT and PCC-RST: For EDU segmentation, we use the RST-DT guidelines, and for relation annotation, we adopt (and slightly modify) the relation set from the PCC-RST (see Section 3.1 for more information on our relation set).

Our UNSC-RST corpus is an addition to the collection of RST-annotated texts, of which, to the best of our knowledge, only one covers texts from the political domain: The GUM corpus, since its v7.0.0 version, includes 15 speeches given in the UN General Assembly (16,720 tokens).<sup>2</sup> In comparison, the UNSC-RST corpus contains more speeches (82 vs. 15 in GUM) and more tokens (56,535 vs. 16,720 in GUM).

Obtaining RST trees automatically is the goal of RST parsing (Nguyen et al., 2021; Kobayashi et al., 2021; Liu and Zeldes, 2023), and RST trees have been used for downstream tasks such as text quality assessment (Skoufaki, 2020), summarization (Altmami and Menai, 2020), sentiment analysis (Kraus and Feuerriegel, 2019), and argument mining (Hewett et al., 2019).

### 2.2 The UNSC Conflicts Corpus

Our RST annotations are done over the same speeches as the Conflict annotations in the UNSC Conflicts corpus (UNSCCon) (Zaczynska et al., 2024). There, Conflicts are defined as verbalized disagreements or critique towards someone present at the UNSC debate (and the term Conflict does not refer to a military or physical conflict). There are different sub-types of Conflict:

- *Direct Negative Evaluations (Direct\_NegEval)* describe Conflicts where the speaker directly directs the critique to another country.  
Example: *This is a claim that takes Russia’s distortion of international law to a new level.* (S/PV.7165, United Kingdom and Northern Ireland)<sup>3</sup>
- *Indirect Negative Evaluations (Indirect\_NegEval)* describe Conflicts where some intermediate entity serving as a proxy is criticized instead of the other country directly. This can be done, for example, by criticizing

<sup>2</sup><https://github.com/amir-zeldes/gum/releases/tag/V7.0.0>.

<sup>3</sup>Examples are taken from UNSC debates on the situation in Ukraine.

a group acting on behalf of another country, or by criticizing a resolution the other country is supporting.

Example: *It is clear where responsibility lies: with the senseless violence of armed separatists and with those who have supported, equipped and advised them.* (S/PV.7165, United Kingdom and Northern Ireland)

- *Challenging* statements accuse another country of not telling the truth (see example below).
- *Corrections* rectify the allegedly false statement.

Example: *To conclude, one of our colleagues said that Kyiv had extended a hand to Moscow and that we had refused to reciprocate.* (*Challenge*)

*But the problem is not with Moscow; it has to do with the fact that Kyiv should have been the one to extend a hand to its people and regions, [...].* (*Correction*) (S/PV.7138, Russian Federation)

### 3 Annotations and Data

In the following, we describe our annotation guidelines, the annotation procedure, and corpus statistics.

#### 3.1 RST Guidelines Expansion

The first step in RST annotation is EDU segmentation. EDUs are sentences or smaller units (mostly clauses). Since in the UNSCon the speeches are already segmented into EDUs for its Conflicts annotation, we directly use their segmentation and refer to [Zaczynska et al. \(2024\)](#) for details on segmentation. The second step in RST annotation consists of choosing discourse relations to link EDUs. The next section describes our modifications to the PCC-RST relations guidelines.

##### 3.1.1 Additional Relations

We use the discourse relation set of ([Stede et al., 2017](#)), and include four additional relations (all taken from RST-DT, except for TOPIC-COMMENT, which is from GUM): SAME-UNIT, CONTRIBUTION, TEXTUAL-ORGANIZATION, and TOPIC-COMMENT. Since the sentence structure in the UNSC speeches is relatively complex (see [Zaczynska et al. \(2024, Table 1\)](#)) we found many cases where the EDU was interrupted by one or more embedded discourse units. To connect interrupted EDUs we use the SAME-UNIT relation. We also

include CONTRIBUTION, which serves to identify the speaker or source of a statement, because for the analysis of Conflicts it can be important to see whether speakers refer to other sources or to themselves (for example, when accusing someone of a false statement, like in *Challenge*-type Conflicts). We use TEXTUAL-ORGANIZATION to make links between different structural elements, such as between the title and the body of the text, or between a section heading and the following text. TOPIC-COMMENT is used for EDUs that do not contribute propositional content to the discourse, including back-channeling, incomplete utterances, and fillers.

##### 3.1.2 Merging Relations

In the guidelines by [Stede et al. \(2017\)](#), REASON and JUSTIFY both describe EDUs that aim to change the attitude of the reader. The difference is that for REASON, the claim is supported by a subjective assessment, while JUSTIFY describes a general basic attitude of the writer. Because this difference seems not relevant for our genre here, we decided to merge both relations and call them REASON.

##### 3.1.3 Rhetorical Questions

A particular challenge was the annotation of *rhetorical questions*, which appear quite frequently in the speeches. In RST-DT, they are labeled as RHETORICAL-QUESTION, which is a sub-type of TOPIC-COMMENT. However, ideally an RST relation should express the purpose of a unit in relation to another one, rather than characterizing a single unit in itself. Since rhetorical questions often have the purpose to emphasize for example a REASON for a claim, or the EVALUATION of a situation or statement, we decided to use these relations, instead of the general TOPIC-COMMENT relation. We only use TOPIC-COMMENT in cases where it is possible to remove the rhetorical question without losing essential information. For more details on the RST relations, we refer to the RST annotation guideline amendment provided in our repository.

#### 3.2 RST-Annotation Procedure

We used the RSTWeb annotation tool for tree building ([Zeldes, 2016](#)). Five annotators were trained for over a month for the first round of RST annotations. Then we conducted parallel annotations for a subset of 32 speeches, with two annotations per speech, based on the guidelines from [Stede et al. \(2017\)](#). For statistical evaluation we use the tool

RST-Tace (Wan et al., 2019), which is based on a qualitative method for comparing RST trees as described in Iruskieta et al. (2015). We computed inter-annotator agreement and found an overall average kappa of 0.44. The kappa score for nuclearity (defining the relative importance of an EDU) is 0.43; for relations it is 0.31; for constituents (the central nucleus) it is 0.43, and for attachment points (the direction of the relation) it is 0.51.

A confusion matrix providing more information about disagreements is given in Appendix A. Note that for the gold annotation we added four relations to the list of relations (see section 3.1.1). Most of the mismatches in the annotations can be related to semantic similarity of the chosen relations. For example, a frequent source of disagreement was LIST vs. CONJUNCTION. Both are essentially enumerating EDUs of the same importance, one using typographical connectors like commas or semicolons, the latter using conjunctions like *and* and *or*. Another frequent disagreement was between E-ELABORATION and ELABORATION. This has also been reported by Hewett (2023). Both relations state that the topic of the discourse is being continued in a more specific way, but for E-ELABORATION, the additional information is only on a single entity.

After we obtained the preliminary annotations for IAA calculation, we proceeded to form the adjudicated gold annotations. Two annotators (one is an author of this paper) annotated the entire corpus of 82 speeches, and continuously discussed progress via chat and in weekly meetings, thus creating the gold annotations according to the updated guidelines.

For the final trees, we decided to make use of the given paragraph breaks within the speech transcriptions. This means we first annotated adjacent EDUs for all paragraphs individually and then completed the tree for the whole speech. This way, we speeded up the annotation process for longer speeches. Another advantage was that it enables us to compare sub-tree structures and discourse relation distributions, as well as to find local most-important EDUs within the paragraphs (see Section 4).

### 3.3 UNSC-RST Corpus Statistics

The UNSC-RST corpus includes 85 speeches and therefore 85 RST trees with 60.87 EDUs per tree on average and 11.32 tokens per EDU on average (56,535 tokens in total). It covers almost all of

the speeches from the UNSCon.<sup>4</sup> The smallest tree has only seven EDUs (S/PV.7138\_spch016, Jordan), whereas the largest one has 194 EDUs (S/PV.7165\_spch019, Ukraine). There are six debates in total, covering two topics: Four debates (61 speeches) on the situation in Ukraine (from 2014), and two debates (24 speeches) on the "Women, Peace, and Security" agenda (both from 2016) dealing with gender aspects in security issues. The corpus includes 578 paragraphs, which are seven paragraphs on average per speech, with a maximum of 20 paragraphs per speech.

## 4 Methods

In this section, we describe the kinds of quantitative and qualitative analyses that we performed; the corresponding results will follow in the next section.

### 4.1 Distribution of Discourse Relations

Inspired by Popoola (2017); Hewett (2023) and others, we first look at the discourse relation distribution. We compare the frequency of RST relations and Conflict annotations per EDU on the leaf nodes (EDUs on the lowest level). In order to compare the distribution of relations between Conflicts, we look at the percentage of RST relations used per Conflict type. PCC-RST divides the set of RST relations into four groups according to their function: (1) *Pragmatic relations* serve to change the attitude of another person; (2) *semantic relations* describe states of affairs in the world; (3) *textual relations* organize the text and make its understanding easier; and (4) *multinuclear relations* enlist two or more EDUs of same importance in a relatively weak rhetorical relation. For our purposes here, we separately build the group of (5) *contrastive relations* that focus on differences or incompatibility of two propositions, often by weighting one as more important than the other. We have not assigned ATTRIBUTIONS to any group because they represent the purely formal action of marking reported speech, without additional rhetorical effect.

Since we are interested in how a Conflict is embedded in the text structure, we also compare the distribution of discourse relations within paragraphs. Thus we compare paragraphs with at least one Conflict annotation to those having no Conflict annotation.

<sup>4</sup>Two speeches were missing in the UNSC-RST at the time we conducted the experiments described in this paper.

We assume that diplomats use more pragmatic RST relations for Conflicts than for Non-Conflicts, because speakers can use pragmatic relations to motivate their criticism of another party, and to strengthen potential coalitions against the criticized position. They can also appeal to the criticized country to change their behavior or to take/refrain from a particular action. The results on relation distribution are in section 5.1.

## 4.2 Analyzing the Tree Structure: Nuclearity Mass Distribution

Besides relation distribution, we inspect the tree structure resulting from the RST annotation. The *central nucleus* (CN) is interpreted as the central statement of the text covered by the tree, and can be reached starting at the top of the tree by following only ‘nucleus’ edges towards the leaf nodes (Mann et al., 1992). Looking at the overall shape of the tree, we can distinguish between “deeper” RST trees that are centered around one core EDU to which there is a single distinctive longest path, and “flatter” trees that have several more or less equally weighted EDUs. Stede (2016) found that for short argumentative texts, deeper trees correlate with more strongly opinionated texts, in comparison to flat trees that can signal more descriptively-oriented text. Making use of the Conflict annotation for the analysis, we were interested in a potential difference between RST trees used for paragraphs with a high proportion of Conflicts versus Non-Conflicts. We look at two levels for the analysis:

**Topics** The UNSC Conflicts corpus includes two topics, each with a different potential for Conflict. The first topic encompasses debates from 2014 about the Ukraine crisis (“Ukraine”), dealing with military conflict in which there are opposing conflicting parties. The second topic encompasses the Women, Peace and Security (“WPS”) agenda, dealing with norm debates. Generally, the Ukraine debates have a more confrontational nature, whereas the WPS debates are largely about reporting on the current situation. Therefore, we expect the Ukraine debates to be more argumentative than the WPS ones.

**Countries** We compare speeches given by the permanent members of the UNSC: China, France, Russian Federation, United Kingdom, and the United States of America. For the Ukraine agenda, we additionally include speeches given on behalf

of Ukraine.

We evaluate two methods to analyze the tree structures described in (Stede, 2016), who used it for the depth of argumentation on a small-scale analysis, and adapt the methods on a larger scale for Conflicts in diplomatic speech. More precisely, we describe two methods for characterizing the depth of an RST tree, both based on the so-called Nuclearity Mass (NM) distribution (Stede, 2016). The first Nuclearity Mass (NM1) value considers solely the number of central nodes, whereas the second Nuclearity Mass (NM2) also takes into account the distance of each node from the root. Central nuclei (CNs) are those EDUs that have zero or one satellite relations on the path from the leaf EDU node to the root of the tree.<sup>5</sup>

- (1) NM1 describes the proportion of CNs to all leaf nodes. For example, the set of leaf nodes in Fig. 1 consists of four EDUs with two CNs. The NM1 value for this tree is therefore 0.5 (2/4).
- (2) NM2 additionally includes the length of the path from the leaf node up to the root ( $l_i$ ). NM2 is the sum of  $l_i$  of the CNs, divided by the sum of all  $l_i$ . In the example, the root node of the subtree comprises EDUs 1-3. The  $l_i$  value for CNs is 13 (4+5+4); the  $l_i$  value for the full subtree is 16 (4+5+4+3). Given the multinuclear relations in this tree (EDUs 1-3), the NM2 value is 0.81 (13/16).

## 5 Results and Discussion

### 5.1 Relation Distribution

In this section, we discuss the overlap of Conflict types and the frequency of RST relations when only considering leaf nodes (Fig. 2 and 3) and inside a paragraph (Fig. 4). Note that in Fig. 2 we did not include relations that indicate mere textual organization (such as SAME-UNIT) or that are too infrequent (less than 10 occurrences both for leaf nodes and paragraphs). We merged the causal relations REASON-N and REASON (to REASON) because they only differ in how they weight two EDUs, i.e. whether the cause is more important than the reason or the other way around. Similarly,

<sup>5</sup>Following Stede (2016), we allow one satellite relation for CN, since we often encounter pairs of EDUs where the satellite elaborates the nucleus but still is strongly connected to the content of the nucleus (i.e., not digressing).

we merge EVALUATION-N and EVALUATION-S (to EVALUATION).

**Attribution:** Looking at ATTRIBUTION relations in Fig. 3, we notice a high proportion of *Challenging* (18.29%) and *Correcting* (6.29%) Conflicts. The high frequency of this relation is to some degree expected since *Challenges* are questioning the truthfulness of statements by another party and therefore are also reporting on what someone has (allegedly) said. *Corrections* are correcting an allegedly false statement, potentially citing a source of information (recall that ATTRIBUTIONS mark reported speech).

**Pragmatic Relations:** In section 4.1, we speculated that diplomats use more pragmatic relations for Conflicts than *Non-Conflicts* because these discourse relations describe the argumentation of the speaker, like justifying a thesis that the author has proposed (EVIDENCE, REASON), or evaluating a state of affairs from the author’s perspective (EVALUATION). In fact, EVALUATION is slightly more often used in *Direct\_NegEval* (2.06%) than for *Non-Conflicts* (1.52%), and EVIDENCE appears more often in *Indirect\_NegEval* (1.13%) than in *Non-Conflicts* (0.93%) (Fig. 3). Nevertheless, Conflicts in general are less often annotated with EVALUATION or other pragmatic relations than *Non-Conflicts* (3.6% pragmatic relations in Conflicts, 5.52% in *Non-Conflicts*) (Fig. 2).

When including the upper levels of the tree (Fig. 4), we see that paragraphs with *NegEval* Conflicts have only slightly more occurrences of RST relations expressing a justification with EVIDENCE and EVALUATION than paragraphs without Conflicts. Nevertheless, REASONS are found more often for *Non-Conflicts* than for Conflicts.

**Contrastive Relations:** Contrastive relations are generally more frequently used in Conflicts than in *Non-Conflicts* (Fig. 2) (4.02% versus 2.5%). Looking at the Conflict types in more detail (Fig. 3), we see that especially *Challenge* and *Correction* have a high proportion of ANTITHESIS and CONTRAST relations, which focus on the difference (CONTRAST) or incompatibility (ANTITHESIS) of two statements, and therefore the co-occurrence is to be expected. For *Direct\_NegEval* we see a peak for CONCESSION, which compares two incompatible states of affairs while regarding the content of one (the nucleus) more important than the other.

**Multinuclear and Semantic Relations:** We observe a high peak for CONJUNCTIONS for Conflicts and *Non-Conflicts*, which marks an enumeration and expresses otherwise little extra meaning. Semantic relations describing, for example, local or temporal CIRCUMSTANCES, causal relations expressing RESULT or PURPOSE appear proportionally more often in Conflicts, especially in cases marked as *Direct\_NegEval*.

Summarizing these results, we discuss possible first interpretations of the rhetorical strategies we can discern from the relation distribution analysis. For a more extensive discussion, we would need more qualitative analysis involving domain experts, to be able to generalize what the relation distribution could implicate for rhetorical strategies used in the UNSC.

Contrary to our hypothesis that Conflicts are more often justified than no Conflicts, in our corpus, pragmatic/justifying relations such as EVALUATION or REASON occur with similar frequencies in texts that do or do not contain Conflicts. On the other hand, we see some semantic relations, such as E-ELABORATION and PURPOSE, more often used with Conflicts than for *Non-Conflicts*. Looking into the speeches, we find that within a Conflict statement, often not only the actions of others are criticized, but especially the ascribed intention of the actors performing the action. These cases are annotated as PURPOSE, which could explain the generally high frequency of this relation.

Further, contrastive relations are more frequently used for Conflicts than for no Conflicts. We saw in a first qualitative study for the WPS debates that diplomats frequently place a positive statement in front of a direct critique that is then contrasted with the latter. Our annotators often used CONTRAST or CONCESSION to relate those two parts, which can indicate a rhetoric strategy to de-emphasize the verbalized critique. Again, these observations will need to be doublechecked with domain experts and tested on more data, but we include them here to exemplify what kind of analysis our corpus potentially enables.

## 5.2 Tree Structure Analysis

### 5.2.1 Nuclearity Mass and Tree Size

At first, we computed the NM1 and NM2 values for complete RST trees, but after some consideration, looked at subtrees within paragraphs instead. The

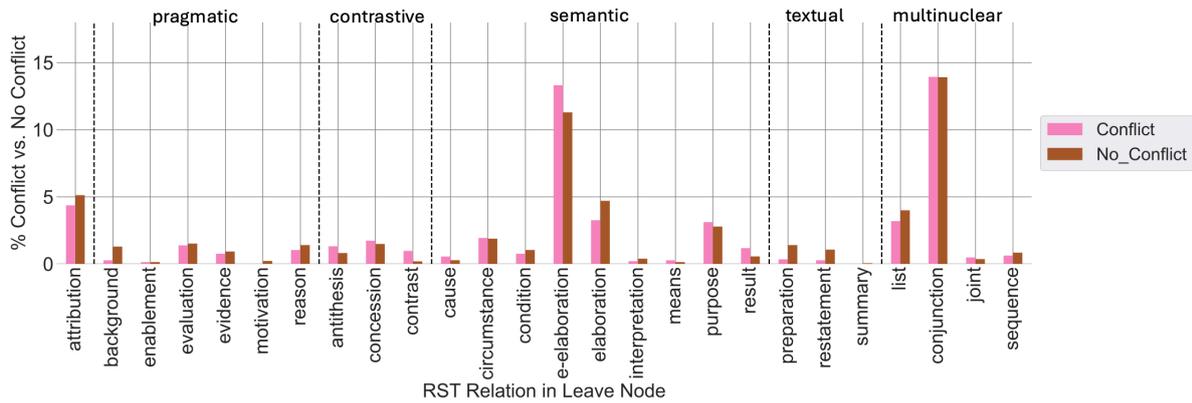


Figure 2: Normalized frequency of RST relations. The relations are grouped by their function.

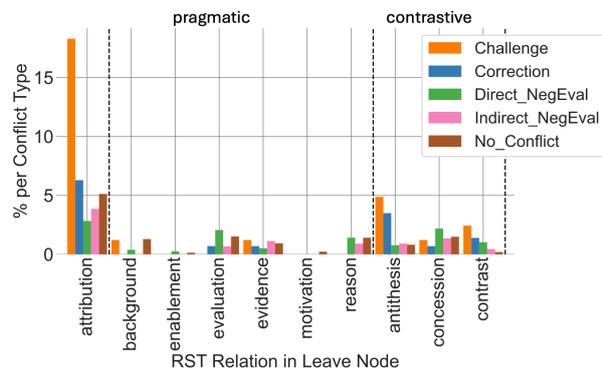


Figure 3: Normalized frequency of RST relations per Conflict type in leave nodes.

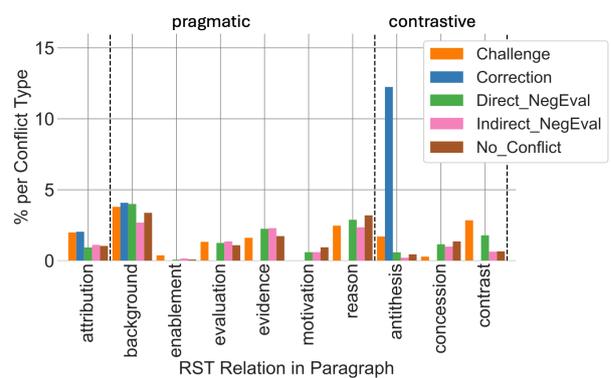


Figure 4: Normalized frequency of RST relations in paragraphs.

reason is that the NM value is sensitive to the size of the tree. In practice, annotators tend to establish a hierarchy between two EDUs, and choose multinuclear relations much less frequently (often for listings). Only multinuclear relations, which assign an equal weight to discourse units, lead to multiple CNs. As a consequence, we observe that the larger the tree, the smaller the NM value. Since the speeches in the UNSC Corpus have a large variety of tree length (see Section 3.3), this observation is especially important for our UNSC-RST.

To quantify this, the standard deviation for number of EDUs per speech/entire tree is 42.67, and for the number of CNs per speech it is six times lower (7.0). Looking at the same values for paragraphs, the standard deviation for EDUs per paragraph is 5.6, and for CNs it is 2.1, which is only 2.7 times lower. Since both NM measures are based on the ratio of leaf nodes to CNs, we decided to continue inspecting subtrees at the paragraph level in order to achieve better comparability of the trees.

## 5.2.2 Results and Discussion Nuclearity Mass

For a paragraph to be labeled as *Conflict*, we define that at least one third of the EDUs in the paragraph should be marked with one of the Conflict types. Otherwise, the paragraphs are marked as *Non-Conflict*. Note that for the analysis of discourse relation distribution in paragraphs (Section 5.1), only one EDU had to contain a Conflict type to be marked as Conflict, since Conflict types are too sparsely distributed to establish a higher threshold.

**Topics:** Broadly comparing the values for both measures NM1 and NM2, we see that they show similar results, but the NM2 values are generally smaller than NM1 values. Looking at Figure 5 on the left, showing the distribution NM values using both measurements, we see that the values for NM1 are higher than for NM2, but both measurements show that the NM distribution is slightly lower for Ukraine than for WPS. The fact that the WPS debates have more discourse units of equal importance is in line with our expectations, as the

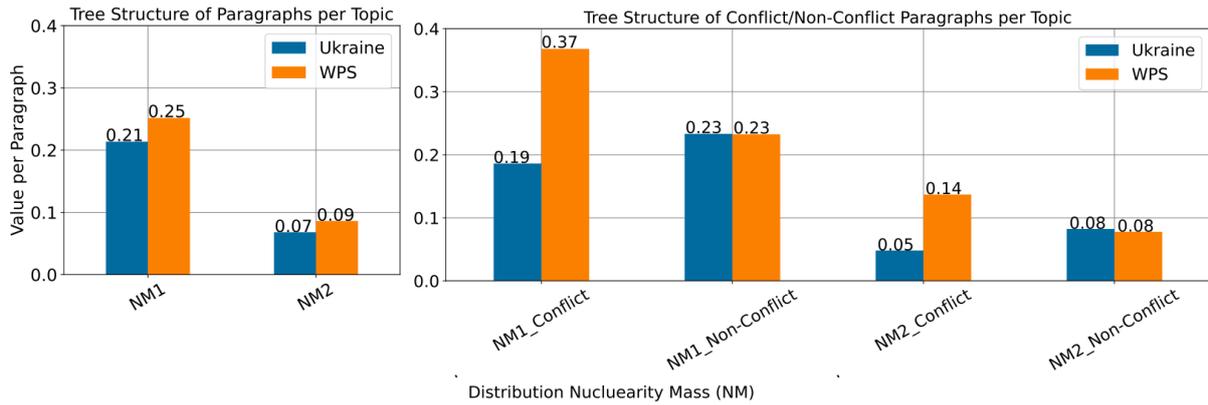


Figure 5: Mean Distribution of NM for Ukraine (194 Conflicts paragraphs, 271 Non-Conflicts) and WPS debates (16 Conflict paragraphs and 97 Non-Conflicts).

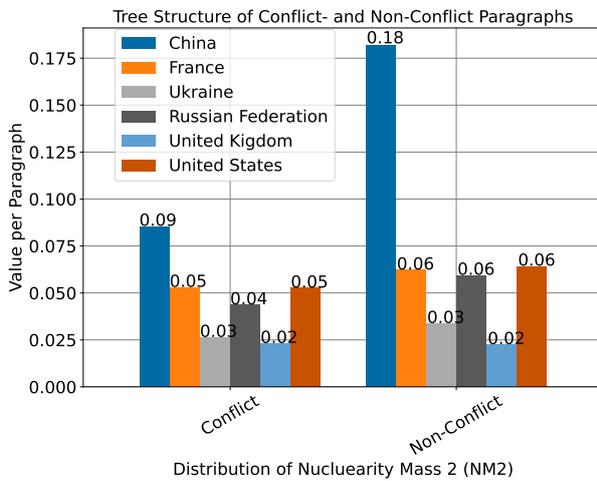


Figure 6: Mean Distribution of NM per Country, comparing Conflict versus Non-Conflict paragraphs.

WPS debates are often about summarizing what has been achieved in terms of gender and security issues and the situation in different countries.

**Topics and Conflicts:** Interestingly, comparing the topics with Conflicts versus *Non-Conflicts* paragraphs, we see that the difference between the topics is only in the Conflict, and that paragraphs with low proportion of Conflict types have similar NM Density values for both topics and both measures NM1 and NM2 (0.23 NM1 and 0.08 NM2 for both topics). One possible explanation would be that the Conflicts in WPS are rhetorically embedded and there is not one central message to which all the discourse units are leading (0.37 for NM1 and 0.14 NM2 for Conflict respectively). For Ukraine, on the other hand, it seems to be the opposite, with smaller values of 0.19 NM1 and 0.05 NM2 in Conflicts for Ukraine, and therefore having deeper tree

structures towards one EDU. Whether this means that the Conflicts in Ukraine are formulated with more intensity must be assessed by political scientists, but it would be a possible conclusion of the tree structures that we find.

**Countries and Conflicts:** Since, as mentioned above, both NM measures show similar values, just on a different scale, we will only look at the NM2 value for the statistics by country (Figure 6, the bar charts for both NM1 and NM2 are in Appendix D). The countries we compare are Ukraine (37 Conflict paragraphs, 36 Non-Conflicts), Russian Federation (29, 40), USA (32, 30), China (4, 18), United Kingdom and Northern Ireland (17, 27), and France (16, 28).

We see that the speeches given by China show the highest distribution of NM2 for both Conflicts and Non-Conflicts, which is insofar interesting as the diplomatic style of the Chinese government until the late 2010s is in fact known as using cooperative rhetoric and avoiding controversy (Yuan, 2023). We also notice a comparably large distance between the average Conflict (0.09) and Non-Conflict (0.18) values in the evaluated Chinese speeches in comparison to other speeches. This might point to a greater style change when expressing critique for the Chinese speeches than for other countries, using more non-argumentative style for *Non-Conflicts* and more argumentative for Conflicts. Nevertheless, we are looking only at four Conflict paragraphs for China, and we would need a larger corpus for greater validity.

All countries have lower NM values than China, with the lowest for both Conflicts and *Non-Conflicts* for Ukraine and the United Kingdom. This indicates an argumentative style that is fo-

cusing on one or a few statements and being more argumentative, in contrast to China. Also in Contrast to China, for Conflicts, the distribution of NM is almost similar to that of *Non-Conflicts*. This may indicate that the countries are not changing their rhetorical style when expressing Conflict as much as might be expected. Also for Conflict and *Non-Conflict*, the highest value for both is that of China, followed by the Russian Federation, France and the United States, and finally by Ukraine and the United Kingdom with the lowest NM values.

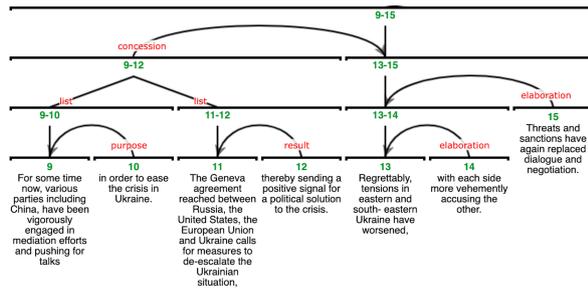


Figure 7: RST Paragraph with EDUs 13-15 being a Conflict (*Direct\_NegEval*) with NM1 0.64 and NM2 0.1 (S/PV.7138, China).

## 6 Conclusion

We present a new corpus with RST annotations on 82 speeches given in the UNSC, aligned with Conflict annotations from the UNSC Conflict Corpus. We report an average inter-annotator agreement of 0.44. By jointly analyzing RST trees and Conflict annotations, we demonstrate how rhetorical analysis can help characterizing the verbalized disagreements or critiques as being more argumentative or having a more narrative style. Comparing paragraphs that contain Conflicts with those that do not, we see that the former on average have rhetorical structures that focus on a central statement, rather than having several statements of same importance. Comparing speeches of six countries in the Council, we only see a larger difference between Conflicts and *Non-Conflicts* for the Chinese speeches. When comparing values between countries, they maintain their rhetorical style, with China always having the flattest, and the United Kingdom the most centralized rhetorical structure.

We see the work presented here as one of the first to use RST to analyze the rhetorical style of diplomats. More generally, we contribute to exploring ways of using RST trees in the analysis of a multi-layer corpus. In future work we want to expand not only the corpus with more topics and speeches, but

also the set of analysis methods. For example, we will have a closer look at patterns of rhetorical relations, and whether some relations co-occur more often than others, which might yield more insights on rhetorical strategies used by diplomats. Based on our presented tree structure analysis, it would also be interesting to compare trees that contain an EDU marked as Conflict as their central nucleus (and thus highlight the criticism) with trees where the Conflicts are hidden in higher parts of the trees (which might serve to weaken it). Our analyses show promising results, and open up a new direction of research, combining Conflict annotations (which are less time-consuming to obtain than RST trees) with manually evaluated and corrected RST parser output, in order to investigate on larger scale in potential future work.

## Limitations

For the analysis, we work with speeches translated into English, which may introduce a bias in the analysis of rhetorical structures, as the annotators pay close attention to linguistic subtleties in order to extract the discourse relationship between text segments. When comparing the rhetorical styles of diplomatic speeches, we need to be aware that the style of individual diplomats can also bring about a change in the strategies we see. In order to analyze this, and rhetorical style in general on a larger scale, we would need more data. The relatively small corpus size is due to the time-consuming process of annotating the RST trees, which took over 5 months. To accelerate the process, we plan to evaluate the performance of RST parsers trained on the latest version of the GUM corpus, which includes political speeches.

## Acknowledgments

We thank our annotators Dietmar Benndorf (main annotator), Louis Cho, Elisa Lübbers, Hugo Meinhof, and Costanza Rasi for their work and valuable feedback on the annotation guidelines. Furthermore, we would like to thank the reviewers for their engaged and helpful comments. The research was supported by the Deutsche Forschungsgemeinschaft (DFG), project (448421482) "Trajectories of Conflict: The Dynamics of Argumentation in the UN Security Council".

## References

- Noof Ibrahim Altmami and Mohamed El Bachir Menai. 2020. [CAST: A Cross-Article Structure Theory for Multi-Article Summarization](#). *IEEE Access*, 8:100194–100211.
- Lynn Carlson and Daniel Marcu. 2001. [Discourse Tagging Reference Manual](#). Technical Report ISI-TR-545, University of Southern California Information Sciences Institute.
- Lynn Carlson, Daniel Marcu, and Mary Ellen Okurovsky. 2001. [Building a Discourse-Tagged Corpus in the Framework of Rhetorical Structure Theory](#). In *Proceedings of the Second SIGdial Workshop on Discourse and Dialogue*.
- Lynn Carlson, Daniel Marcu, and Mary Ellen Okurovski. 2002. [RSTDiscourse Treebank LDC2002T07](#). Linguistic Data Consortium. Web Download.
- Freya Hewett. 2023. [APA-RST: A Text Simplification Corpus with RST Annotations](#). In *Proceedings of the 4th Workshop on Computational Approaches to Discourse (CODI 2023)*, pages 173–179, Toronto, Canada. Association for Computational Linguistics.
- Freya Hewett, Roshan Prakash Rane, Nina Harlacher, and Manfred Stede. 2019. [The utility of discourse parsing features for predicting argumentation structure](#). In *Proceedings of the 6th Workshop on Argument Mining*, pages 98–103, Florence, Italy. Association for Computational Linguistics.
- Mikel Iruskietia, Iria da Cunha, and Maite Taboada. 2015. [A qualitative comparison method for rhetorical structures: identifying different discourse structures in multilingual corpora](#). *Language Resources and Evaluation*, 49(2):263–309.
- Naoki Kobayashi, Tsutomu Hirao, Hidetaka Kamigaito, Manabu Okumura, and Masaaki Nagata. 2021. [Improving Neural RST Parsing Model with Silver Agreement Subtrees](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1600–1612, Online. Association for Computational Linguistics.
- Mathias Kraus and Stefan Feuerriegel. 2019. [Sentiment analysis based on rhetorical structure theory: Learning deep neural networks from discourse trees](#). *Expert Systems with Applications*, 118:65–79.
- Yang Janet Liu and Amir Zeldes. 2023. [Why Can't Discourse Parsing Generalize? A Thorough Investigation of the Impact of Data Diversity](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3112–3130, Dubrovnik, Croatia. Association for Computational Linguistics.
- William C. Mann, Christian M.I.M. Matthiessen, and Sandra A. Thompson. 1992. [Rhetorical structure theory and text analysis](#). In William C. Mann and Sandra A. Thompson, editors, *Pragmatics & Beyond New Series*, volume 16, page 39. John Benjamins Publishing Company.
- William C. Mann and Sandra A. Thompson. 1988. [Rhetorical Structure Theory: Toward a functional theory of text organization](#). *Text - Interdisciplinary Journal for the Study of Discourse*, 8(3):243–281. Publisher: De Gruyter Mouton Section: Text & Talk.
- Thanh-Tung Nguyen, Xuan-Phi Nguyen, Shafiq Joty, and Xiaoli Li. 2021. [RST Parsing from Scratch](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1613–1625, Online. Association for Computational Linguistics.
- Olu Popoola. 2017. [Using Rhetorical Structure Theory for Detection of Fake Online Reviews](#). In *Proceedings of the 6th Workshop on Recent Advances in RST and Related Formalisms*, pages 58–63. Association for Computational Linguistics.
- Victoria L. Rubin and Tatiana Vashchilko. 2012. [Identification of Truth and Deception in Text: Application of Vector Space Model to Rhetorical Structure Theory](#). In *Proceedings of the Workshop on Computational Approaches to Deception Detection*, pages 97–106, Avignon, France. Association for Computational Linguistics.
- Mirco Schönfeld, Steffen Eckhard, Ronny Patz, and Hilde van Meegdenburg. 2019. [The UN Security Council debates 1995-2017](#). *CoRR*, abs/1906.10969.
- Sophia Skoufaki. 2020. [Rhetorical Structure Theory and coherence break identification](#). *Text & Talk*, 40(1):99–124. Publisher: De Gruyter Mouton.
- Manfred Stede. 2016. [Towards assessing depth of argumentation](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3308–3317, Osaka, Japan. The COLING 2016 Organizing Committee.
- Manfred Stede, Maite Taboada, and Debopam Das. 2017. [Annotation Guidelines for Rhetorical Structure](#). *Unpublished manuscript*.
- Kun Sun, Rong Wang, and Wenxin Xiong. 2021. [Investigating genre distinctions through discourse distance and discourse network](#). *Corpus Linguistics and Linguistic Theory*, 17(3):599–624. Publisher: De Gruyter Mouton.
- Shujun Wan, Tino Kutschbach, Anke Lüdeling, and Manfred Stede. 2019. [RST-Tace A tool for automatic comparison and evaluation of RST trees](#). In *Proceedings of the Workshop on Discourse Relation Parsing and Treebanking 2019*, pages 88–96, Minneapolis, MN. Association for Computational Linguistics.

- Sandra Williams and Richard Power. 2008. [Deriving Rhetorical Complexity Data from the RST-DT Corpus](#). In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco. European Language Resources Association (ELRA).
- Shaoyu Yuan. 2023. [Tracing China's diplomatic transition to wolf warrior diplomacy and its implications](#). *Humanities and Social Sciences Communications*, 10(1):837.
- Karolina Zaczynska, Peter Bourgonje, and Manfred Stede. 2024. [How Diplomats Dispute: The UN Security Council Conflict Corpus](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 8173–8183, Torino, Italia. ELRA and ICCL.
- Amir Zeldes. 2016. [rstWeb - A Browser-based Annotation Interface for Rhetorical Structure Theory and Discourse Relations](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, pages 1–5, San Diego, California. Association for Computational Linguistics.
- Amir Zeldes. 2017. [The GUM Corpus: Creating Multilayer Resources in the Classroom](#). *Language Resources and Evaluation*, 51(3):581–612.

## A Appendix: Inter-Annotator Agreement

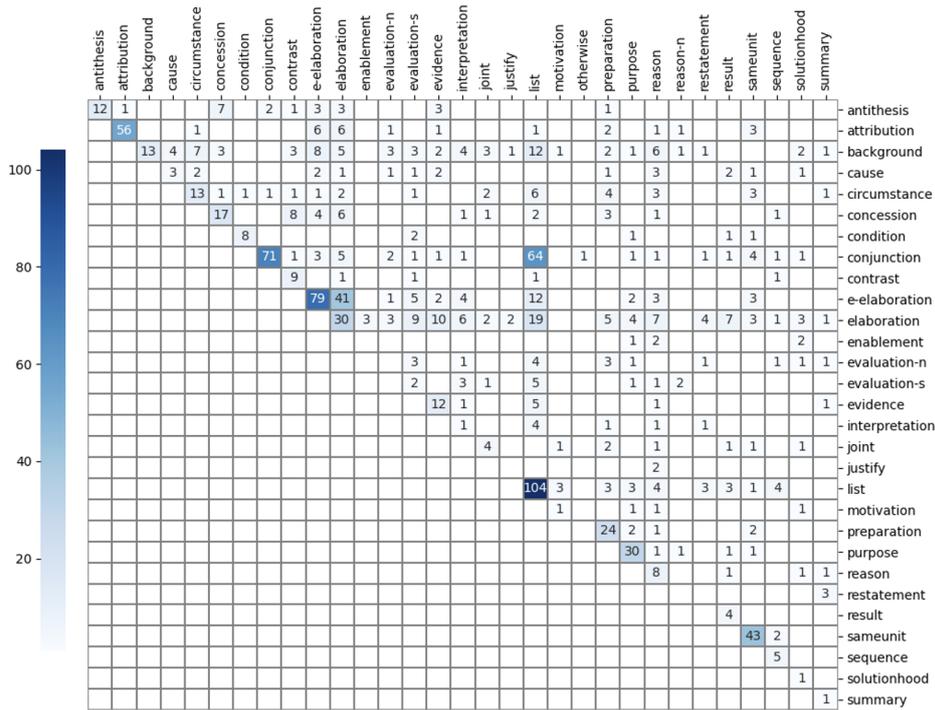


Figure 8: Confusion matrix with RST relations for two parallel annotations per RST-tree.

## B Appendix: Example RST Trees with different tree shapes

We included two example trees from the UNSC Conflicts corpus, where the first one has a clearly-identifiable central nucleus ("We trust that Russia will take notice of its isolation."). The second tree shows a tree with a higher distribution of NM with several EDUs having a multinuclear relations toward the top of the tree, and several points that are perceived as being equally important to the author of the text. For the upper tree in Figure 9, the average values per paragraph are 0.27 NM1 and 0.046 NM2; for the lower tree they are 0.64 and 0.15.

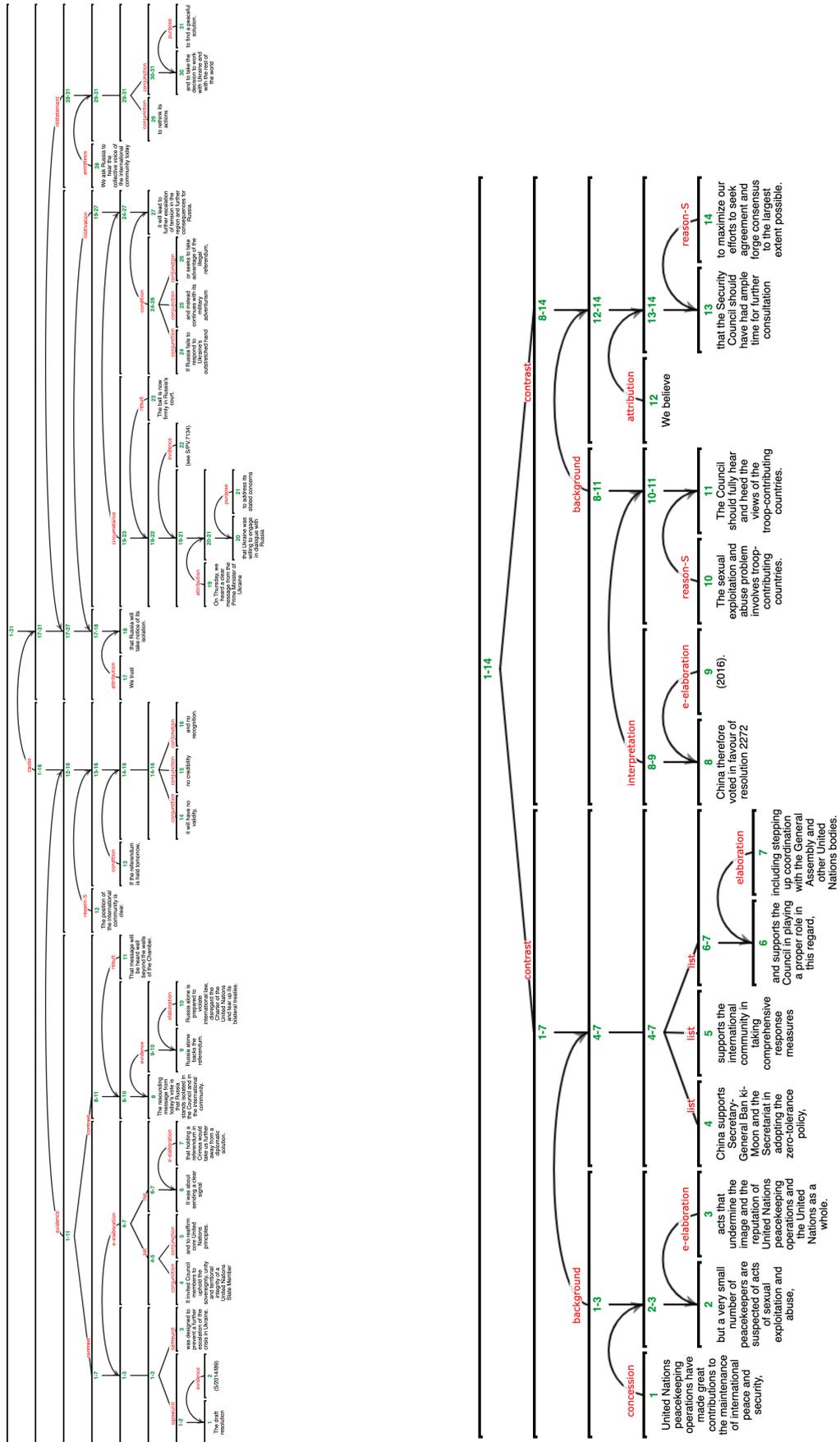


Figure 9: RST Example trees (S/PV.7138\_spch006 by United Kingdom, and S/PV.7643\_spch008 by China below) with different nuclearity mass distribution.

### C Statistics for RST Relation Distribution Bar Charts

	Challenge	Correction	Direct NegEval	Indirect NegEval	Non-Conflict
paragraph #EDUs	1,054	49	12,864	3,314	12,299
leaf nodes #EDUs	82	143	776	441	3,550

Table 1: Number of EDUs per Conflict Type

### D Nuclearity Mass per Country for both Measures NM1 and NM2

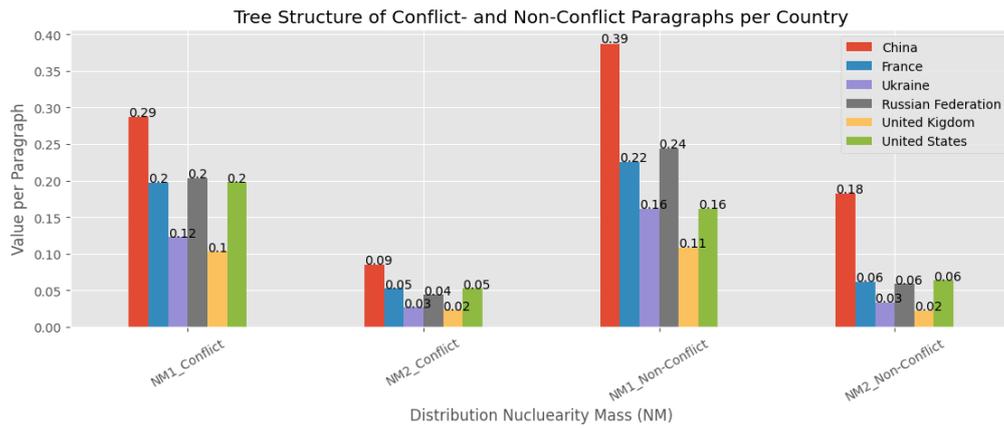


Figure 10: Mean Distribution of NM per Country, comparing Conflict versus Non-Conflict paragraphs.

# Elaborative Simplification for German-language Texts

Freya Hewett<sup>1,2</sup>, Hadi Asghari<sup>1,3</sup>, Manfred Stede<sup>2</sup>

<sup>1</sup>AI & Society Lab, Humboldt Institute for Internet and Society, Berlin, Germany  
firstname.lastname@hiig.de

<sup>2</sup>Applied Computational Linguistics, University of Potsdam, Germany  
lastname@uni-potsdam.de

<sup>3</sup>Faculty of Electrical Engineering and Computer Science, TU Berlin, Germany

## Abstract

There are many strategies used to simplify texts. In this paper, we focus specifically on the act of inserting information or *elaborative simplification*. Adding information is done for various reasons, such as providing definitions for concepts, making relations between concepts more explicit, and providing background information that is a prerequisite for the main content. As all of these reasons have the main goal of ensuring coherence, we first conduct a corpus analysis of simplified German-language texts that have been annotated with Rhetorical Structure Theory (RST). We focus specifically on how additional information is incorporated into the RST annotation for a text. We then transfer these insights to automatic simplification using Large Language Models (LLMs), as elaborative simplification is a nuanced task which LLMs still seem to struggle with.

## 1 Introduction

There are many strategies used to simplify texts. Sentences can be shortened, split or paraphrased, complex words replaced with synonyms, and information can be reordered, dropped or inserted (Amancio and Specia, 2014; Alva-Manchego et al., 2019). In this paper, we focus specifically on the act of inserting information.

Inserting information is done for various reasons: providing definitions for concepts, making relations between concepts more explicit, and providing background information that is a prerequisite for the main content. These all should contribute to decreasing complexity and therefore ideally ensuring coherence; the semantic or pragmatic relationships that link units in a discourse to other units (Das and Taboada, 2018). Readers need to recognise these relationships in order to make sense of the text, so a more coherent text should increase comprehension and also allow readers to recognise the communicative function of the text (cf. Nussbaumer, 1993).

In our study, we focus on German-language texts and aim to transfer insights from a detailed corpus analysis to automatic simplification models, to improve their ability in inserting information and therefore their overall ability at simplification. We use a corpus of parallel newspaper articles that have been annotated with Rhetorical Structure Theory (RST). RST annotations provide information about how segments in a text are related to each other within semantic or pragmatic relations such as *cause*, *background*, or *contrast* (Mann and Thompson, 1988). Our corpus analysis examines how inserted information in simplified texts can affect the coherence, and also what purpose the additional information has.

In order to utilise these discourse structure annotations for the task at hand, we first add a new layer of annotations by labelling the transformations that are applied to the original sentences to create the simplified sentences. One of these labels is ‘Insert complementary information’ which we focus on in more detail. We examine the role that this inserted information plays in the overall RST annotation.

We then transfer these insights to automatic simplification using Large Language Models (LLMs), exploring the use of different prompts.

In summary, our contributions are: we extend the APA-RST corpus (Hewett, 2023) to include transformation labels. We show results of an extensive corpus analysis, showing how new information is inserted in text simplification, and how this affects the coherence. We then explore models for document-level text simplification for German using the insights from our corpus analysis, with results comparable with the state of the art.

In Section 2 we present an overview on work that has looked at the insertion of new information in simplified text. In Section 3 we present our annotations of alignment labels and fine-grained inserted information categories, before presenting our RST analysis. Section 4 gives details on our

models and experiments with them, and we reflect on our results and possible avenues for future work in Section 5. We publish our annotations and models at <https://github.com/fhewett/GermanElabSimplification>.

## 2 Related work

Srikanth and Li (2021) introduce the term ‘elaborative simplification’ to describe content addition in text simplification. The elaborative content added consists of ‘definitions, explanations or clarifications to improve readability’ with effective elaborations providing background information ‘in a contextual manner’. They focus on this contextual aspect, annotating 1,300 instances of elaborative simplification in the Newsela corpus (Xu et al., 2015), categorising them according to the level of contextual specificity. They experiment with GPT-2, fine-tuning it on the simplest texts in Newsela and their annotated elaborations. Their best-performing model has the four sentences preceding an elaboration in a simplified text as input, and generates an elaboration as output, with the level of context specificity as determined by the gold annotation. Wu et al. (2023) use these annotated instances and add Questions Under Discussion (QUD), to show which questions elaborations answer. They find that the most common purpose of the elaboration is to explain a concept, followed by elaborations explicitly describing the cause of consequence of an event. They use GPT-3 for zero-shot elaboration generation, experimenting with including an automatically generated or human annotated QUD in the prompt or not. The results show that manually written QUDs produce the best elaborations. These studies build on ideas proposed by Alva-Manchego et al. (2020), who list explanation generation as an area of future work (albeit in the context of sentence-level simplification), stating that it involves elaborating ‘on the concept in a natural way that keeps the text grammatical, is meaning preserving, and is simple’. Additionally, the well-established evaluation metric for automatic simplification, SARI, rewards ‘addition operations’ (Xu et al., 2016).

Another related area of text simplification is conceptual complexity, defined as accounting for ‘the background knowledge necessary to understand mentioned concepts as well as the implicit connections that the reader has to access between the mentioned concepts in order to fully understand a

text’ (Hulpuş et al., 2019).

Our work is also related to the field of factuality (evaluation) of language model outputs: Devaraj et al. (2022) create a taxonomy of factual errors in automatic simplification, including ‘Information Insertion’ which is described as inserting ‘irrelevant or erroneous content’. They differentiate between these insertion errors and useful insertions, such as ‘defin[ing] jargon or provid[ing] explanatory content’. In the field of automatic summarisation, Maynez et al. (2020) differentiate between intrinsic and extrinsic hallucinations, where the latter refers to ‘adding information not directly inferable from the input document’. They find that ‘over 90% of extrinsic hallucinations were erroneous’ i.e. are ‘neither faithful nor factual’. Maynez et al. (2020) also find factual hallucinations to be ‘acceptable if they lead to better summaries that are factual with respect to the document and the associated background knowledge’. This last point is particularly relevant to the task of simplification.

In various guidelines on *Leichte Sprache* (LS) – a highly simplified rule-based version of German – inserting factual information is allowed and also even desirable, in order to increase the level of comprehension on the one hand, and to allow readers to potentially learn new information on the other hand (Maaß, 2015). The guidelines state that translators of LS are allowed to provide explanations, additional remarks, and (concrete) examples, in order to make abstract concepts or difficult words more comprehensible. Maaß (2015, p.130) does however state that translators, after adding these definitions, explanations and examples, should make sure that the text still has an argumentative flow. Bredel (2016) state that additional explanations in texts in LS can hinder the flow of the text and potentially also cause problems on the text level. These aspects are the specific focus of the current study, i.e. what happens to the structure and coherence of the text overall when elaborative simplification is used.

Other corpus studies which focus on the transformation operations between non-simplified and simplified text often define an operation for inserting information. This category encompasses sub-categories such as inserting *eliciting information*, *complementary external information*, *spurious information*, *pre-requisite information*, *concrete examples of abstract concepts or phenomena* (Amancio and Specia, 2014; Alva-Manchego et al., 2019; Sun et al., 2021; Laban et al., 2023). This

category has also been used in German-language corpus studies: Stodden et al. (2023) manually align parallel texts with a category for additional information and Jablotschkin et al. (2024) find that phrases such as ‘for example’ or ‘that means’ feature heavily in simplified texts and are used for *explaining difficult words, making abstract concepts more concrete and connecting the sentences of a text explicitly*.

### 3 Corpus analysis

The main corpus we work with is the APA-RST. The corpus consists of German-language newspaper articles, which are classified as being at B1 and A2 level, according to the Common European Framework of Reference for Languages (CEFR), which is a scale from A1 (beginner language learner) to C2 (native speaker). There are 75 parallel articles in the corpus, with 25 at each level (original<sup>1</sup>, B1 and A2), covering various topics such as politics, culture and sport. The articles have been annotated with RST and manually aligned at sentence level; further information can be found in the original publication (Hewett, 2023). Due to the relatively small sample size, we extend our analysis to label 200 instances of the ‘APA’ subcorpus of DEplain (Stodden et al., 2023) which features a larger number of newspaper articles from the same publisher as APA-RST. This subcorpus has been aligned at the sentence level, between the versions B1 and A2.

#### 3.1 Adding transformation labels

Two annotators added transformation labels to the sentence alignments in the APA-RST, i.e. a label to describe how the original content was transformed for the simplification. We determined our labels by first selecting a subset of the most relevant labels from previous work (cf. Section 2). We then annotated a few texts and refined the definitions and added or removed labels. Our final label set consisted of **Paraphrase** (the content is the same, but the wording and/or structure are different), **Simple split** (original sentence has been split into two or more sentences, the structure and vocabulary are similar), **Complex split** (a split combined with a paraphrase), **Join** (content from two or more original sentences is combined in one simplified sentence), **Drop extra information** (sentences are

<sup>1</sup>These articles do not have a language level but are assumed to be at C1/C2 level.

Label	OR⇒B1	B1⇒A2
Paraphrase	15%	46%
Simple split	1%	9%
Complex split	23%	13%
Join	4%	3%
Drop extra info	34%	13%
Insert complementary info	19%	9%
Implicit	2%	4%
Identical	2%	3%

Table 1: Distribution of transformation labels. Note that for OR⇒B1 78% of the sentences are dropped, for B1⇒A2 14% are dropped. The distribution of the labels amongst the remaining 22% and 86% are shown here.

fairly similar, but some content has been dropped for the simplification), **Insert complementary information** (the simplified version contains content that is not explicit in the original), **Implicit** (content is included implicitly in original), and **Identical** (sentences are identical). Often the majority of sentences could be described as being paraphrases, and so the label **Paraphrase** was only to be used when no other category was suitable. The inter-annotator agreement, calculated using Cohen’s kappa, is .62 for the labels from original to B1 and .72 for B1 to A2, which compares to related work (.62 for five transformation categories in Laban et al. 2023).

The distribution of our labels can be seen in Table 1. For the rest of the study, we focus on the labels **Insert complementary information** and **Implicit**. Although these do not constitute the largest categories of transformations in a simplification, we choose to focus on them as choosing the right complementary information to insert requires high-level reasoning and is linked to the ‘hallucinatory’ nature of texts produced by LLMs.

#### 3.2 Categories

We built a small typology of categories of inserted information, based on the transformation labels and their descriptions that were outlined in Section 2. Our categories and their descriptions can be seen in Table 2. We label all sentences that have the alignment label **Insert complementary information** or **Implicit**. In addition to this, we focus on the DEplain alignments labelled with **Additional**. We exclude any of the DEplain sentences which do not match with our alignment transformation guidelines, i.e. if a sentence is labelled as **Additional**, but would be labelled as a different category according to our guidelines, we exclude it. Note that

Name of category	Description	Example	%
Example	Provide an example to make a concept clearer.	For example, coloured pencils from the same company cost more in some shops than in others.	4.2%
Background	Provide information that is a prerequisite for understanding the rest of the text.	In the Spanish region of Catalonia, many people voted in favour of independence from Spain in 2017.	33.1%
Relation	Make a relation more clear/explicit between concepts.	The new virus variant emerged for the first time in South Africa. (Next sentence: All people who have returned from certain South African countries in the last few days should now take a PCR test).	32.2%
Definition	Provide a definition or summary of a concept.	Pub is the English word for a <i>Lokal</i> .	15.1%
Additional	Provide information that is new but is not necessarily required for understanding the main points.	Marcel Sabitzer won the vote last year.	15.5%

Table 2: The names, descriptions and distribution of our fine-grained labels for inserted information.

the APA texts often include glossaries in the simplifications, providing definitions on concepts and words. We do not include these in our analysis, as we focus on coherence within the main text.

The largest categories of inserted information are **Background** and **Relation**, which are both specific to the context of the text that is being simplified. **Examples** are the rarest kind of inserted information; we note however that this is not to say that examples are rare in the texts overall, it is often the case that the examples are present in the original texts and therefore do not constitute *inserted* examples. We note that additional information that seems to have no purpose other than providing more (non-prerequisite) knowledge also occurs (**Additional**), but that generally there is a balance between succinctness and level of simplification.

### 3.3 RST analysis

We look at the RST trees and the overall structure of the texts in APA-RST, and consider the individual properties of the inserted information, such as the position, the RST relation, the nuclearity status, and how this relates to the fine-grained category (i.e. the type of inserted information, as outlined in the previous section). Adding definitions and prerequisite information is done to contribute to making a text easier to understand, i.e. by making relations between concepts and facts more explicit and reducing the background knowledge needed to understand a text. However, adding this new information changes the structure and flow of texts, and also changes the way adjacent statements relate to one another. Analysing the RST annotations could help shed light on how the structure of texts

change and how new information is used to ‘facilitate connections between content in the original text’ (Srikanth and Li, 2021).

**Relation.** Overall we find that when the function of the inserted information is annotated as ‘relation’, i.e. making the link between two concepts more explicit, the inserted information is part of an RST relation broadly belonging to the causal category, such as *cause*, *motivation* or *evidence*. This can for example be seen in Figure 1a, where segments 6 and 7 are inserted information which have been annotated as ‘relation’; they serve as the consequence of the causal relation, which is left more implicit in the original and the B1 text. This inserted information also makes the contrast relation, which connects a large amount of segments in the text, even more apparent, as it evens out the amount of sentences on each side of the contrast relation (2 vs. 2 in the A2 text, 3 vs. 1 in the B1 text).

**Background.** The inserted ‘background’ information is often at the beginning of the text; either directly at the beginning, as in:

**After a week of lockdown in Austria, the government started discussing the Corona situation on Monday.** (*N elaboration*, 2-29-11-21-b1)<sup>2</sup>

Or after an initial sentence that has been paraphrased from the original article. In some cases this summarising background sentence at the beginning of the original articles is suitable to start a simplification with, and in other cases it is necessary to add

<sup>2</sup>The whole texts can be viewed here: <https://github.com/fhewett/apa-rst>. Sentences in **bold** represent inserted information.

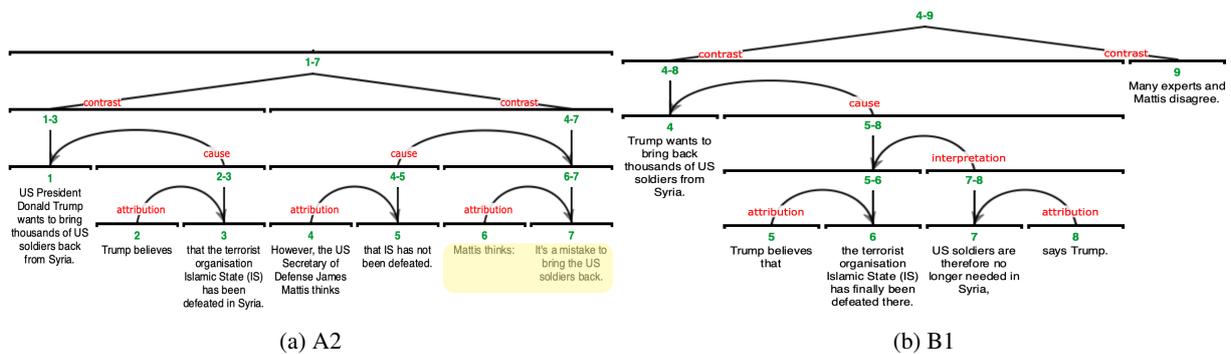


Figure 1: Extracts of the RST annotations for the text 1-21-2-18. The new information is highlighted in yellow. The trees were created using rstWeb (Zeldes, 2016).

information before this first sentence. This background information is often elaborated upon in the article and therefore often has the relation *elaboration* or *background*. In the A2 versions, the content added in the B1 versions is expanded upon with more additional content, to make relations even clearer or to reduce the amount of presupposed background knowledge:

**Because the hailstones were so large, they caused a lot of damage.** (*N evidence*, 3-21-2-18-b1)

This is expanded in the A2 text with two additional sentences preceding it:

**When it hails, icy stones fall from the sky. Normally the hailstones are as small as peas.** (*S background*, *S concession*, 3-21-2-18-a2)

This indicates that when creating simplified texts at different levels, the same content that has been added for a more complex level can be expanded upon for a less complex level (as opposed to adding new content which covers a different topic than the previously added content).

**Definitions.** When definitions are added to the text directly (as opposed to glossary entries, which are displayed outside of the text), they are often used for conversions, or for translations:

**That [23%] is almost a quarter more expensive than last year.** (*S elaboration*, 3-29-11-21-b1)

Inserting new information does create more "distance" between some entities:

In New York, the city in the US, a painting has been sold at auction for around

45 million dollars. **That is around 40 million euros.** The picture originates from the Italian painter Sandro Botticelli. (*S e-elaboration*, 5-freitag-28-1-22-a2)

In this text, the information about the equivalent euros amount is added, and the third sentence then goes on to talk about the painting again (i.e. the entity introduced in the first sentence). It is not clear if this added distance makes comprehension more difficult. It seems that, at least in the articles published by the APA, longer definitions are not favoured in the main text, instead being given in a separate glossary. On the one hand, this ensures that the added definition does not cause too much distance between information on the same entity, on the other hand, it requires the reader to move between the main text and the additional glossary, interrupting a normal reading flow.

We note that there are no clear trends regarding the local (the importance of a segment within a segment-level relation) or global nuclearity (the importance of a segment within the overall tree) of the inserted information, indicating that it has many roles within a text.

Inserted **Examples** do not occur in the APA-RST, and as **Additional** inserted information may in fact be undesirable in a simplification (the information is unnecessary and increases the length of the text), we do not go into detail on this category.

### 3.4 Summary of corpus analysis

Our transformation labels show that the insertion of information does occur at both simplification levels, and whilst not as common as dropping information or splitting sentences, it still is frequent, particularly in simplification of original texts to B1.

Our fine-grained categories show that **Background** and **Relation** are the most common types

Prompt ID	Prompt text
Basic	Can you please summarise and simplify the following text to a B1/A2 level in German? Write a maximum of $N$ sentences.
Background	Basic + add 1-2(B1)/2-3(A2) sentences at the beginning to give the user an overview of the topic. The text should have a clear introduction and information should be presented in a logical order.
Relation	Basic + add more contextual information to make the text easier to understand.

Table 3: The different prompts we use in our experiments.  $N$  is changed dynamically to reflect the amount of sentences in the reference simplification, and B1 or A2 used depending on the test set.

of inserted information, indicating that effective text simplification also involves conceptual simplification, i.e. decreasing the amount of background knowledge needed by the reader and therefore making relations more explicit. These transformations are more contextual than simply providing a definition.

Our RST analysis shows that background information is often at the beginning of a text, and often has the relation *elaboration* or *background*. Definitions that are added to the text could create ‘distance’ between related concepts, i.e. they add information that only attaches to one segment in the annotation, which may be why definitions only occur fairly rarely. In other texts, summarising sentences are used at the beginning or end of a sub-tree, so before the topic is changed slightly. When comparing simplifications from B1 to A2, the inserted content expands on the content that has already been inserted for the B1 text. Inserted content which makes a relation more clear often has a causal relation, so is making a cause or a consequence more explicit.

#### 4 Automatic simplification models

We use Meta-Llama-3-8B-Instruct for our experiments as it is one of the most capable open-weight LLMs at the time of writing and performs well in benchmarks.<sup>3</sup> Additionally, LLMs that have been trained using strategies such as instruction-tuning and RLHF (as is the case for Llama-3) have been found to perform well in the task of automatic sentence simplification (Kew et al., 2023). We use Meta-Llama-3-8B-Instruct out-of-the-box, and also use this base model to fine-tune on B1 texts and A2 texts. We then explore using different prompts which are influenced by the findings from

<sup>3</sup>More information can be found on the model card on HuggingFace: [https://github.com/meta-llama/llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md)

our corpus analysis.

#### 4.1 Experimental setup

For fine-tuning, we use the same kind of texts found in the APA-RST, but in an extended version, and with no annotations.<sup>4</sup> The original texts are aligned with versions at B1 and A2. We use 2000 articles for training, and 296 for testing. We remove headlines and any glossary entries. We use the Basic prompt in Table 3 for fine-tuning; we include the word ‘summarise’ in the prompt as the simplified texts in our corpus are approximately a third of the length of the original texts. We separately fine-tune a model with A2 texts (FT-A2) and with B1 texts (FT-B1). As we use the 25 texts from APA-RST for our prompting experiments for inference, these texts are neither in the training nor the test set. Information on hyperparameters can be found in the Appendix A.2.

At inference, in addition to a basic prompt, we try out two other prompts (per model) which target the aspects **Background** and **Relation**.

We focused on these two categories as they were found to be most prominent in our corpus analysis. We leave experiments with the other categories for future work, but note that examples which are inserted in the simplification (i.e. the category **Example**) were rare in our corpus analysis and that additional information (i.e. the category **Additional**) could be difficult to evaluate and is potentially also not desirable even in a gold simplification, as it increases the complexity of a text and introduces potentially unnecessary information.

The prompts can be seen in Table 3. We use the texts from APA-RST as part of the prompts, for in-context learning. We used the following template for the **Background** and **Relation** prompts<sup>5</sup>:

<sup>4</sup>A version of this dataset is also used by Rios et al. (2021) and Stodden et al. (2023).

<sup>5</sup>The exact format can be found in our repository: <https://github.com/fhewett/GermanElabSimplification>

Model	Prompt	Test set	SARI ↑	FRE ↑	M.P. ↑	S ↑	C ↑	F ↑	Avg. ↑
Baseline	Basic <sup>A2</sup>	A2	41.2	59.4	<b>.89</b>	.38	<b>.96</b>	<b>.84</b>	<b>.77</b>
FT-A2	Basic <sup>A2</sup>	A2	<b>44.0</b>	<b>70.6</b>	.49	<b>.82</b>	.56	.64	.63
Baseline	Basic <sup>B1</sup>	B1	42.3	56.8	<b>.85</b>	.4	<b>.9</b>	<b>.9</b>	<b>.76</b>
FT-B1	Basic <sup>B1</sup>	B1	<b>42.4</b>	<b>60.0</b>	.75	<b>.55</b>	.6	.75	.66

Table 4: Comparing Llama-3 out-of-the-box and fine-tuned. The test set consists of 296 articles. The mean FRE score for the reference simplifications is 63.2 for the B1 texts, 69.1 for the A2. FT stands for fine-tuned. The right hand side shows the results of the manual evaluation, done on the outputs from each model for 10 texts. M.P. stands for meaning preservation, S for simplification, C for coherence, F for factuality; the score represents the percentage of ‘yes’ answers.

Model	Prompt	Test set	SARI ↑	FRE ↑	M.P. ↑	S ↑	C ↑	F ↑	Avg. ↑
FT-A2	Basic <sup>A2</sup>	A2	44.0	70.6	.48	.8	.58	<b>.63</b>	.62
FT-A2	Background	A2	44.2	<b>70.8</b>	.51	.8	<b>.59</b>	.54	.61
FT-A2	Relation	A2	<b>44.5</b>	70.7	<b>.55</b>	<b>.95</b>	.57	.55	<b>.65</b>
FT-B1	Basic <sup>B1</sup>	B1	42.4	60.0	<b>.75</b>	.55	.6	<b>.75</b>	<b>.66</b>
FT-B1	Background	B1	42.6	<b>64.7</b>	.47	<b>.79</b>	<b>.63</b>	.32	.55
FT-B1	Relation	B1	<b>43.0</b>	64.0	.58	.68	.47	.68	.61

Table 5: Results for prompting experiments. The test set consists of 296 articles. The mean FRE score for the reference simplifications is 63.2 for the B1 texts, 69.1 for the A2. FT stands for fine-tuned. The right hand side shows the results of the manual evaluation, done on the outputs from each model for 10 texts. M.P. stands for meaning preservation, S for simplification, C for coherence, F for factuality; the score represents the percentage of ‘yes’ answers.

**system** You are a helpful assistant and help the user to understand texts.  
**user** {basic prompt} {original article}  
**assistant** {text without inserted information}  
**user** Thank you, that is good, but {additional prompt}  
**assistant** {text with inserted information}  
**user** Great, {additional prompt} {next original article}  
**assistant**

To determine the wording for the prompts in Table 3, we first prompt Llama-3, asking it to tell us which of two texts are easier to understand and why; one text is an A2/B1 text from APA-RST, the other text is the same but with the inserted information removed (and small changes made to keep the text coherent). An example of this can be found in Appendix A.3. The overall approach was inspired by the concept of prompt chaining, where splitting up a task into subtasks potentially allows for more controllability (cf. Wu et al. 2022); as we want to encourage elaborative simplification, our approach involves making this explicit in the prompt by differentiating between simplification with and without elaboration.

## 4.2 Results

Model	Prompt	Test set	Elab
Baseline	Basic <sup>A2</sup>	A2	.04
FT-A2	Basic <sup>A2</sup>	A2	.64
FT-A2	Background	A2	.65
FT-A2	Relation	A2	<b>.67</b>

Table 6: Manual evaluation of the models’ ability to perform elaborative simplification. The ‘Elab’ column refers to percentage of output texts in which the annotators said that elaborative simplification was performed.

We perform automatic evaluation using the metrics SARI (Xu et al., 2016) and Flesch Reading Ease (FRE). SARI is often used to evaluate simplification<sup>6</sup>, and FRE is used to evaluate readability. FRE takes the average sentence length and average amount of syllables per word into account. The mean FRE score for the reference simplifications in our dataset is 63.2 for the B1 texts and 69.1 for the A2 texts. SARI is based on n-grams added, kept, and deleted between a reference simplification and a generated simplification. We report these metrics to be comparable with related work, but note that they have been criticised for being insuffi-

<sup>6</sup>We are aware of D-SARI (Sun et al., 2021), a version of SARI specifically adapted for document-level simplification evaluation; however we choose to use SARI to be comparable with other related work on document-level simplification (Stodden et al., 2023; Laban et al., 2023).

Model/Prompt	Simplification	Baseline simplification (FT-A2/B1, basic prompt)
FT-A2 Background	<b>Salzburg is a city in Austria. It is located in the Alps. Salzburg is famous for its festival.</b> [...] That’s why Salzburg is the best city to visit in the coming year.	According to the travel book publisher "Lonely Planet", the city of Salzburg in Austria is the best city to visit in the coming year. [...]
FT-A2 Relation	<b>There are 14 mountains in Nepal that are very high. They are called eight-thousanders.</b> The highest mountain is Mount Everest. The Nepalese Nirmal Purja has climbed all 14 mountains. [...]	Nirmal Purja is a mountaineer from Nepal. He has set a record. He has climbed all 14 eight-thousanders in just 7 months. [...]
FT-B1 Background	<b>Ursula Stenzel is a politician from the FPÖ. She is a city councillor in Vienna.</b> [...]	Vienna City Councillor Ursula Stenzel (FPÖ) has not withdrawn after her appearance at a rally organised by the far-right <i>Identitären</i> . [...]
FT-B1 Relation	[...] <b>This is a problem because cars emit a lot of carbon dioxide. This is harmful for the environment.</b> The Austrian Transport Club (VCÖ) is therefore calling for more buses and trains. [...]	In Austria, car traffic has risen sharply since 2010. [...] The VCÖ is calling for a denser public transport network with more frequent train and bus connections.

Table 7: Examples of texts generated with different models and different prompts, compared to the basic prompt. The texts have been translated from German. The desired inserted information is in bold. We note that the FT-A2 **Relation** output contains a factual error, which is reflected in our manual evaluation.

cient measures of the quality of a simplification (cf. [Alva-Manchego et al. 2021](#)).

We extract 35 input texts and manually evaluate the outputs of our different models and prompts. We annotate the model outputs manually according to four criteria: meaning preservation, simplicity, coherence, factuality. Each criterion is a binary yes/no question. In addition to this, for a subset of 20 of these input texts we additionally annotate if the A2 models performed elaborative simplification. We only include the A2 models in this second evaluation as we use reference annotations to guide the evaluation and the majority of the instances in our corpus analysis were from A2 texts, due to the structure of DEplain. In total, three annotators evaluated 260 output texts. For 60 of these texts we have double annotations. The inter-annotator agreement for these texts across all criteria is .37 calculated using Cohen’s kappa or .8 using the F1 score.

**Llama-3 out-of-the-box vs. fine-tuned.** As can be seen in Table 4, our fine-tuned models only slightly outperform Llama-3 out-of-the-box (referred to as baseline) for the B1 texts, but for A2 texts the improvement is more pronounced, particularly in terms of readability, as reflected by the FRE score. Our results are higher than ([Rios et al., 2021](#)), who report a highest SARI score of 32.9 using APA data, and compare to ([Stodden et al., 2023](#)), who report a highest SARI score of 44.6 when simplifying from B1 to A2 (not from standard to A2/B1, as we do in this study). We note that this improvement is rather due to the improvements that LLMs

have made, rather than our method. The manual evaluation shows that the baseline model produces coherent, factual texts that cover the main points of the article, but are not necessarily written in a simpler way. As our main goal is simplification, we use our fine-tuned models for our prompting experiments.

**Prompting experiments.** As can be seen in Table 5, our prompts do result in slightly higher SARI and FRE scores. However, according to our manual evaluation, the prompts lead to a drop in factuality, meaning preservation and coherence. Overall, our prompts do lead to more simplification, and most importantly for this study, more elaborative simplification (cf. Table 6). Table 7 shows some examples where our prompts have had the intended effect, as compared to the basic prompt. The last example in Table 7 contains a factual error, which is a typical example of the nature of the factual errors we observed. The insertion of irrelevant or non-factual information is particularly problematic in the context of text simplification, where target users of a simplification will typically have difficulties comprehending the input text and may be less able to discern if the inserted information is factual or not (cf. [Devaraj et al. 2022](#)).

## 5 Conclusion and outlook

We have presented an in-depth analysis of elaborative simplification in German-language texts, using RST annotations and more fine-grained categories. We have experimented with using these insights to improve an LLM’s ability to produce

elaborative simplifications. Our fine-tuned model and our different prompts do encourage the model to insert additional information (see Table 6), increase the level of simplification, and also result in marginal improvements on the SARI and FRE scores. However, the coherence and factuality seem to be adversely affected, indicating that these outputs contain repetitions or so-called hallucinations. This confirms results from related work, where conservative models may produce output texts that preserve the meaning of the input text, but fail to produce simplifications (cf. Cripwell et al. 2024).

As our analysis showed, not all simplified texts contain additional information and certainly not all types of additional information. In most cases, just one type is necessary, i.e. for texts of a political nature, more background knowledge and the relation between the entities in the text may be more important for understanding the text. Future work could investigate on selecting a prompt dependent on the input text.

Adding new information is not trivial; as can be seen in Figure 1, making relations more explicit, for example, can also slightly change the content of a text. In Figure 1b, segment 9 leaves some room for interpretation, as ‘disagreeing’ is not specific, whereas segments 4 to 7 in Figure 1a make this ‘disagreement’ very concrete. By keeping content more open and vague, it is easier to stay ‘factual’, showing that it is a fine line between making relations explicit and staying factual. Overall, elaborative or additive simplification remains a challenging sub-task of automatic simplification.

As shown by our manual evaluation, factuality and meaning preservation seem to represent separate requirements. We therefore advocate for factuality being included as a separate and additional evaluation criterion for text simplification, as up until now faithfulness and factuality seem to have been used interchangeably in the simplification literature, and simplifications are often (manually) evaluated for their meaning preservation (i.e. faithfulness). Our experiments have been limited to fine-tuning and prompting approaches, but experiments which alter the training/fine-tuning paradigm and loss function could also be promising, as at the moment ‘most summarization [and simplification] systems are trained to maximize the log-likelihood of the reference summary at the word-level, which does not necessarily reward models for being faithful’ (Maynez et al., 2020).

## Acknowledgments

We would like to thank Birte Lübbert and Irina Kühnlein for their support with annotating our models’ outputs. We are grateful to the anonymous reviewers for their helpful feedback. This research was funded by a grant from the German Ministry of Education and Research (BMBF).

## References

- Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2019. [Cross-sentence transformations in text simplification](#). In *Proceedings of the 2019 Workshop on Widening NLP*, pages 181–184, Florence, Italy. Association for Computational Linguistics.
- Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2020. [Data-Driven Sentence Simplification: Survey and Benchmark](#). *Computational Linguistics*, 46(1):135–187.
- Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2021. [The \(Un\)Suitability of Automatic Evaluation Metrics for Text Simplification](#). *Computational Linguistics*, 47(4):861–889.
- Marcelo Amancio and Lucia Specia. 2014. An Analysis of Crowdsourced Text Simplifications. In *Proceedings of the 3rd Workshop on Predicting and Improving Text Readability for Target Reader Populations (PITR) @ EACL 2014*, pages 123–130.
- Ursula Bredel. 2016. *Leichte Sprache: theoretische Grundlagen, Orientierung für die Praxis*. Sprache im Blick. Dudenverlag, Berlin.
- Liam Cripwell, Joël Legrand, and Claire Gardent. 2024. Evaluating Document Simplification: On the Importance of Separately Assessing Simplicity and Meaning Preservation. In *3rd Workshop on Tools and Resources for People with READING Difficulties (READI)*, pages 1–14, Torino, Italy.
- Debopam Das and Maite Taboada. 2018. [Signalling of Coherence Relations in Discourse, Beyond Discourse Markers](#). *Discourse Processes*, 55(8):743–770.
- Ashwin Devaraj, William Sheffield, Byron Wallace, and Junyi Jessy Li. 2022. [Evaluating factuality in text simplification](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7331–7345, Dublin, Ireland. Association for Computational Linguistics.
- Freya Hewett. 2023. [APA-RST: A text simplification corpus with RST annotations](#). In *Proceedings of the 4th Workshop on Computational Approaches to Discourse (CODI 2023)*, pages 173–179, Toronto, Canada. Association for Computational Linguistics.

- Ioana Hulpuş, Sanja Štajner, and Heiner Stuckenschmidt. 2019. [A spreading activation framework for tracking conceptual complexity of texts](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3878–3887, Florence, Italy. Association for Computational Linguistics.
- Sarah Jablotschkin, Elke Teich, and Heike Zinsmeister. 2024. [DE-Lite – a New Corpus of Easy German: Compilation, Exploration, Analysis](#). In *Proceedings of the Fourth Workshop on Language Technology for Equality, Diversity, Inclusion*, pages 106–117, Malta. Association for Computational Linguistics.
- Tannon Kew, Alison Chi, Laura Vásquez-Rodríguez, Sweta Agrawal, Dennis Aumiller, Fernando Alva-Manchego, and Matthew Shardlow. 2023. [BLESS: Benchmarking large language models on sentence simplification](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13291–13309, Singapore. Association for Computational Linguistics.
- Philippe Laban, Jesse Vig, Wojciech Kryscinski, Shafiq Joty, Caiming Xiong, and Chien-Sheng Wu. 2023. [SWiPE: A dataset for document-level simplification of Wikipedia pages](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10674–10695, Toronto, Canada. Association for Computational Linguistics.
- Christiane Maaß. 2015. *Leichte Sprache: das Regelbuch*. Number 1 in *Barrierefreie Kommunikation*. Lit, Münster.
- William C. Mann and Sandra A. Thompson. 1988. [Rhetorical Structure Theory: Toward a functional theory of text organization](#). *Text - Interdisciplinary Journal for the Study of Discourse*, 8(3):243–281.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. [On Faithfulness and Factuality in Abstractive Summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Markus Nussbaumer. 1993. [Textbegriff und Textanalyse](#). In Peter Eisenberg and Peter Klotz, editors, *Sprache gebrauchen – Sprachwissen erwerben*, pages 63–84. Klett, Stuttgart.
- Annette Rios, Nicolas Spring, Tannon Kew, Marek Kostrzewa, Andreas Säuberli, Mathias Müller, and Sarah Ebling. 2021. [A New Dataset and Efficient Baselines for Document-level Text Simplification in German](#). In *Proceedings of the Third Workshop on New Frontiers in Summarization*, pages 152–161. Association for Computational Linguistics.
- Neha Srikanth and Junyi Jessy Li. 2021. [Elaborative Simplification: Content Addition and Explanation Generation in Text Simplification](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 5123–5137, Online. Association for Computational Linguistics.
- Regina Stodden, Omar Momen, and Laura Kallmeyer. 2023. [DEplain: A German parallel corpus with intralingual translations into plain language for sentence and document simplification](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16441–16463, Toronto, Canada. Association for Computational Linguistics.
- Renliang Sun, Hanqi Jin, and Xiaojun Wan. 2021. [Document-level text simplification: Dataset, criteria and baseline](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7997–8013, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. [AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts](#). In *CHI Conference on Human Factors in Computing Systems*, pages 1–22, New Orleans LA USA. ACM.
- Yating Wu, William Sheffield, Kyle Mahowald, and Junyi Jessy Li. 2023. [Elaborative Simplification as Implicit Questions Under Discussion](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5525–5537, Singapore. Association for Computational Linguistics.
- Wei Xu, Chris Callison-Burch, and Courtney Napoles. 2015. [Problems in Current Text Simplification Research: New Data Can Help](#). *Transactions of the Association for Computational Linguistics*, 3:283–297.
- Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. [Optimizing Statistical Machine Translation for Text Simplification](#). In *Transactions of the Association for Computational Linguistics*, volume 4, pages 401–415.
- Amir Zeldes. 2016. [rstWeb - A Browser-based Annotation Interface for Rhetorical Structure Theory and Discourse Relations](#). In *Proceedings of NAACL-HLT 2016 System Demonstrations*, pages 1–5, San Diego, CA.

## A Appendix

### A.1 Limitations

In our study we have worked with one fairly homogeneous dataset. Different target groups and different genres will require different kinds of elaborative simplification. For example, [Wu et al. \(2023\)](#), find that definitions are the most common form of elaboration; the target group of their dataset is children.

As parts of our dataset are available online, we do not know if the data was part of the dataset used to pre-train Llama-3.

## A.2 Hyperparameters

We use an NVIDIA V100S with 32GB VRAM for training and inference. Our hyperparameters can be found in Table 8. Note we also ran inference with a temperature of 0.4; the evaluation metrics were lower and so we only include the evaluation of models with this lower temperature.

temperature	0.0001
batch size per device	1
gradient accumulation steps	4
learning rate	3e-5
no. epochs	1
learning rate scheduler type	cosine
adam $\beta_1$	0.9
adam $\beta_2$	0.95

Table 8: Hyperparameters

## A.3 Determining wording for prompts

To determine the wording for the **Background** and **Relation** prompts, we give the following input text and replace the {text with/out inserted information} with either a text from our corpus analysis that has inserted information from the category **Background** or **Relation**, respectively.

**system** You are a helpful assistant and help the user to understand texts.

**user** Can you tell me which text is simpler? Text 1: {text without inserted information} or Text 2: {text with inserted information}

**assistant**

Example {text with background}, the first sentence in bold is removed for the {text without background}:

**Energy has become much more expensive in the past year.** Many households are struggling to pay their energy bills. This is why the Austrian government has decided to introduce a so-called energy cost equalisation scheme. Almost all Austrian households will receive a one-off payment of 150 euros. Households in need will receive an additional 150 euros.

This applies, for example, to the unemployed and people who receive benefits or a very low pension. In this way, the government wants to prevent households from falling into hardship in winter. (3-freitag-28-1-22-b1)

# Examining Gender and Power on Wikipedia Through Face and Politeness

Adil Soubki<sup>□</sup>, Shyne Choi<sup>□</sup>, Owen Rambow<sup>●</sup>

<sup>□</sup>Department of Computer Science, <sup>●</sup>Department of Linguistics

<sup>■</sup>Institute for Advanced Computational Science, Stony Brook University  
asoubki@cs.stonybrook.edu, {shyne.choi, owen.rambow}@stonybrook.edu

## Abstract

We propose a framework for analyzing discourse by combining two interdependent concepts from sociolinguistic theory: face acts and politeness. While politeness has robust existing tools and data, face acts are less resourced. We introduce a new corpus created by annotating Wikipedia talk pages with face acts and we use this to train a face act tagger. We then employ our framework to study how face and politeness interact with gender and power in discussions between Wikipedia editors. Among other findings, we observe that female Wikipedians are not only more polite, which is consistent with prior studies, but that this difference corresponds with significantly more language directed at humbling aspects of their own face. Interestingly, the distinction nearly vanishes once limiting to editors with administrative power.

## 1 Introduction

Brown and Levinson (1987) (henceforth B&L) introduce an influential theory of politeness based on the concept of face, which they claim to be culturally universal. In this theory, face – i.e. the public image one seeks to claim – is a two-sided coin. Agents attend to their desire to have their wants appreciated, which they call positive face, as well as a complementary desire to act unimpeded and maintain freedom, which they call negative face. The face of every agent is ensnared with that of every other agent – agents cannot have their desires appreciated if they cannot appreciate the desires of others. As a result, utterances can raise (+) or threaten (-) the positive (Pos) or negative (Neg) face of the speaker (S) or hearer (H).

A face threat or face raising is not a property of particular linguistic choices, but of communicative intent. If I want to request information from you, then I necessarily need to threaten your negative face, since, if I am successful in communicating

my request to you, I will oblige you to answer and thus I will restrict your choice of actions. In B&L’s theory, discourse participants may choose among various strategies for minimizing threats to face. These strategies are linguistic strategies (for example, using hedges), and the choice of strategy depends on many factors such as cultural conventions and the discourse situation (who is talking to whom under what circumstances).

Work related to NLP has concentrated on studying linguistic manifestations of politeness (Walker et al., 1997; Danescu-Niculescu-Mizil et al., 2013) while largely disregarding the notion of face act. While B&L are frequently cited, the deep insight of their theory comes from a complexity which has been ignored. Their theory is not simply about politeness, but about how politeness, situated in the context of rational action, manifests from a combination of performing face acts to achieve certain goals and using mitigation strategies to lessen the impact of face-threatening acts. Danescu-Niculescu-Mizil et al. (2013) use politeness markers inspired by B&L strategies as features for a system which predicts perceived politeness without modeling face acts. Dutt et al. (2020) predict face acts in isolation from perceived politeness. In this paper, we re-examine the Wikipedia Talk Pages Corpus (Danescu-Niculescu-Mizil et al., 2012; Chang et al., 2020) and demonstrate how bringing face acts and politeness together provides deeper insight.

We do this by producing an annotation of face acts on the corpus and training a new model to label utterances. We then use this tool, along with prior systems which produce judgements of perceived politeness, to label roughly 1.3 million sentences from Wikipedia talk pages. To our knowledge, we are the first to apply an annotation grounded in politeness theory to a text corpus of this scale.

The paper is structured as follows. We start with a review of relevant literature (§2) and present our

theoretical framework (§3). We then turn to producing an annotation of face acts on the Wikipedia Talk Pages Corpus and building a tagger using this new dataset (§4). Our framework is then applied by bringing this new tagger together with existing tools to re-analyze the corpus, paying special attention to gender and power (§5). We end by reporting our conclusions along with a discussion of future work (§6).

All of the code written, datasets prepared, and experimental observations made in the course of this research will be made available on [GitHub](#).<sup>1</sup>

## 2 Related Work

The theory of politeness of B&L has found applications in many fields including sociology, psychology, and linguistics. Google Scholar lists nearly 38,000 citations. Curiously, in NLP there has not been much work building explicitly on B&L. [Danescu-Niculescu-Mizil et al. \(2013\)](#) concentrate on one type of face-threatening act (FTA), namely the negative face-threatening act of a request, and investigate the strategies used for doing this FTA. To do this, they use crowd sourcing to rate the requests on a politeness scale. They develop a model which predicts the politeness of these requests and use it to study the interactions between users on Wikipedia and StackExchange. [Ziems et al. \(2023\)](#) show that fine-tuning on the data of [Danescu-Niculescu-Mizil et al. \(2013\)](#) substantially outperforms zero-shot approaches.

The face acts (FAs) themselves are the object of [Dutt et al. \(2020\)](#). In addition to developing a dataset annotated with FAs, they present a FA classifier based on a neural architecture they devise on top of BERT, which achieves 69% F-measure (0.60 macro). As the data involves participants convincing others to donate to a charity, they also use this corpus to investigate the relationship between face acts and persuasion by predicting if a participant chose to donate. This corpus, which we refer to as the “CMU Face Acts Corpus” (or “CMU Corpus” for short) in this paper, is the direct inspiration for our annotation effort on the Wikipedia data. We differ from their annotation scheme in some important details; we present our annotation in §4. In prior work, we investigated the interaction of intention (through dialog act tagging) and face acts in the CMU Corpus ([Soubki and Rambow, 2024](#)).

There has been an explosion work in compu-

tational social science in general, in which NLP tools are used to extract relevant signals from large amounts of data in order to study a social phenomenon, such as changing attitudes towards certain topics as expressed on social media. For an overview, see ([Edelmann et al., 2020](#)). In the area of studying how gender and power shape written dialogs, there has been some work in NLP. Working with corporate emails, ([Prabhakaran et al., 2014](#)) find that gender differences become exaggerated when looking at individuals with greater social power; specifically, among people with power, women behave *more* differently from men than when comparing people without power.

Finally, turning to the study of politeness and gender outside of NLP, there have been some studies based on manual analysis of collected data, for example ([Herring, 1994](#); [Tannen, 1994](#); [Kunsmann, 2013](#)). For space reasons, we discuss only one example in more detail. [Kendall \(2005\)](#), using a framing approach following ([Goffman, 1974](#)), finds that women in power who “downplay status differences (...) are exercising and constituting their authority by speaking in ways that accomplish work-related goals while maintaining the faces of their interlocutors”. In the terminology of B&L (which [Kendall \(2005\)](#) does not use), women perform similar face acts to men but use strategies to mitigate the effects, which results in women in power appearing more polite than men in power.

## 3 Theoretical Framework

In this section we provide a brief summary of relevant concepts from politeness theory as it relates to our work. Our goal in this paper is to explore how face acts contribute to the perception of politeness. For B&L, “face” refers to the public self-image of agents, and it is a universal component of human interaction. It consists of two complementary facets ([Brown and Levinson, 1987](#), §3.1, p. 61). (1) negative face: “the basic claim to territories, personal preserves, rights to non-distraction – i.e. to freedom of action and freedom from imposition.” (2) positive face: “the positive consistent self-image or ‘personality’ (crucially including the desire that this self-image be appreciated and approved of) claimed by interactants.”

A face act is an intentional communicative act which inherently interacts with the face of the speaker and/or addressees ([Brown and Levinson, 1987](#), §3.2, p. 65). Face acts can threaten (-) or

---

<sup>1</sup><https://github.com/cogstates/wikiface>

Face Act	Mnemonic	Sample Discourse Goals
HNEG-	IMPOSITION	Requests, commands, questions, offers, promises, ...
HPOS-	DISAGREEMENT	Criticism, insults, disapproval, ...
HNEG+	PERMISSIVENESS	Granting permission, making exceptions, ...
HPOS+	AGREEMENT	Seeking common ground, group cohesion, ...
SNEG-	INDEBTEDNESS	Thanking, accepting offers or thanks, commitments, ...
SPOS-	APOLOGIES	Confessions, embarrassment, ...
SNEG+	AUTONOMY	Refusing requests, asserting freedoms, ...
SPOS+	CONFIDENCE	Self-promotion, signaling virtue, ...

Table 1: Face acts with mnemonic label and examples of discourse goals.

affirm (+) the face; they can be about the speaker’s face (S) or the hearer’s (H); and they can be about positive (Pos) or negative (Neg) face. This gives us eight possible face acts, shown in Table 1, where we also provide a short mnemonic names which we will use in this paper, as the terminology of B&L can be unintuitive.

Face acts are part of a larger sequence of choices a speaker makes. First, the speaker chooses a discourse goal or goals (which may form a hierarchy) which will be realized in a speech act (Austin, 1962); then they determine which face acts contribute to the discourse goals; they then choose a strategy to realize this face act, in conformance with the cultural norms of their community which are mutually known by them and the hearer in the communicative context (age, gender, power differential of the discourse participants); and finally, they produce the utterance, which the hearer will perceive as more or less polite, given the discourse goal of the speaker, the communicative context, and the mutually known cultural norms. We see that the notion of “strategy” plays a crucial role in the mediation between face act performance and perceived politeness, and B&L devote a large portion of their study to strategies. Unfortunately, there are no corpora annotated for face act strategies.<sup>2</sup>

We emphasize that face acts do not imply perceived politeness (§B). Consider the following examples from the Wikipedia corpus.

[1] B: Why open a peer review when we are looking for someone to do the GA review?

<sup>2</sup>Danescu-Niculescu-Mizil et al. (2012) use a notion of “strategy” which is defined by a grouping of lexical items that are assumed to affect the hearer’s perception of politeness. They can be considered a simple approximation of the notion in B&L, and in fact helps in predicting politeness. We have chosen not to use these “strategies” (though they are straightforward to determine, as they are based exclusively on word matching), since we would like to address the issue in a more principled manner in the future.

A: Why request a second GA, 3 days after the first one failed?

[2] A: Hi Plange, any reason why this category is named differently to the others?

Both utterances are HNEG-/IMPOSITION face acts, because they impose on the hearer the obligation to respond. However, (1) rejects the previous question by B and challenges B, while (2) is just a request for information, so that (1) is perceived as more impolite than (2).

It is possible for a single utterance to perform multiple face acts at once. For example, (1) could also be seen as DISAGREEMENT, since it entails a critique of B’s actions. However, Dutt et al. (2020) observed multi-labeled acts in only 2% of their data, leading them to consider a single label per utterance. We make this simplification as well in the work presented in this paper.

## 4 Face Act Tagging

In this section we outline the data, modeling techniques, and evaluation measures used in developing our face act tagger for Wikipedia talk pages.

### 4.1 Dataset

On Wikipedia, talk pages are used by editors to coordinate changes and improvements to the encyclopedia.<sup>3</sup> A variety of social and power dynamics are at play in these conversations which can range from discussions of bureaucratic process to heated, and sometimes personal, conflicts. The Wikipedia Talk Pages Corpus (Danescu-Niculescu-Mizil et al., 2012) collects 125,292 exchanges between 38,462 editors resulting in a total of 391,294 posts for analysis. Unlike the CMU Face Acts Corpus, where participants are on mostly level ground, editors can hold administrative privileges or greater notoriety

<sup>3</sup>[https://en.wikipedia.org/wiki/Wikipedia:Talk\\_page\\_guidelines](https://en.wikipedia.org/wiki/Wikipedia:Talk_page_guidelines)

within the community, resulting in interactions with large social distance. Additionally, some editors self-identify gender on their user page.<sup>4</sup> This is desirable in our case as it allows us to study how these social factors interact with face and politeness.

There can be nested replies in talk pages which allow for situations where an utterance is not a reply to the preceding utterance. We do not attempt to correct for these cases and sort first to preserve reply structure and then by the time of the post.

## 4.2 Annotation

Similar to the CMU Corpus, we use the criteria outlined by B&L, which serves as our reference. The CMU Corpus annotation guidelines, as the authors noted, contain some departures from politeness theory. In particular, the CMU Corpus annotates both thanking and complimenting as AGREEMENT. In contrast, B&L analyze thanking and complimenting as INDEBTEDNESS and IMPOSITION, respectively. We choose to remain faithful to B&L, and in fact assert this to be a critical piece of the theory. Consider a compliment such as *you have a lovely smile*. How is it that a compliment can be taken so poorly by the addressee if the speaker is not risking anything? They are often very risky social acts because the speaker assumes they are among the people their addressee wishes to be complimented by; a very imposing assumption. Thanking, on the other hand, can be seen as an exchange of currency. Similar to writing an IOU, the speaker offers a token of their freedom to the addressee. We note that we expect future versions of face act annotations to annotate multiple face acts at once, which may resolve this difference between the CMU Corpus annotation style and ours.

We randomly selected 200 conversations from the WikiTalks data for manual annotation. As the posts contain multiple sentences, each with the possibility of their own face act, we segment the sentences prior to annotation using spaCy (Honnibal and Johnson, 2015). To reduce errors in segmentation, we scrubbed hypertext tags and masked any remaining urls. This resulted in 1850 sentences. We will refer to these basic units of annotation as “utterances” in the following sections. Two of the authors annotated the 1850 utterances for face acts. We examined 100 utterances labeled by both annotators and computed a Cohen’s Kappa of 0.69 which indicates moderate to substantial agreement.

<sup>4</sup>[https://en.wikipedia.org/wiki/Wikipedia:User\\_pages](https://en.wikipedia.org/wiki/Wikipedia:User_pages)

## 4.3 Modeling

We model face act tagging as a text classification task. Given a sequence of  $n$  utterances  $S = [t_1, t_2, \dots, t_n]$ , we wish to assign a label  $y \in Y$  where  $Y$  represents a set containing the 8 possible face acts and one additional label for no face act. Recently, many classification tasks have achieved stronger results using parameter efficient fine-tuning methods of larger models rather than full fine-tuning smaller ones (Hu et al., 2022; Dettmers et al., 2024). We adopt this approach and use Llama-3-8B (AI@Meta, 2024) and LoRA with Int8 quantization (Dettmers et al., 2022) for fine-tuning.<sup>5</sup> Details of the configuration are given in Appendix A.

## 4.4 Data Representation

While fine-tuning approaches unify many aspects of the model design, they present challenges when it comes to determining effective input and output representations.

We provide the models an input which contains an utterance prefixed with the Wikipedia username of the discourse participants,<sup>6</sup> along with previous utterances as context. Each utterance is followed by a newline character. We give an example with two lines of context, though in our experiments we use more, as discussed just below.

```
[Input]
Jossi: I will.
Jossi: Just play nice, that is all I ask.
Kelly: What’s that supposed to mean?

[Output]
hpos-
```

The target output is a distribution where the highest probability is given to the correct label for the final utterance of the input text, in this case HPOS-(DISAGREEMENT). We experimented with different output formats, and found they do not make much of a difference. In our experiments we noticed context to be a critical factor with the optimal size varying by model. Llama 3 performed best with a size of four, for a total of five utterances. As there are no previous turns for the first four turns in each dialog, those examples are provided in a similar format containing only three, two, one or no lines of context.

<sup>5</sup>Our choice of Llama-3 was informed by a preliminary set of experiments in which a variety of pre-trained models and methods were examined on single seed runs.

<sup>6</sup>We note that the Wikipedia usernames shield the actual identity of the discourse participant, and that the Wikipedia username is public.

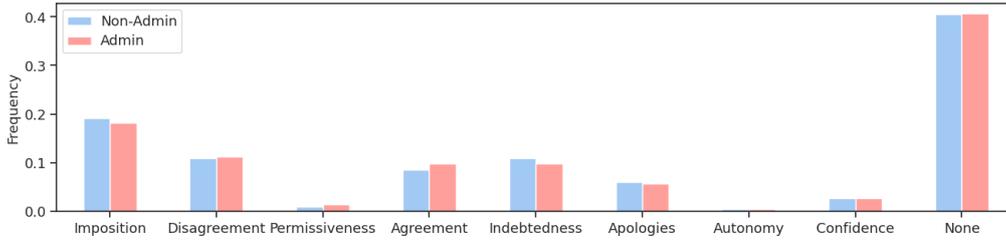


Figure 1: Frequency of face acts for admins and non-admins.

#### 4.5 Experimental Setup And Evaluation

We perform all experiments using five-fold cross validation and the evaluation metrics are averaged across all five folds. We evaluate model performance using F-measure for each of the nine classes as well as micro and macro F-measure aggregated over all labels. We performed hyperparameter tuning, and report metrics only for the best model.

#### 4.6 Results

The results of these experiments are reported in Table 2. We achieve a micro-averaged F1 of 0.68 (average across five folds). Since the task is, with the exception of some nuances (§4.2), identical to the CMU Face Acts Corpus we also tried continued training on the CMU Face Acts Corpus, but this did not improve performance. We suspect this is due to the difference in genre and slight change in annotation procedure, which results in a different distribution of labels between the two datasets.

### 5 Application and Analysis

We apply our new face act tagger along with the politeness scores provided by ConvoKit (Chang et al., 2020) to study the interactions of face and polite-

<b>Micro</b>	0.68
<b>Macro</b>	0.51
<b>IMPOSITION</b>	0.73
<b>DISAGREEMENT</b>	0.56
<b>PERMISSIVENESS</b>	0.40
<b>AGREEMENT</b>	0.58
<b>INDEBTEDNESS</b>	0.80
<b>APOLOGIES</b>	0.56
<b>AUTONOMY</b>	0.04
<b>CONFIDENCE</b>	0.14
<b>NONE</b>	0.76

Table 2: Mean F1 across all folds of our annotation.

ness over the entire Wikipedia Talk Pages Corpus. Our face act tagger is trained using our entire annotation (§4.2) before applying it to the Wikipedia data. This produces roughly 1.3 million sentences labeled with face acts and perceived politeness. We note that the politeness scores are obtained for the entire turn, as this is what the perceived politeness model is trained on, while face acts are tagged by sentence to allow for greater granularity.

In our analysis of politeness we investigate how polite (magnitude) editors are perceived to be by looking at their scores and how often that occurs (frequency) by considering the proportion of utterances in the top 25% of politeness scores. For face acts, we compare the overall distribution (frequency) of labels. Statistical significance is calculated using the Mann-Whitney U test. This analysis was also performed on only the human annotated portion of the data and the trends remained consistent. We report results on the entire corpus.

#### 5.1 Admin Differences

On Wikipedia, editors with administrative status wield significant power in the community including the ability to block or unblock users by IP address and delete or restore pages. This increased status is known to be recognized in the community (Danescu-Niculescu-Mizil et al., 2012; Burke and Kraut, 2008; Leskovec et al., 2010) which endows editors with these powers through public elections. We note that politeness theory anticipates speakers with greater social power than their addressee to more often select strategies that reduce ambiguity and lengthiness. This means opting to perform face threatening acts more often (as opposed to avoiding them all together) and mitigating them through the trade-offs of strategies less often, which one would expect to correspond with a perception of being less polite overall.

We divide utterances by their politeness score into the polite utterances (top 25%), neutral (next 50%) and impolite (bottom 25%). When compar-

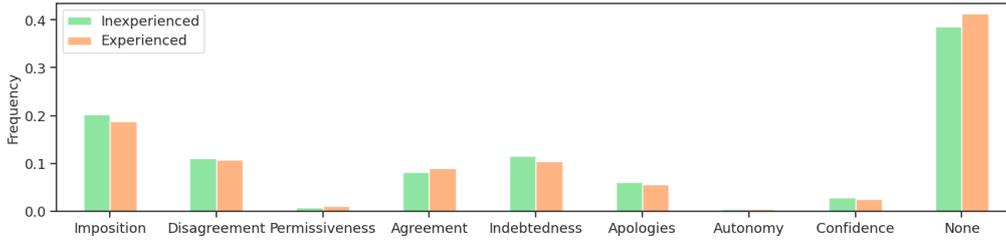


Figure 2: Frequency of face acts by editor experience

ing politeness between admins and non-admins we see the same trend as observed by Danescu-Niculescu-Mizil et al. (2013). Utterances produced by editors with administrative privileges (“admins”) are not more often impolite, however they are significantly ( $p < 0.001$  using the Mann-Whitney U test) less frequently polite, with a mean score difference of 3. Additionally the frequency by which admins produce polite posts is also significantly ( $p < 0.001$ ) lower resulting in messages which are deemed polite 5% less often compared to non-admin editors.

When looking at the distributional differences in face acts by adminship (Figure 1) this decrease in politeness corresponds with small, but salient variations. Admins are significantly ( $p < 0.001$ ) less likely to express INDEBTEDNESS (e.g. thanking, accepting offers) and APOLOGIES (e.g., admitting mistakes, confessions). Though admins produce more utterances labeled AGREEMENT (e.g. appreciation, seeking common ground, group cohesion), their AGREEMENT utterances are significantly ( $p < 0.001$ ) less often (−4% absolute) perceived as polite compared to AGREEMENT utterances by non-admins. Similarly, while non-admins do more IMPOSITION (e.g. issuing commands, making requests), their IMPOSITION utterances are significantly ( $p < 0.05$ ) more often (+3% absolute) taken politely compared to IMPOSITION utterances by admins.. This shows, as we anticipated, that face acts do not imply politeness, contrary to possible intuition.

## 5.2 Experience Differences

We explore whether the experience and productivity of the editor is another means to achieve increased social power without the explicit additional privileges the “admin” title confers. To investigate this we categorize users by the number of edits they have made and label users in the top and bottom quartiles “experienced” and “inexperienced”, respectively.

	Politeness
<b>Experienced Admin</b>	0.34 <sup>†</sup>
<b>Experienced Non-Admin</b>	0.36 <sup>†</sup>
<b>Inexperienced Admin</b>	0.38 <sup>†</sup>
<b>Inexperienced Non-Admin</b>	0.40 <sup>†</sup>

Table 3: Mean politeness scores for difference admin types. All differences are found to be significant using the Mann-Whitney U test with  $p < 0.001$ .

	Inexperienced	Experienced
<b>Impolite</b>	0.07	0.07
<b>Polite</b>	0.35 <sup>‡</sup>	0.28 <sup>‡</sup>

Table 4: Proportion of turns classified as (im)polite by editor experience level. ‡ indicates significance with  $p < 0.0001$  using the Mann-Whitney U test.

We observe similar trends in politeness among experienced editors (Table 4) to that of admins, with turns by experienced editors being labeled polite 7% less often relative to inexperienced editors. When looking at the differences in face acts (Figure 2) we note that there are ways in which newcomers behave like experienced Wikipedians such as a willingness to the face act DISAGREEMENT. However, like admins, experienced users are significantly ( $p < 0.001$ ) less likely to express INDEBTEDNESS or APOLOGIES. Unlike when comparing by admin status, we find that experienced admins are significantly ( $p < 0.001$ ) less likely to interact with face all together (more labeled NONE).

We now investigate how experience interacts with admin status. As expected, experience is correlated ( $r = 0.37$ ) with adminship with nearly half of all admins landing in the top quartile of editors by edit count. We find admins in the top quartile by edit count are significantly ( $p < 0.001$ ) less polite than the bottom quartile. Additionally, intersecting experience with admin status (Table 3) finds a spectrum. Experienced admins are the least polite but

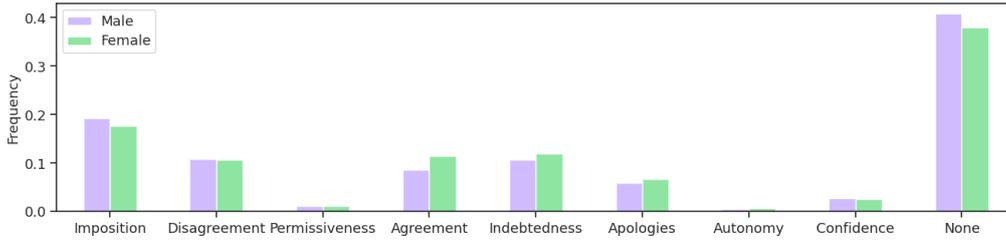


Figure 3: Frequency of face acts by gender.

experienced non-admins are less polite than inexperienced admins. This indicates that these factors are additive in their contribution to social power.

### 5.3 Gender Differences

Some editors self-identify their gender on their user page allowing us to study communicative differences along this axis as well. Prior work found female Wikipedians to be generally more polite (Danescu-Niculescu-Mizil et al., 2013) which is consistent with studies in several domains. We also observe this, with utterances by women scoring more polite (+5,  $p < 0.001$ ), more often (+7%,  $p < 0.0001$ ).

When comparing the distribution of face acts (Figure 3) we see several disparities that the politeness scores alone do not convey. In general, the NONE category is lower for women, i.e. female Wikipedians are more likely, and perhaps more willing, to interact with face in their utterances. When doing so, they humble their own positive face (APOLOGIES, e.g. admitting mistakes, making confessions, accepting compliments) and their own negative face (INDEBTEDNESS, e.g. thanking, accepting apologies) more often than men. This self-deference is accompanied by fewer impositions on their addressee’s face (IMPOSITION, e.g. requests, commands, insults, criticism) and more attention to the hearer’s own wants (AGREEMENT, e.g. seeking common ground, showing respect). Unlike when looking at admins, these AGREEMENT utterances are less frequently judged to be impolite. These trends have been observed in various prior studies (Lakoff, 1973; Prabhakaran and Rambow, 2017; Herring, 1994).

### 5.4 Intersectional Differences

We have seen that male Wikipedians are less polite, more distant with regards to face, and more likely to express IMPOSITION (§5.3). Similarly, much of the same is true when comparing admins to non-admins (§5.1). How do these factors interact? As

	Male	Female
<b>Non-Admin</b> <sup>‡</sup>	0.37	0.43
<b>Admin</b>	0.34	0.35
<b>Inexperienced</b> <sup>‡</sup>	0.41	0.43
<b>Experienced</b> <sup>‡</sup>	0.34	0.42

Table 5: Mean politeness scores by experience and admin status compared across gender. ‡ indicates significance with  $p < 0.0001$  using the Mann-Whitney U test when comparing across gender.

mentioned in §2, previous work in other domains has found gender differences to become exaggerated in the communication patterns of individuals with power. One might expect a similar trend to hold on Wikipedia.

When comparing politeness across both gender and administrative status (Table 5), we find that this does not appear to be the case. While women admins are more polite (magnitude) than male admins, the difference is not significant ( $p > 0.1$ ). Meanwhile, their non-admin counterparts are significantly more polite than non-admin men (+6,  $p < 0.0001$ ). Among non-admin editors, women produce utterances in the top quartile of politeness 10% more often than men, while this reduces to just 1% when comparing admins across genders.

Overall the distribution of face acts (Figure 4) between male and female admins is similar to that of non-admins (the red lines for admins and blue lines for non-admins in Figure 4 are in the same direction), except that the difference between men and women is reduced (the red lines are shorter than the blue lines). There is one striking exceptions: among non-admins, men make many more IMPOSITION (e.g., making requests, issuing commands) face acts than women, but this difference disappears for admins (and in fact women perform IMPOSITION utterances slightly more frequently than men). We note that IMPOSITION is the face act that becoming an admin specifically entitles the

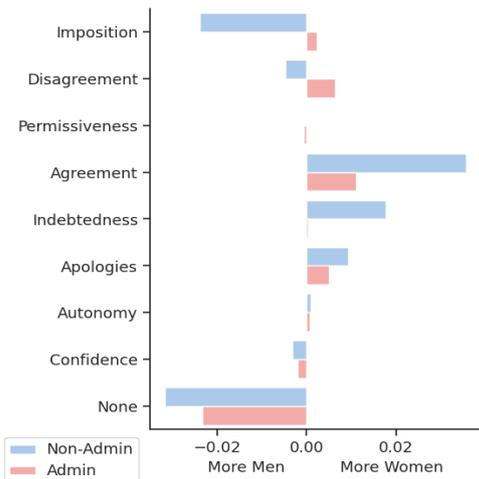


Figure 4: Differences between relative usage of face acts by gender, broken down by non-admins (blue) and admins (red); lines to the right (left) indicate that women (men) perform the face act more often

editor to perform: admins have the right to request changes (and that changes be undone). We speculate that female admins specifically make use of their socially sanctioned power, while men perform IMPOSITION acts even when having no specific admin authority. In summary, admin privileges maintain but substantially lessen the previously observed gender differences in politeness and face. Put differently, female admins behave more like men (whether admins or not), which we also saw in the politeness scores (Table 5).

We now turn to the intersection of gender and experience. Here, we see a strikingly different result. For all conditions (non-admin, admin, inexperienced, experienced), women are more polite. However, we see from Table 5 that men become more impolite as they become experienced, while this is not the case for women: there is no significant change in their politeness as they become experienced. The only exception is for women who become admins (who are, often, experienced), who behave as men do. Put differently, experience and the official power designator of “admin” do not function in the same way across gender: for men, both result in less politeness, but for women, only the “admin” title does.

When looking at face acts (Figure 5), we see that for some categories the differences between men and women are reduced with experience (the orange bars are shorter than the green bars). However, a notable exception is for INDEBTEDNESS, for which we see a large increase in the difference

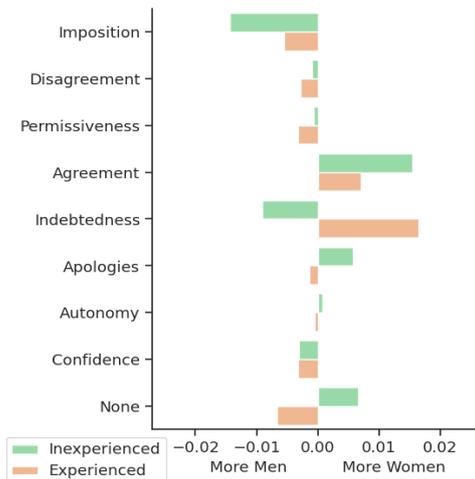


Figure 5: Differences between relative usage of face acts by gender, broken down by inexperienced (green) and experienced (orange) users; lines to the right (left) indicate that women (men) perform the face act more often

between men and women, and in fact a flip in which gender performs it more often. When looking at the absolute numbers (not shown in the table), we can see why: women do not change the frequency of their INDEBTEDNESS utterances at all as they gain experience, while men decrease their frequency of INDEBTEDNESS utterances from 12.3% to 9.8% of their utterances. This decrease is a major contributor to the decrease in politeness among experienced men (but not among experienced women). We extend our previous interpretation by speculating that experienced women do not feel they have a socially sanctioned position of power, and/or men experience a decrease in social distance towards other Wikipedians as they become more experienced, while women do not.

## 6 Conclusion and Future Work

We identify an optimized method for training face act taggers using fine-tuning on LLMs, contribute a new corpus annotated for face acts, and make available a pre-trained model for use on Wikipedia. Through several methods of analysis we demonstrate the usefulness of examining perceived politeness in combination with face acts by reporting a number of findings based on their interaction. In future work we plan to allow multiple face acts per utterance (including for the same segment), and to incorporate the strategy (as conceived of by B&L) more explicitly into our modeling framework.

## Limitations

The principal scientific limitation of this work is that we could only consider three aspects of the larger model of B&L: face acts, the communicative setting (gender and power), and perceived politeness. The major missing elements in the full framework include intention, communicative intention, social norms, and strategies. We intend this paper to be a first step towards a fuller implementation of an explicit cognitive theory of communication which involves all of the mentioned elements.

The experiments for this work were performed using computational resources that are not, in general, freely available. In part due to these computational requirements, but also a result of minimal data, we were not able to evaluate the techniques on additional languages and acknowledge the limitations this places on extending our results to other cultures. We also note along similar lines that while [Brown and Levinson \(1987\)](#) claim their theory of politeness to be culturally universal, this claim has been contested – most notably for eastern cultures ([Al-Duleimi et al., 2016](#)). As discussed in detail above, taking utterances to have a single face act or intent is a critically limiting assumption which lends some uncertainty to our conclusions.

We note that while many of the linguistic differences observed were consistent across multiple rounds of analysis and significant using the Mann-Whitney U test, the effect sizes were generally small. The conclusions should be interpreted with that in mind.

## Ethics Statement

Despite an analysis of the errors, we cannot verify the safety of this system in any user-oriented context and therefore do not recommend such uses without further study. While we do not produce any datasets directly from human annotations, we do use several datasets which were, to the best of our knowledge, compiled ethically. As the primary object of study in this work is the relationship between politeness and language, we do not anticipate broad risks to its application.

## Acknowledgements

This material is based upon work supported by the Defense Advanced Research Projects Agency (DARPA) under the CCU (No. HR001120C0037, PR No. HR0011154158, No. HR001122C0034) program. Soubki has received additional support

from the National Science Foundation (NSF) under No. 2125295 (NRT-HDR: Detecting and Addressing Bias in Data, Humans, and Institutions). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the NSF or DARPA.

We thank both the Institute for Advanced Computational Science and the Institute for AI-Driven Discovery and Innovation at Stony Brook for access to the computing resources needed for this work. These resources were made possible by NSF grant No. 1531492 (SeaWulf HPC cluster maintained by Research Computing and Cyberinfrastructure) and NSF grant No. 1919752 (Major Research Infrastructure program), respectively.

We would also like to thank our anonymous reviewers for their perceptive comments, which improved this work.

## References

- AI@Meta. 2024. [The llama 3 herd of models](#).
- Hutheifa Y. Al-Duleimi, Sabariah Hj Md Rashid, and Ain Nadzimah Abdullah. 2016. A critical review of prominent theories of politeness. *Advances in Language and Literary Studies*, 7:262–270.
- J. L. Austin. 1962. *How to do things with words*. Oxford University Press.
- Penelope Brown and Stephen C Levinson. 1987. *Politeness: Some Universals in Language Usage*. Cambridge University Press.
- Moira Burke and Robert Kraut. 2008. [Taking up the mop: Identifying future wikipedia administrators](#). In *CHI '08 Extended Abstracts on Human Factors in Computing Systems*, CHI EA '08, page 3441–3446, New York, NY, USA. Association for Computing Machinery.
- Jonathan P. Chang, Caleb Chiam, Liye Fu, Andrew Wang, Justine Zhang, and Cristian Danescu-Niculescu-Mizil. 2020. [ConvoKit: A toolkit for the analysis of conversations](#). In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 57–60, 1st virtual meeting. Association for Computational Linguistics.
- Cristian Danescu-Niculescu-Mizil, Lillian Lee, Bo Pang, and Jon Kleinberg. 2012. Echoes of power: Language effects and power differences in social interaction. In *Proceedings of WWW*, pages 699–708.
- Cristian Danescu-Niculescu-Mizil, Moritz Sudhof, Dan Jurafsky, Jure Leskovec, and Christopher Potts. 2013.

- A computational approach to politeness with application to social factors. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 250–259, Sofia, Bulgaria. Association for Computational Linguistics.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Gpt3.int8(): 8-bit matrix multiplication for transformers at scale. In *Advances in Neural Information Processing Systems*, volume 35, pages 30318–30332. Curran Associates, Inc.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- Ritam Dutt, Rishabh Joshi, and Carolyn Rose. 2020. Keeping up appearances: Computational modeling of face acts in persuasion oriented discussions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7473–7485, Online. Association for Computational Linguistics.
- Achim Edelmann, Tom Wolff, Danielle Montagne, and Christopher A. Bail. 2020. Computational social science and sociology. *Annual Review of Sociology*, 46(1):61–81.
- Erving Goffman. 1974. *Frame Analysis: An Essay on the Organization of Experience*. Northeastern University Press, Boston, MA.
- Susan C Herring. 1994. s. In *Cultural performances: Proceedings of the third Berkeley women and language conference*, pages 278–294.
- Matthew Honnibal and Mark Johnson. 2015. An improved non-monotonic transition system for dependency parsing. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1373–1378, Lisbon, Portugal. 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*.
- Damjan Kalajdzievski. 2023. A rank stabilization scaling factor for fine-tuning with lora.
- Shari Kendall. 2005. Creating gendered demeanors of authority at work and at home. In Janet Holmes and Miriam Meyerhoff, editors, *The Handbook of Language and Gender*. Blackwell.
- Peter Kunsmann. 2013. Gender, status and power in discourse behavior of men and women. *Linguistik Online*, 5(1).
- Robin Tolmach Lakoff. 1973. Language and woman’s place. *Language in Society*, 2:45 – 79.
- Jure Leskovec, Daniel P. Huttenlocher, and Jon M. Kleinberg. 2010. Governance in social media: A case study of the wikipedia promotion process. In *Proceedings of the Fourth International AAAI Conference on Weblogs and Social Media*.
- Vinodkumar Prabhakaran and Owen Rambow. 2017. Dialog structure through the lens of gender, gender environment, and power. *Dialogue Discourse*, 8:21–55.
- Vinodkumar Prabhakaran, Emily E. Reid, and Owen Rambow. 2014. Gender and power: How gender and gender environment affect manifestations of power. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1965–1976, Doha, Qatar. Association for Computational Linguistics.
- Adil Soubki and Owen Rambow. 2024. Intention and face in dialog. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9143–9153, Torino, Italia. ELRA and ICCL.
- Deborah Tannen. 1994. *Talking from 9 to 5: Women and Men in the Workplace: Language, Sex and Power*. Avon Books, New York.
- Marilyn A. Walker, Janet E. Cahn, and Stephen J. Whitaker. 1997. Improvising linguistic style: Social and affective bases for agent personality. In *Proceedings of the First International Conference on Autonomous Agents, AGENTS ’97*, page 96–105, New York, NY, USA. Association for Computing Machinery.
- Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. 2023. Can large language models transform computational social science? ArXiv preprint arXiv:2305.03514.

## A Configuration Details for Experiments

For all experiments we fine-tune Llama-3-8B on each of the five cross-validation folds with a batch size of 1 and no gradient accumulation steps. The AdamW optimizer is configured with a learning rate of  $2e-5$ , weight decay of 0, and epsilon of  $1e-8$ . As the cross-validation preparation does not contain a development set to conserve data, we train for a fixed 10 epochs. We configure LoRA with  $\alpha$  of 16, dropout of 0.1, and  $r$  of 64. Since  $r$  is somewhat large, we observed slightly better results using rank-stabilization which scales adapters during forward passes by a factor of  $\alpha/\sqrt{r}$ , instead of the typical  $\alpha/r$  (Kalajdzievski, 2023). These parameters were arrived at through a run of hyperparameter tuning experiments.

## B Supplementary Correlation Analysis

This analysis was performed based our model (§4) output on the Wikipedia Talk Pages Corpus. Aside from INDEBTEDNESS (e.g. thanking, commitments, accepting offers), DISAGREEMENT (e.g. criticism, insults, disapproval), and NONE (avoiding face altogether) the correlations have fairly low magnitude (absolute value less than 0.1).

	<b>Politeness</b>	<b>Impoliteness</b>
<b>IMPOSITION</b>	0.01	0.05
<b>DISAGREEMENT</b>	-0.11	0.18
<b>PERMISSIVENESS</b>	-0.01	0.01
<b>AGREEMENT</b>	0.03	-0.04
<b>INDEBTEDNESS</b>	0.31	-0.25
<b>APOLOGIES</b>	0.04	-0.07
<b>AUTONOMY</b>	0.00	-0.01
<b>CONFIDENCE</b>	-0.01	-0.01
<b>None</b>	-0.17	0.06

Table 6: Pearson’s correlation coefficients between politeness scores and face acts.

# ReALM: Reference Resolution As Language Modeling

Joel Ruben Antony Moniz\*<sup>1</sup>, Soundarya Krishnan\*<sup>2</sup>, Melis Ozyildirim<sup>3</sup>,

Prathamesh Saraf, Halim Cagri Ates, Yuan Zhang, Hong Yu<sup>4</sup>

{<sup>1</sup>joelmoniz, <sup>2</sup>skrishnan22, <sup>3</sup>melisozyildirim, <sup>4</sup>hong\_yu}@apple.com

Apple

## Abstract

Reference resolution is an important problem, one that is essential to understand and successfully handle contexts of different kinds. This context includes both previous turns and context that pertains to non-conversational entities, such as entities on the user’s screen or those running in the background. While LLMs have been shown to be extremely powerful for a variety of tasks, their use in reference resolution, particularly for non-conversational entities, remains underutilized. This paper demonstrates how LLMs can be used to create an effective system to resolve references of various types, by showing how reference resolution can be converted into a language modeling problem, despite involving forms of entities like those on screen that are not traditionally conducive to being reduced to a text-only modality. We demonstrate large improvements over an existing system with similar functionality across different types of references, with our smallest model obtaining absolute gains of over 5% for on-screen references. We also benchmark against GPT-3.5 and GPT-4, with our smallest model achieving performance comparable to that of GPT-4, and our larger models substantially outperforming it.

## 1 Introduction

Human speech typically contains ambiguous references such as "they" or "that", whose meaning is obvious (to other humans) given the context. Being able to understand context, including references like these, is essential for a conversational assistant that aims to allow a user to naturally communicate their requirements to an agent, or to have a conversation with it (Luger and Sellen, 2016; Ljungholm, 2021). In addition, enabling the user to issue queries about what they see on their screen is a crucial step in ensuring a true hands-free experience in voice assistants. For instance, consider

the following interactions between a user and an agent shown in Table 1.

Table 1: Sample Interactions between a user and an agent.

Speaker	Dialogue
User	Show me pharmacies near me
Agent	Here is a list I found.
Agent	... (list presented)
User (eg 1)	Call the one on Rainbow Rd.
User (eg 2)	Call the bottom one.
User (eg 3)	Call this number (present onscreen)

Here, it is immediately apparent that it would not be possible for the Agent to understand or complete the user’s query without the ability to use and comprehend context. It also stands to reason that there are multiple types of context that are necessary to handle user queries: conversational context, on-screen context, and background entities.

Recent Large Language Models (LLMs) (Stammach et al., 2022; Touvron et al., 2023; Santhanam et al., 2022; Dettmers et al., 2023) have often enabled end-to-end experiences, perhaps even obviating the need of a traditional multi-stage pipeline that includes reference resolution (Khatri et al., 2018). There are, however, still several real-world cases where a pipeline is valuable, perhaps even essential, and an end-to-end approach falls short. First, when a framework runs completely on-device (for example, for privacy and efficiency reasons) on a system such as a smartphone that has relatively limited computing power, due to the low-power nature of the system and latency constraints, using a single, large, end-to-end model is infeasible: using a single LLM for this task would usually require the use of a large model with long prompts for true end-to-end experiences (Wei et al., 2022). Second, consider the case when the model has to integrate with APIs, has to consume information from components upstream, or has to provide

\* Equal contribution

information to be consumed downstream: while in these cases it is possible to have an end-to-end approach having the LLM write API calls (Patil et al., 2023; Qin et al., 2023), this often requires a large language model and a complete overhaul of existing pipelines, which might be cumbersome or completely infeasible. Third, the use of a focused model would allow for an existing reference resolution module to be swapped with improved versions in a transparent way, while providing improved ability to hill-climb and improved interpretability, by virtue of the system being modular. Finally, for the task under consideration in this paper, reference resolution does not include solely conversational references, but also includes the ability to reference an on-screen and/or a background entity that is part of what the user currently perceives in their interaction with a device, but has not been a part of the conversational history that results from their direct interaction with the virtual agent in question. There thus continues to be utility in exploring "traditional" NLP tasks such as reference resolution, despite some of the larger language models being able to handle them implicitly. In this work, we thus advocate the use of (relatively) smaller language models, but fine-tuned for specifically and explicitly for the task of reference resolution.

Along similar lines, relying on language modeling alone (Bajaj et al., 2022; Patra et al., 2022; Zheng et al., 2023) has recently shown great promise in being able to handle a variety of tasks (Wang et al., 2018, 2019; Hendrycks et al., 2020; Wei et al., 2021; Chung et al., 2022), such as causal reasoning, linguistic acceptability, question answering, textual entailment and even coreference resolution: Using Language Models (LMs) does exceedingly well on tasks that can be modeled in a sequence-to-sequence fashion. However, the biggest challenge with adopting this technique for the general reference resolution task in the context of a voice assistant lies in resolving references to entities on the screen and using their properties, in other words, getting the LM to, informally speaking, "see". In particular, it is non-obvious how to encode entities on a screen in a manner that is conducive to being resolved by an LM, while also being consistent enough with how conversational entities are encoded to enable the LM to successfully perform reference resolution on both types of entities.

In this work, we propose reconstructing the screen using parsed entities and their locations

to generate a purely textual representation of the screen that is visually representative of the screen content. The parts of the screen that are entities are then tagged, so that the LM has context around where entities appear, and what the text surrounding them is (Eg: call the business number). To the best of our knowledge, this is the first work using a Large Language Model that aims to encode context from a screen.

## 2 Related Work and Motivation

While traditional reference resolution systems have explored conversational and visual/deictic references in great depth (Kottur et al., 2018; Schwartz et al., 2019; Kang et al., 2019), resolving on-screen references is a domain that has been relatively under-explored. However, as shown above, conversational agents on a mobile device need to understand references to the screen, and to support such experiences, to be truly natural. On screen references differ from visual and deictic references for several reasons: they tend to be more structured and highly textual, which enables the use of a lighter model to treat this as a text-only problem without a visual component; further, user queries around on-screen elements often tend to be more action-oriented rather than QA based; finally, they use synthetic screens rather than natural real-world images, which are much easier to parse, but whose distribution completely differs from that on which larger pre-trained image-based systems (such as CLIP (Radford et al., 2021)) tend to be trained. Further, jointly being able to perform conversational and on-screen reference resolution has been even less explored, with prior work often focusing on images and graphics (Willemsen et al., 2023), or UI elements (You et al., 2024).

Vision transformers (Dosovitskiy et al., 2020; Touvron et al., 2021; Liu et al., 2021; Yu et al., 2021) and other pre-trained models have recently gained prominence as a popular first step in tasks that require visual understanding. However, these tend to be trained on natural, real-world images rather than screenshots of on-screen layouts, which have a very different distribution. In addition, these can be extremely expensive to (pre-)train, requiring a very large number of images and several hundred GPU hours (or more). Further, they tend to not perform as well on images heavily embedded with text, and dedicated textual understanding approaches (Xu et al., 2020, 2021; Hwang et al.,

2021a,b; Hong et al., 2022) tend to heavily rely on multiple modules such as bounding box detection and OCR while also relying on good image quality. Joint vision+text models are also substantially more expensive with respect to parameters and computational cost. Finally, these models would need to parse text to be able to perform function (Eg: “call the business number” needs to extract the number associated with the business landline from the raw image), a process which can be complex and compute intensive when bearing in mind that the underlying text and its location on the screen has been referred by the system, and as a consequence can be relatively easily extracted without large, complex models.

The most closely related work which we are aware of, and which we consequently use as our baseline, is that of Ates et al. (2023), an extension of Bhargava et al. (2023) which deals purely with on-screen references; however, it suffers from several drawbacks, which we address in this work. First, these approaches rely on a dedicated “Category module” to deal with type-based references. This module often requires manually on-boarding entities every time a new type is created (a common occurrence in voice assistants, as the supported functionality of the assistant is expanded over time). In addition, such modules often treat each type as distinct, with the similarity of different types ignored. This, in turn, leaves on the table the potential positive transfer that could have happened between semantically related classes (such as “phone number” and “contact”) when data is added for one of those classes. This approach is thus difficult to scale to new entity types and use cases. Second, these systems rely on the use of hand-crafted rule-based textual overlap features, which require heavy feature engineering and tend not to be robust. In addition, these heuristics often do not account for semantic similarity, and are not able to encode real-world understanding or commonsense reasoning. Finally, these methods effectively classify how related each entity is to the query in question independently of all other entities and later threshold them, whereas our current approach directly picks out the most relevant option (or options), while also allowing for no entities to be relevant. Our approach thus additionally has the advantage of removing the reliance on a set threshold, while also providing all the functionality supported in the previous approaches.

### 3 Task

We formulate our task as follows: Given relevant entities and a task the user wants to perform, we wish to extract the entity (or entities) that are pertinent to the current user query. The relevant entities are of 3 different types:

1. **On-screen Entities:** These are entities that are currently displayed on a user’s screen
2. **Conversational Entities:** These are entities relevant to the conversation, which predominantly include those that come from a previous turn. For example, let’s say that the first turn of the user is “Call Mom.”, which is an unambiguous turn that uses a contact called Mom. Shortly after, if the user says “Text her”, the reference “her” needs to be resolved to the contact for “Mom” that was brought up in the previous turn; this contact is thus a conversational entity. Another example might involve an interaction in which the user requests for a list of places or alarms to choose from (or the agent presents one for a user turn such as “Show me pharmacies near me”); each item in this list then becomes a conversational entity for subsequent turns.
3. **Background Entities:** These are relevant entities that come from background processes that might not necessarily be a direct part of what the user sees on their screen or their interaction with the virtual agent; for example, an alarm that starts ringing or music that is playing in the background.

We pose the task of reference resolution as a multiple choice task for the LLM, where the intended output is a single option (or multiple options) from the entities shown on the user’s screen. In some cases, the answer could also be “None of these”, in which case the model needs to predict “0”.

To evaluate this task, we check if the predicted set of options matches the ground truth set; in other words, we allow the model to output the relevant entities in any order, i.e. if the Ground Truth is entities 8, 7, and 4, then we accept any permutation of these three correct entities while evaluating the performance of the model.

Note that as in Ates et al. (2023); Bhargava et al. (2023), we assume that entities along with their types come in from an upstream system (for example, through a mechanism involving entity pullers which are able to extract entities in a high recall

manner or through a donation from a device app, as in Aas et al. (2023)).

## 4 Datasets

Our datasets comprise data that was either synthetically created, or created with the help of annotators. Each data point contains the user query and a list of entities, along with the ground-truth entity (or set of entities) that are relevant to the corresponding user query. Each entity, in turn, contains information about its type and other properties such as the name and other textual details associated with the entity (the label and time of an alarm, for example). For data points where relevant on-screen context exists, this context is available in the form of the bounding box of the entity, and the list of objects surrounding it along with properties of these surrounding objects such as their types, textual contents and locations. Note that our data collection follows that of Bhargava et al. (2023); Ates et al. (2023); we present an overview here and direct the interested reader to the aforementioned papers for a more detailed description. Note also that each dataset below is somewhat representative of one of our tasks of interest (with our synthetic data bucket being used for both conversational and background entity resolution).

Table 2: Dataset Sizes (Train Set and Test Set)

Dataset	Train	Test
Conversational	2.3k	1.2k
Synthetic	3.9k	1.1k
On-screen	10.1k	1.9k

### 4.1 Conversational Data

In this case, data is collected for entities that are relevant to the user interaction with the agent. To do this, annotators are shown sample conversations between a user and an agent with synthetic lists of entities provided, and asked to provide queries that unambiguously reference an arbitrarily picked entity in the aforementioned synthetic list. Annotators might thus be provided with a synthesized list of businesses or alarms and asked to refer to a particular entity within that list.

For example, the annotator might be shown a list of businesses that are synthetically constructed, and then asked to refer to a specific one in the list provided; for instance, they might say “Take me to

the one that’s second from the bottom” or “Call the one on Main Street”.

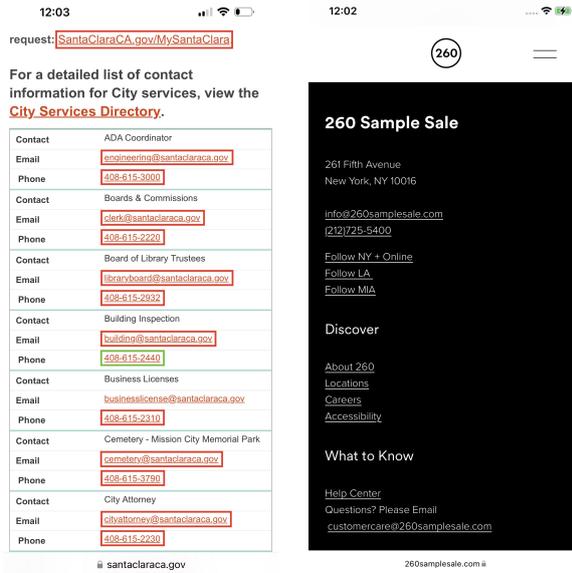
### 4.2 Synthetic Data

Another approach to obtain data is to rely on synthetic data from templates. This approach is particularly useful for type-based references, when the user query and the entity type are sufficient to resolve the reference, and descriptions are not relied upon. Note that the synthetic nature of this dataset does not preclude it from containing datapoints in which multiple entities can be resolved to a given reference: for example, for the query “play it”, “it” can be resolved to all entities of both the types “music” and “video”.

The pipeline used to generate the synthetic data comprises of two parts: a set of templates and a list accompanying each template. The first part, a “language template”, contains different variations of queries that can be used for targeted cases, with slots present that can be filled pragmatically from those defined in the “slot list”. The second, a “slot list” accompanying the aforementioned template, includes mentions and other possible slot values (often comprising of named entities that aren’t mentions, or other slots that can take a large number of possible values such as date-times) if necessary. The slot list also contains the ground truth entity (or entities) that the mentions listed, when filled into the language template, could resolve to.

The data generation pipeline then takes the language template and slot list, and uses them to generate the possible queries. It does this by substituting corresponding values from the slot lists into the language templates to obtain fully formed user queries. The corresponding synthetic data is formed by using these queries and the ground truth entities present in the slot list, and adding in entities of other types into the data to serve as random negatives.

For example, a given language template might consist of phrases like “share [mention] with [name]” and “send [mention] to [name] please”. The corresponding slot list might have “[mention]” mapping to “this address” and “that address”, “[name]” mapping to various person names, and the ground truth entity tagged as “email address” and “physical address”. The pipeline then generates queries like “share that address with Mom” with “email address” and “physical address” entities marked as possible ground truth entity types, and entities of other types marked as negative.



(a) Screenshot example used in first annotation project (b) Screenshot example used in second annotation project

Figure 1: Sample screenshots used in the annotation of on-screen data. The data was annotated in a two-step process, as described in Section 4.3.

### 4.3 On-screen Data

As in Bhargava et al. (2023), screen data were collected from various web pages where phone number, e-mail and/or physical address information exist. Our on-screen data annotation comprised of a two-phase process. The first phase was used to obtain queries based on the screens shown, and the second one was for identifying the entities and mention for the given query. In the first grading project, annotators were given a screenshot (Figure 1a) with green and red boxes, and were asked to classify the green boxed data into one of the entities such as phone number, email address, etc. Then, annotators were then asked to provide three unique queries for the green boxed data.

In the second annotation project (Figure 1b), queries collected in the first step were shown to annotators one by one with their corresponding screenshots (but this time, without the bounding boxes), and with all the screen entities as a list. The annotators were asked if the query contains a mention to one of the given visual entities, and if the query sound natural. They were also asked to provide the entities from the list that were referred to in the given query, and to tag the part of the query referring that entity.

## 5 Models

We compare our proposed model ReALM, described in detail in Section 5.3 below, with two baseline approaches: one based on the reference resolver proposed in MARRS (Section 5.1), which is non-LLM based, and one based on ChatGPT (both GPT-3.5 and GPT-4; Section 5.2).



(a) Conversational User Turns (b) Onscreen Capture

Figure 2: Technical diagrams representing user turns with a conversational assistant in (a), and a user screen in (b). Shaded rectangles represent various elements shown on the screen detectable by screen parser-extractors.

### 5.1 MARRS

As a baseline, we compare against the system proposed in Ates et al. (2023), in turn a variation of Bhargava et al. (2023), both of which are non-LLM based approaches. While the latter approach focuses on on-screen entities, MARRS extends it to conversational and background entities as well. For our baseline comparison, we trained a re-implementation of this system with the datasets described in Section 4, which includes conversation, on-screen and synthetic data. Note that in contrast to our approach, which uses a generic off-the-shelf LLM, this baseline we compare against was specifically designed for the task of reference resolution.

### 5.2 ChatGPT

As another baseline, we run the GPT-3.5 (Brown et al., 2020; Ouyang et al., 2022) and GPT-4 (Achiam et al., 2023) variants of ChatGPT, as available on 2024-01-24, with in-context learning. As in our setup, we aim to get both variants to predict a list of entities from a set that is available. In the case of GPT-3.5, which only accepts text, our input consists of the prompt alone; however, in the

case of GPT-4, which also has the ability to contextualize on images, we provide the system with a screenshot for the task of on-screen reference resolution, which we find helps substantially improve performance. Note that our ChatGPT prompt and prompt+image formulation are, to the best of our knowledge, in and of themselves novel. While we believe it might be possible to further improve results, for example, by sampling semantically similar queries up until we hit the prompt length, this more complex approach deserves further, dedicated exploration, and we leave this to future work.

### 5.3 Our Approach

In this section, we provide examples of conversational and onscreen reference resolution tasks, followed by how we prompt the model to resolve the same.

We use the following pipeline for fine-tuning an LLM (a FLAN-T5 model (Chung et al., 2022)) in our case. We provide the parsed input to our model, and finetune it. Note that unlike for the baseline, we do not run an extensive hyperparameter search on the FLAN-T5 model, sticking to the default fine-tuning parameters.

Select which among the following entities, if any, are required to understand the user request below. Output 0 if none of the entities are relevant.  
 User request: Call the one on Rainbow St  
 User Entities:  
 0. None  
 1. Type: Local Business | Name: Walgreens | Address: 225 Rainbow St, San Jose CA 94088  
 2. Type: Local Business | Name: CVS | Address: 105 E El Camino Real, Sunnyvale, CA 94087  
 3. Type: Local Business | Name: Qwark | Address: 1287 Hammerwood Ave, Sunnyvale, CA  
 Relevant entity:

Select which among the following entities, if any, are required to understand the user request below. Output 0 if none of the entities are relevant.  
 User request: Save the phone number at the bottom-right Screen:  
 Your New home!  
 Steven Realtors Inc.  
 Trusted by over 5 million  
 Proud homeowners  
 Contact Us  
 Monday - Saturday -  
 Friday Sunday  
 {{1. (206) 198 1999}} {{2. (206) 198 1699}}  
 Relevant entity:

Each data point consisting of a user query and the corresponding entities is converted into a sentence-wise format that we can feed to an LLM

for training. Examples of the input before and after processing are shown in Appendix Sections A and C, with examples of how we convert entities of different types into text shown in Appendix B. Note that the entities are shuffled before being sent to the model so that the model does not overfit to particular entity positions.

With respect to the output that the model predicts, empirically, we find that the model is consistently able to predict a valid integer (or list of integers), without deviating and outputting any other text. In addition, we observe that the model also respects general output constraints (such as not predicting a ‘0’ that represents ‘None of These’ at the same time as one or more other entities) as well as those constraints enforced by the input (such as ensuring all predicted entity indices actually exist on the input side). The one exception that we observe is that, on occasion, we find that the model predicts the same entity twice (successively) in its output list. The only post-processing heuristic we apply is thus to convert the model’s predictions into a set of unique entities.

#### 5.3.1 Conversational References

For the sake of this work, we assume conversational references to be of two types: type-based and descriptive. Type-based references are heavily reliant on using the user query in conjunction with the types of the entities to identify which entity (of a set of entities) are most relevant to the user query in question: for example, if the user says “play this”, we know that they are referring to an entity like a song or a movie, as opposed to a phone number or an address; “call him” likewise refers to a contact or possibly a phone number, as opposed to an alarm. Descriptive references, in contrast, tend to use a property of the entity to uniquely identify it: “The one in Times Square” for example might help uniquely refer to one among a set of addresses or business. Note that it is often the case that references might rely on both types and descriptions to unambiguously refer to a single object: consider the examples “play the one from Abbey Road” vs “directions to the one on Abbey Road”, both of which rely on both the entity type and description to identify a song in the first case and address in the second. In our proposed approach, we simply encode the type and various properties of the entity. We show our detailed encoding scheme in Appendix B.

### 5.3.2 Onscreen References

For onscreen references, as in [Bhargava et al. \(2023\)](#), we assume the presence of upstream data detectors that are able to parse screen text to extract entities. These entities are then available along with their types, bounding boxes and a list of non-entity text elements surrounding the entity in question.

To encode these entities (and thereby, the relevant parts of the screen) into the LM in a manner that involves text alone, we use the novel algorithm given in Algorithm 1. Intuitively, we assume the location of all entities and their surrounding objects to be representable by the center of their respective bounding boxes. We then sort these centers (and thereby, the associated objects) from top-to-bottom (i.e., vertically, along the y-axis), and then use a stable sort to sort from left-to-right (i.e., horizontally, along the x-axis). Next, all objects that are within a margin are treated as being on the same line, and are separated from each other by a tab; objects further down outside the margin are placed on the next line, and this is repeatedly, effectively encoding the screen in a left-to-right, top-to-bottom fashion in plain text.

## 6 Results

Table 3: Model Accuracy for Different Datasets. A prediction is correct if the model correctly predicts all relevant entities, and incorrect otherwise. **Conv** refers to the Conversational Dataset, **Synth** to the Synthetic one, **Screen** to the Onscreen one and **Unseen** to a conversational dataset pertaining to a held-out domain.

Model	Conv	Synth	Screen	Unseen
<b>MARRS</b>	92.1	99.4	83.5	84.5
<b>GPT-3.5</b>	84.1	34.2	74.1	67.5
<b>GPT-4</b>	97.0	58.7	90.1	98.4
<b>ReALM-80M</b>	96.7	99.5	88.9	99.3
<b>ReALM-250M</b>	97.8	99.8	90.6	97.2
<b>ReALM-1B</b>	97.9	99.7	91.4	94.8
<b>ReALM-3B</b>	97.9	99.8	93.0	97.8

We present our results in Table 3. Overall, we find that our approach outperforms the MARRS model in all types of datasets. We also find that our approach is able to outperform GPT-3.5, which has a significantly larger number of parameters than our model by several orders of magnitude. We also find that our approach performs in the same ballpark as the latest GPT-4 despite being a much lighter

---

### Algorithm 1: Onscreen Parse Construction with Turn Object Injection

---

**Data:** List of turn objects

**Result:** Onscreen parse

```

1 onscreen_parse ← Empty list of onscreen
  parse elements;
  // Step 0: Get all text boxes
  present in the screen
2 for each turn object t, index i do
  // Step 1: Get unique
  surrounding objects
3 surrounding_objects ← Set of
  surrounding objects for t;
  // Step 2: Insert turn objects
  into the set
4 surrounding_objects ←
  surrounding_objects ∪ {[i.t]};
  // Step 3: Sorting the centers of
  all surrounding objects
5 sorted_objects ← Sort objects in
  surrounding_objects by center (Top →
  Bottom, Left → Right);
  // Step 4: Determine vertical
  levels
6 margin ← Margin for considering objects
  at the same level;
7 levels ← List of vertical levels;
8 for each object o in sorted_objects do
9   same_level ← List of objects at the
  same level as o;
10  for each object other in
  sorted_objects do
11    if o is not the same as other and
  |o.center_top −
  other.center_top| ≤ margin
  then
12    same_level ←
  same_level ∪ {other};
13  levels ← levels ∪ {same_level};
  // Step 5: Construct onscreen parse
14 for each level l in levels do
15   level_parse ← Empty string;
16   for each object obj in l do
17     level_parse ←
  level_parse + "\t" + obj;
18  onscreen_parse ←
  onscreen_parse + "\n" + level_parse;
19 return onscreen_parse;

```

---

(a) Semantic Understanding	(b) Summarisation
<b>User Request:</b> Call the evening Number	<b>User Request:</b> Remind me to get printouts before the tax deadline
<b>Screen:</b> {{1. 9 AM - 5 PM}} {{2. 901.969.3120}} {{3. 5 PM - 9 PM}} {{4. 901.969.3391}}	<b>Screen:</b> Tax Deadlines 2023 {{1. Feb 15}} Reclaim your tax exemption from withholding {{2. April 18}} First-quarter estimated tax payment due
<b>Model Output:</b> 4	<b>Model Output:</b> 2
(c) World Understanding	(d) Commonsense Reasoning
<b>User Request:</b> Take me to the one in Washington	<b>User Request:</b> Save the link to the breakfast Recipe
<b>Screen:</b> Indian Embassy {{1. 1701 El Camino Real, Mountain View 94040}} {{2. 333 Dexter Ave N, Seattle 98109}} {{3. 8295 Tournament Drive, Memphis, TN 38125}}	<b>Screen:</b> IMAGE Strawberry Granola {{1. Recipe link}} IMAGE Lavender boba tea {{2. Recipe link}}
<b>Model Output:</b> 2	<b>Model Output:</b> 1

Table 4: Qualitative examples that demonstrate the ability of ReALM to adapt to complex use-cases.

(and faster) model. We especially wish to highlight the gains on onscreen datasets, and find that our model with the textual encoding approach is able to perform almost as well as GPT-4 despite the latter being provided with screenshots.

Additionally, we also experiment with models of different sizes. We see that while performance in general improves across all dataset families with an increase in model size, the difference is most pronounced for the onscreen datasets, which alludes to the task being more complex in nature. Interestingly, and contrary to an otherwise consistent trend of larger models performing better, we find that performance on our Unseen dataset, which contains a held-out domain, first decreases with an increase in model size before increasing again. We hypothesize that this is due to the double-descent phenomenon (Nakkiran et al., 2019).

## 6.1 Analysis

### GPT-4 $\approx$ ReALM $\gg$ MARRS for new use-cases:

As a case study, we explore zero-shot performance of this model on an unseen domain: Alarms (we show a sample data point in Appendix Table 12). The last column in Table 3 compares the performance of all approaches and baselines on this unseen test set. We find that all of the LLM-based approaches outperform the FT model for this test set. Among the two, we find that the performance of ReALM and GPT-4 are very similar for the un-

Table 5: User Request for Setting or Home Device

<b>User Request:</b> Can you make it brighter?
<b>Entities Shown to User:</b> 1. Type: Settings 2. Type: UserEntity   homeAutomationAccessoryName
<b>GPT-4 Prediction: 1 Ground Truth:</b> 1, 2

seen domain. Additionally, Table 4 shows completely new experiences enabled by ReALM due to the LLM’s superior ability to perform complex understanding of natural language.

### ReALM $>$ GPT-4 for domain-specific queries

We find that due to finetuning on user requests, ReALM is able to understand more domain-specific questions. Consider Table 5. GPT-4 incorrectly assumes the reference to be about only a setting, whereas the ground truth consists of a home automation device in the background as well, and GPT-4 lacks the domain knowledge to be able to recognise that the device would also be relevant to this reference. ReALM, in contrast, doesn’t suffer from this due to being trained on domain-specific data.

## 7 Conclusion and Future Work

In this work, we demonstrate how large language models, which are typically trained on text alone, can be also be adapted to perform reference resolution to items in an extra-linguistic context. We

do this by encoding entity candidates as natural text; we demonstrate how entities that are present on the screen can be passed into an LLM using a novel textual representation that effectively summarizes the user’s screen while retaining relative spatial positions of these entities. Our proposed system is thus able to resolve references in a variety of human-computer interaction settings, such as those involving on-screen, conversational and background entities; we note, however, that our proposed approach focuses primarily on anaphoric and deictic references, and we leave the extension of our system to handle other types of references, such as bridging references, to future work.

In addition, we show that ReALM outperforms previous approaches, and performs roughly as well as the current state-of-the-art LLM, GPT-4, despite consisting of far fewer parameters, even for on-screen references despite being purely in the textual domain. It also outperforms GPT-4 for domain-specific user queries, thus making ReALM an ideal choice for a practical reference resolution system that can exist on-device without compromising on performance.

While our approach is effective in encoding the position of entities on the screen, we find that it may not be able to resolve complex user queries that rely on nuanced positional understanding. We thus believe that exploring more complex approaches such as splitting the screen into a grid and encoding these relative spatial positions into text, while challenging, is a promising avenue of future exploration. In addition, in contrast to a critical assumption of our proposed system, not all on-screen entities are textual. While extending this paper to cover on-screen images, graphics and UI elements is beyond the scope of this work, this is certainly another extension that merits further investigation.

## Ethics Statement

While LLMs can generate unexpected output including potentially harmful text, our system offers the ability to constrain decoding or use simple post-processing to ensure this does not happen. Note however that practically we find very little hallucination or even text that deviates from the format that the models were finetuned on, and thus do not constrain the decoding of the LLM.

## Acknowledgements

The authors would like to thank Nidhi Rajshree, Stephen Pulman, Leon Liyang Zhang, Jiarui Lu, Jeff Nichols, Shruti Bhargava, Dhivya Piraviperumal, Junhan Chen and the anonymous reviewers for their help, suggestions, and feedback.

## References

- Cecilia Aas, Hisham Abdelsalam, Irina Belousova, Shruti Bhargava, Jianpeng Cheng, Robert Daland, Joris Driesen, Federico Flego, Tristan Guigue, Anders Johannsen, et al. 2023. Intelligent assistant language understanding on device. *arXiv preprint arXiv:2308.03905*.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Halim Cagri Ates, Shruti Bhargava, Site Li, Jiarui Lu, Siddhardha Maddula, Joel Ruben Antony Moniz, Anil Kumar Nalamalapu, Roman Hoang Nguyen, Melis Ozyildirim, Alkesh Patel, et al. 2023. Marrs: Multimodal reference resolution system. In *Proceedings of The Sixth Workshop on Computational Models of Reference, Anaphora and Coreference (CRAC 2023)*, pages 51–58.
- Payal Bajaj, Chenyan Xiong, Guolin Ke, Xiaodong Liu, Di He, Saurabh Tiwary, Tie-Yan Liu, Paul Bennett, Xia Song, and Jianfeng Gao. 2022. Metro: Efficient denoising pretraining of large scale autoencoding language models with model generated signals. *arXiv preprint arXiv:2204.06644*.
- Shruti Bhargava, Anand Dhoot, Ing-Marie Jonsson, Hoang Long Nguyen, Alkesh Patel, Hong Yu, and Vincent Renkens. 2023. Referring to screen texts with voice assistants. In *ACL*.
- 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.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias

- Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Teakgyu Hong, Donghyun Kim, Mingi Ji, Wonseok Hwang, Daehyun Nam, and Sungrae Park. 2022. Bros: A pre-trained language model focusing on text and layout for better key information extraction from documents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10767–10775.
- Wonseok Hwang, Hyunji Lee, Jinyeong Yim, Geewook Kim, and Minjoon Seo. 2021a. [Cost-effective end-to-end information extraction for semi-structured document images](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3375–3383, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Wonseok Hwang, Jinyeong Yim, Seunghyun Park, Sohee Yang, and Minjoon Seo. 2021b. [Spatial dependency parsing for semi-structured document information extraction](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 330–343, Online. Association for Computational Linguistics.
- Gi-Cheon Kang, Jaeseo Lim, and Byoung-Tak Zhang. 2019. Dual attention networks for visual reference resolution in visual dialog. 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 2024–2033.
- Chandra Khatri, Behnam Hedayatnia, Anu Venkatesh, Jeff Nunn, Yi Pan, Qing Liu, Han Song, Anna Gottardi, Sanjeev Kwatra, Sanju Pancholi, et al. 2018. Advancing the state of the art in open domain dialog systems through the alexa prize. *arXiv preprint arXiv:1812.10757*.
- Satwik Kottur, José MF Moura, Devi Parikh, Dhruv Batra, and Marcus Rohrbach. 2018. Visual coreference resolution in visual dialog using neural module networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 153–169.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022.
- Alice Ljungholm. 2021. Voice interaction vs screen interaction when controlling your music-system. In *Proceedings of the 21st Student Conference in Interaction Technology and Design*, pages 103–108.
- Ewa Luger and Abigail Sellen. 2016. "like having a really bad pa" the gulf between user expectation and experience of conversational agents. In *Proceedings of the 2016 CHI conference on human factors in computing systems*, pages 5286–5297.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. 2019. Deep double descent: Where bigger models and more data hurt. In *International Conference on Learning Representations*.
- 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.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*.
- Barun Patra, Saksham Singhal, Shaohan Huang, Zewen Chi, Li Dong, Furu Wei, Vishrav Chaudhary, and Xia Song. 2022. Beyond english-centric bitexts for better multilingual language representation learning. *arXiv preprint arXiv:2210.14867*.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- 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](#).
- Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2022. [COLBERTv2: Effective and efficient retrieval via lightweight late interaction](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3715–3734, Seattle, United States. Association for Computational Linguistics.
- Idan Schwartz, Seunghak Yu, Tamir Hazan, and Alexander G Schwing. 2019. Factor graph attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2039–2048.
- Dominik Stammach, Maria Antoniak, and Elliott Ash. 2022. [Heroes, villains, and victims, and GPT-3: Automated extraction of character roles without training data](#). In *Proceedings of the 4th Workshop of Narrative Understanding (WNU2022)*, pages 47–56,

- Seattle, United States. Association for Computational Linguistics.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. 2021. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*.
- 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*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *Transactions on Machine Learning Research*.
- Bram Willemsen, Livia Qian, and Gabriel Skantze. 2023. Resolving references in visually-grounded dialogue via text generation. In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 457–469.
- Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, Min Zhang, and Lidong Zhou. 2021. LayoutLMv2: Multi-modal pre-training for visually-rich document understanding. 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 2579–2591, Online. Association for Computational Linguistics.
- Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. 2020. Layoutlm: Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1192–1200.
- Keen You, Haotian Zhang, Eldon Schoop, Floris Weers, Amanda Swearngin, Jeffrey Nichols, Yinfei Yang, and Zhe Gan. 2024. Ferret-ui: Grounded mobile ui understanding with multimodal llms. *arXiv preprint arXiv:2404.05719*.
- Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. 2021. Vector-quantized image modeling with improved vqgan. *arXiv preprint arXiv:2110.04627*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.

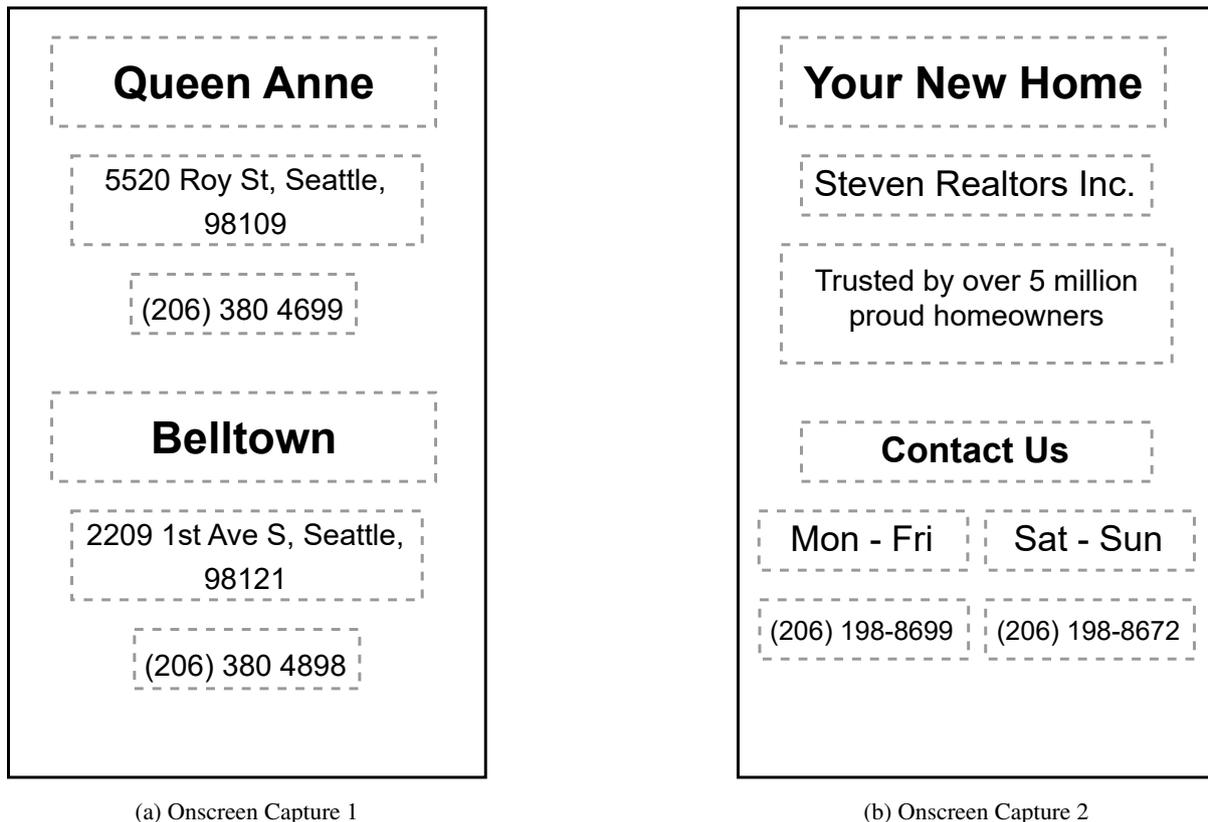


Figure 3: Technical diagrams representing user screens. Shaded rectangles represent various elements shown on the screen detectable by screen parser-extractors.

## A Encoding onscreen entities

First, we show sample representations of what a screen grab might look like, as parsed and exposed to the system. We show these representations in Figure 3

We now describe some other strategies of encoding on-screen elements that we explored.

- Clustering:** We explored a clustering-based approach wherein we performed a spatial clustering of the various surrounding objects present in the screen. We did this to establish semantic clusters wherein a user could refer to nearby bounding boxes (such as the contact information) by a particular title. The detailed approach is given in Algorithm 2, and a sample encoding is shown in Table 6. The biggest drawback of the approach was that the prompt length often explodes as the number of entities in a cluster increases, as each of the objects in the cluster would have every other object in its surrounding objects.
- Onscreen Grab:** To mitigate this issue, we employed a second approach (similar to our

final approach), wherein we parsed the screen as in our final approach, the only difference being that we didn't annotate the turn objects within the parse itself, but provided the turn objects as a list instead (see Table 7).

- Onscreen Grab with Injected Turn Objects:** Finally, the exact algorithm employed in our final approach is given in 1, and a sample encoding is shown in Table 8.

We show an ablation in Figure 4, in which we

Table 6: Clustering-based encoding

<b>User Request:</b> Get me directions to the branch in Queen Anne
<b>Entities Shown to User:</b> 1. Type: Postal Address   Value: 5520 Roy St, Seattle 98109   surr_objects: Queen Anne, (206) 380 4699 2. Type: Phone Number   Value: (206) 380 4699   surr_objects: Queen Anne, 5520 Roy St, Seattle 98109 3. Type: Phone Number   Value: (206) 380 4898   surr_objects: Belltown, 2209 1st Ave S, Seattle 98121 4. Type: Postal Address   Value: 2209 1st Ave, Seattle 98121   surr_objects: Belltown, (206) 380 4898
<b>Ground Truth:</b> 1

Table 7: Onscreen Grab encoding

<b>User Request:</b> Save the phone number at the bottom-right
<b>Screen:</b> Your New home! Steven Realtors Inc. Trusted by over 5 million Proud homeowners Contact Us Monday - Saturday - Friday Sunday (206) 198 1699 (206) 198 1999
<b>Entities Shown to User:</b> 1. Type: Phone Number   Value: (206) 198 1999 2. Type: Phone Number   Value: (206) 198 1699
<b>Ground Truth:</b> 1, 2

show the performance of the various encoding approaches described above (and some other hill-climbing efforts).

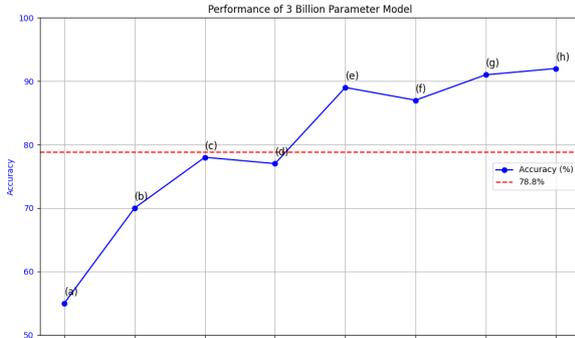


Figure 4: Performance improvements with each experiment – (a) Baseline Finetuned LLM, (b) Obtaining screen elements through OCR, (c) Obtaining screen elements through UI elements and Clustering (d) Adding an extra newline between the instruction and user request, (e) Onscreen Grab, (f) Onscreen Grab with injected turn objects, (g) Onscreen Grab with injected turn object + needing lines to be separated by at least Margin, (h) Separating elements in the same line by a tab

We show the algorithm used to encode onscreen

Table 8: Injected Onscreen Encoding (Final Approach)

<b>User Request:</b> Save the phone number at the bottom-right
<b>Screen:</b> Your New home! Steven Realtors Inc. Trusted by over 5 million Proud homeowners Contact Us Monday - Saturday - Friday Sunday {{1. (206) 198 1999}} {{2. (206) 198 1699}}
<b>Ground Truth:</b> 1, 2

**Algorithm 2:** Surrounding Object Clustering and Prompt Generation

**Data:** List of MDF turn objects

**Result:** Updated turn objects with surrounding object prompts

```

1 for each MDF turn object t do
  // Step 1: Get unique
  // surrounding objects
2 surrounding_objects ← Set of
  // unique surrounding objects for t;
  // Step 2: Spatially cluster
  // surrounding object bounding
  // boxes
3 clusters ←
  DBScan(surrounding_objects,
4 rect_distance);
  // Step 3: Predict the cluster
  // for turn object
5 t_cluster ← Predicted cluster for t;
6 for each surrounding object s in
  surrounding_objects do
7   if s belongs to cluster t_cluster
  then
  // Step 4: Process
  // non-overlapping
  // surrounding objects
8   if no string overlap between t
  and s then
9     Add s to the prompt under
  key 'surrounding_object';
  // Step 5: Provide global
  // positioning information
10 t.distance_from_top ← Compute
  // distance from the top for t;
11 t.distance_from_left ← Compute
  // distance from the left for t;
12 return prompt;

```

entities, described in Section 5.3.2, in Algorithm 1.

Table 9: Entity Domains and their Representations

Entity Type	After
alarm	Type: Alarm   time: 08:06 PM; label: brush hair; status: Off
app	Type: App   clock
book	Type: Book
date time	Type: DateTime   1   1   2021
email address	Type: EmailAddress   membership@ipsa.org
flight number	Type: FlightNumber
general text	Type: GeneralText
home device	Type: UserEntity   heater
home room	Type: UserEntity   Db Bedroom
local business	Type: LocalBusiness   PostalAddress: 15 Broad St, Albany 31701   Ameris Bank   list_position: 13
media album	Type: MediaItem   MediaType: MediaType_Album   Mellon Collie
package	Type: Package
painting	Type: Painting
person	Type: Person   Sebastian
phone number	Type: PhoneNumber   955 545 060
photo	Type: Photo
physical address	Type: PostalAddress   GeographicArea: 814 Elmwood Ave, NY, 14222
plant animal	Type: PlantAnimal
setting	Type: Setting   dark mode
tracking number	Type: TrackingNumber
url	Type: Uri   NY.gov

## B Entity Representations

In Table 9, we show some examples of various domains and their representations, as fed into the LLM.

Table 10: Sample input with single ground truth

---

**User Request:** Call the one on Rainbow St.

---

**Entities Shown to User:**

1. Type: Local Business | Name: Walgreens | Address: 225 Rainbow St, San Jose CA 94088
2. Type: Local Business | Name: CVS | Address: 105 E El Camino Real, Sunnyvale, CA 94087
3. Type: Local Business | Name: Qwark | Address: 1287 Hammerwood Ave, Sunnyvale, CA 94089

---

**Ground Truth:** 1

---

Table 11: Sample input with multiple ground truths

---

**User Request:** Save the address.

---

**Entities Shown to User:**

1. Type: Postal Address | Value: 225 Rainbow St, San Jose CA 94088
2. Type: Email Address | Value: contactus@cv.com
3. Type: URL | Value: cvspharmacies.com/usa

---

**Ground Truth:** 1, 2, 3

---

## C Sample Inputs

In this section, we show examples of how inputs into the model have been encoded, in the form of a visual representation in Tables 10, 11 and 12.

Table 12: User Request for Alarms

---

**User Request:** Switch off the one reminding me to pick up didi.

---

**Entities Shown to User:**

1. Type: Alarm | open laptop
2. Type: Alarm | text Lauren to shower
3. Type: Alarm | pick up didi
4. Type: Alarm | forget this

---

**Ground Truth:** 3

---

# Dialog Flow Induction for Constraining LLM-Based Chatbots

Stuti Agrawal\*, Pranav Pillai\*, Nishi Uppuluri\*, Revanth Gangi Reddy, Zoey Li,  
Gokhan Tur, Dilek Hakkani-Tur, Heng Ji

University of Illinois Urbana-Champaign

{stutia3, nu4, ppillai3, revanth3, shal2}@illinois.edu

{gokhan, dilek, jih}@illinois.edu

## Abstract

LLM-driven dialog systems are used in a diverse set of applications, ranging from healthcare to customer service. However, given their generalization capability, it is difficult to ensure that these chatbots stay within the boundaries of the specialized domains, potentially resulting in inaccurate information and irrelevant responses. This paper introduces an unsupervised approach for automatically inducing domain-specific dialog flows that can be used to constrain LLM-based chatbots. We introduce two variants of dialog flow based on the availability of in-domain conversation instances. Through human and automatic evaluation over various dialog domains, we demonstrate that our high-quality data-guided dialog flows<sup>1</sup> achieve better domain coverage, thereby overcoming the need for extensive manual crafting of such flows.

## 1 Introduction

The widespread use of Large Language Models (LLMs) (OpenAI et al., 2023) for chatbots, highlighted by their human-like conversational abilities across many topics, faces challenges in specialized domains due to their tendency to go off-topic. This generalization capability, while a strength, necessitates the development of more effective control mechanisms to ensure chatbots remain within the desired domain of conversation, especially in specialized fields such as healthcare or legal advice. Controlling LLM-based chatbots can be effectively managed through dialog flows or schemas<sup>2</sup> (Bohus and Rudnicky, 2009; Mosig et al., 2020), which structure conversations along predefined paths of dialog actions, acting as directed graphs where nodes represent actions by the user or bot, and

\*denotes equal contribution

<sup>1</sup>Code is available at <https://github.com/gangiswag/dialog-flows>

<sup>2</sup>We use the terms *flows* and *schemas* interchangeably. Our definition of dialog schemas follows Mosig et al. (2020) to be analogous to task specifications, different from task slots.

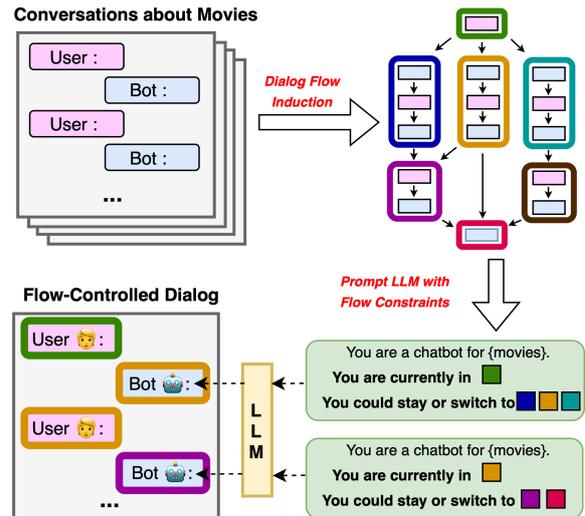


Figure 1: Figure demonstrating how automatically induced domain-specific dialog flows can be used to constrain chatbots to produce domain-focused responses.

edges are the transitions between actions. This structure helps steer the conversation, keeping it within relevant topics, and also enables chatbots to adapt to new tasks or domains without prior training (Zhao et al., 2023).

However, the construction of precise dialog flows is challenging (Huang et al., 2020), given the diversity of dialog in different domains. The most prevalent approaches (Mehri and Eskenazi, 2021; Zhao et al., 2023) use schemas that are carefully handcrafted by the dialog system developers. The design of dialog schemas thus has significant manual overhead for developers, resulting in scalability and coverage limitations (Zhang et al., 2020).

This paper introduces an unsupervised method to generate domain-specific dialog flows, exploiting GPT-4’s knowledge to systematically create detailed dialog flows reflecting conversational patterns in various domains. We begin by prompting GPT-4 to produce a structured representation of dialog interactions between users and bots, and then further refine this through self-reflective feedback

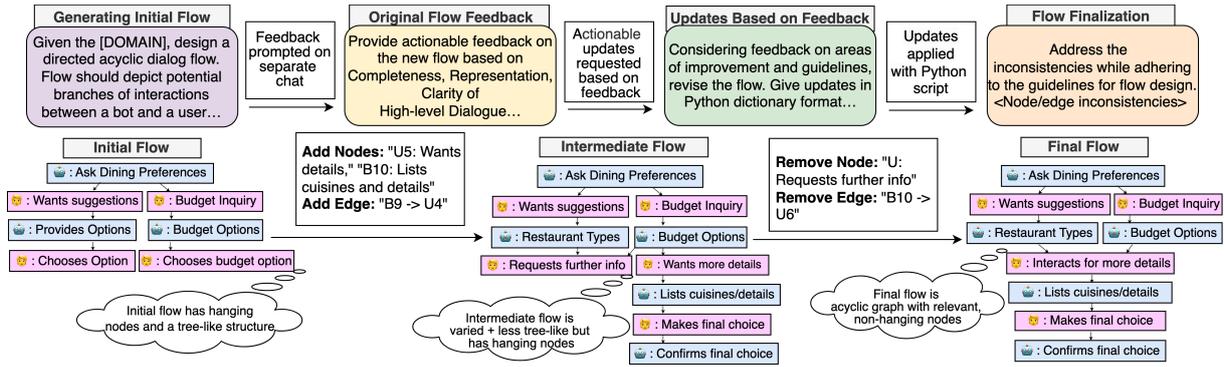


Figure 2: Figure showing the process for intrinsic flow induction. An initial flow is first generation which is further refined with feedback, update, and clean-up stages. Detailed prompts for each stage are provided in the appendix.

based on a set of predefined criteria (see figure 2).

Further, when we have domain-specific conversations, our approach automatically identifies distinct user and bot dialog actions within these conversations (see figure 3). These dialog actions, along with selected conversations that exemplify each action, are used to condition the GPT-4 prompt to ensure the dialog flows are grounded using actual domain instances. This approach enables the automated creation of structured dialog flows, facilitating the development of effective domain-specific chatbots that adhere to their domain’s conversational boundaries. Our main contributions are:

- This paper introduces an approach for automatically constructing dialog flows for various domains in an unsupervised manner.
- The proposed method uses a multi-step framework, that can further leverage domain-specific dialog instances, leading to a graph-like flow illustrating the structure of conversations in the domain.

## 2 Dialog Flow Induction

A dialog flow is a flowchart comprising nodes which can be a user or bot dialog action, and edges that denote logical flow or transitions between these actions. Dialog flows are tailored to different domains. Figure 2 shows an excerpt of a dialog flow, with more detailed examples in the appendix. In this section, we detail our approach for automatically inducing the dialog flow for a given conversation domain. Specifically, we induce two variants of dialog flows, namely *intrinsic flows* (in §2.1) or *data-guided flows* (in §2.2) depending on whether sample conversations in the domain are available.

### 2.1 Intrinsic Dialog Flow

When domain-specific conversation data is unavailable, we propose to induce dialog flows using the *intrinsic* domain-related knowledge of LLMs and their understanding of conversational principles. Our intrinsic flow induction process starts with GPT-4 creating an initial flow based on the domain’s name. Next, GPT-4 self-evaluates the flow based on predetermined guidelines, to provide concrete actionable feedback for improvement. Using this feedback, GPT-4 then suggests a set of edits, which are automatically applied to the initial flow. Finally, automated checks are run to identify inconsistencies in the flow, which GPT-4 then handles in the end clean-up stage. Figure 2 shows the overall intrinsic flow induction process, with more details on each step provided below.

**Initial Flow Generation:** The flow induction starts with prompting GPT-4 with a specific generation prompt to create a dialog flow, as shown in Figure 2. Along with the domain name, the prompt includes details on the intended structure of the dialog flow. After the initial flow is generated, it undergoes further refinement as detailed next.

**Flow Feedback and Updates:** The initial flow often suffers from low coverage along with ambiguous or repetitive action labels for bot and user nodes. We address these by leveraging GPT-4 for self-assessment (Bai et al., 2022) and refining the dialog flow based on the feedback. The refinement process starts by obtaining GPT-4 feedback based on the following aspects:

- **Representativeness:** Both the bot and user actions should be relevant to the domain, and should not be vague or generic.

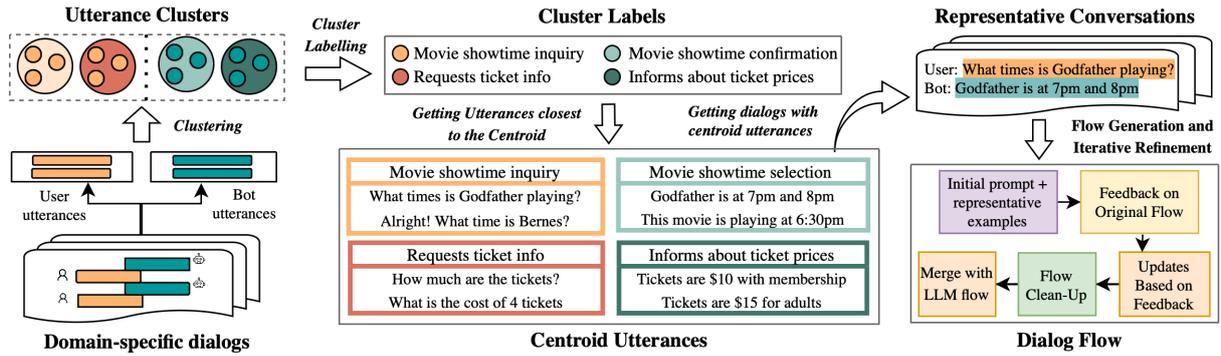


Figure 3: Figure showing the methodology for inducing dialog flows using a data-guided approach. Representative examples from the domain conversation instances are used to condition the GPT-4 prompts.

- **Coverage:** Ensuring the flow captures a broad range of conversational possibilities relevant to the domain.
- **Clarity of Dialog Action:** Each node should reflect a clear and meaningful dialog action.
- **Optimality:** Eliminate redundancy, ensuring no nodes depict overlapping dialog actions.

Based on the shortcomings identified by the self-reflective feedback, GPT-4 is then prompted to output a set of concrete updates to be made to the flow, which can include nodes or edges to add, remove, or edit. To control for the extent to which the flow changes, the updates are performed with an automated Python script rather than directly prompting GPT-4 to apply the updates<sup>3</sup>.

**Flow Finalization:** Finally, the dialog flow undergoes a clean-up stage where trivial inconsistencies, such as dangling non-terminal nodes, bot-bot or user-user connections, are identified. These are passed as input to GPT-4 along with a final prompt, to ensure the flow is structurally correct.

## 2.2 Data-Guided Dialog Flow

The intrinsic dialog flow induction approach, while expansive in its scope, relies predominantly on the model’s inherent knowledge of the typical interactions and transitions that could occur within the specified conversation domain. However, when dialog instances within the given domain are provided, the intrinsic flow can be updated to include actual conversational patterns. We call this approach *data-guided* flow induction, which aims to mirror real-world dialog dynamics. Specifically, the approach conditions the GPT-4 flow generation prompt with

<sup>3</sup>We hypothesize that this provides the ability to heuristically control different aspects of the dialog flow, such as depth, breath, density of edges, etc.

representative examples in the form of action labels and sample conversations for the domain, which help ground the flow to real-life conversation data. Figure 3 gives an overview of data-guided flow induction process, with more details provided below.

**Identifying Representative Examples:** Given dialog instances for a domain, the following steps identify the user and bot actions, along with sample conversations that are representative of the domain.

- **Clustering and Labeling:** The user and bot utterances from dialogs in the domain are clustered separately using SentenceBert (Reimers and Gurevych, 2019) embeddings. Next, GPT-4 is prompted to label each cluster with a dialog action by providing it with the utterances closest to each centroid.
- **Cluster Merging:** Next, we merge clusters that exhibit significant overlaps in terms of action intent, based on the cosine similarity between the labels. This reduces the redundancy in the action labels by grouping clusters with similar actions.
- **Picking sample conversations:** Finally, the conversations that include utterances corresponding to the cluster centroids are picked as the representative dialog instances to include in the GPT-4 prompt for flow generation. This ensures that the conversations encompass a wide spectrum of dialog actions and user intents specific to the domain.

**Flow Generation:** As shown in Figure 3, the flow induction follows a similar generation process as the intrinsic dialog flow. Firstly, the representative action labels and sample conversations for the domain are included in the initial flow generation prompt. Next, the feedback, update, and clean-up steps are applied to result in a dialog flow.

**Merging with Intrinsic Flow:** The intrinsic flow approach creates broad, expansive dialog flows, but can still fall short of reflecting domain-specific patterns from real-world conversations. On the other hand, solely relying on the domain dialog instances can hurt extensiveness, as they can have limited variability. Hence, we adopt a hybrid approach for the data-guided flow by merging the intrinsic flow with the flow induced solely from domain-specific data. This capitalizes on the extensive scope of the intrinsic flow with the detailed focus from domain data. This merging step is achieved by prompting GPT-4 to identify and retain distinctive features from the intrinsic flow, while removing redundant elements. We call this final flow, *data-guided* flow.

### 3 Experiments

We perform both human and automatic evaluations to assess the induced dialog flows.

#### 3.1 Datasets

Open-domain dialog can involve a single conversation touching upon different domains, such as movies, sports, music, etc. Hence, for simplicity, we consider domains from task-oriented dialog in our experimental settings, wherein the domains are distinct and correspond to the end user task, such as movie tickets, flight booking, restaurant reservations, etc. We consider a dialogs across various task-oriented domains, comprising 24 domains<sup>4</sup> from MetaLWoz (Shalyminov et al., 2019) and 5 domains from MultiWOZ (Budzianowski et al., 2018). For the data-guided flow induction, for each domain, we utilized 80% of the data as domain-specific instances available for training, with the remaining 20% reserved for evaluating coverage of the bot-bot transitions (described later in §3.3).

#### 3.2 Human Evaluation of Flow Quality

The evaluators (five undergraduate computer science students) were tasked with examining data-guided and intrinsic flows across the 24 different domains from MetaLwoz. The evaluators were given detailed guidelines (provided in the appendix), and were instructed to assess each flow on a scale of 1 to 5 for *domain coverage*, *conclusiveness* and *coherence*.

Table 1 shows numbers from human evaluation of the data-driven and intrinsic dialog flows. The

<sup>4</sup>We excluded domains that had ambiguous or generic names, such as Play Times, Catalogue, Agreement Bot, etc.

	Intrinsic	Data-driven
Domain Coverage	90.7	<b>93.0</b>
Conclusiveness	<b>87.8</b>	87.7
Coherence	84.5	<b>84.8</b>

Table 1: Results from human evaluation (in %) of different aspects of the induced dialog flows

Dataset	Intrinsic	Data-driven
MetaLWoz	31.6	<b>33.1</b>
MultiWOZ	39.9	<b>43.0</b>

Table 2: Bot-Bot transition coverage (in %) for the proposed variants of dialog flows on the MetaLWoz (Shalyminov et al., 2019) and MultiWOZ (Budzianowski et al., 2018) datasets. Detailed domain-wise numbers are provided in Table 3 in the appendix.

numbers (expanded to a scale of 20-100) are averaged over all the domains, with flows for each domain being annotated by 5 evaluators. We can see that the data-driven flow, on account of leveraging domain-specific dialog instances, improves over the intrinsic flow on domain coverage. Further, both dialog flows have similarly high scores for conclusiveness and coherence, implying our unsupervised approach, by leveraging GPT-4, can automatically induce high-quality dialog flows. We employed Randolph’s kappa to evaluate the multi-rater agreement. Our findings revealed a kappa value of 0.32, indicating a fair level of agreement across the board. Specifically, the domain coverage metric exhibited the highest kappa value of 0.46, signifying moderate agreement.

#### 3.3 Automatic Evaluation of Flow Coverage

Next, we automatically evaluated the domain coverage of different dialog flows, by measuring the coverage on capturing bot-to-bot transitions within the domain conversations in the test set. We leveraged Mistral-7B-Instruct (Jiang et al., 2023) to classify bot utterances into the most appropriate node in the dialog flow. We then examined whether the next bot utterance mapped to the directly succeeding node in the dialog flow. Essentially, this metric measures the percentage of bot-bot transitions in domain conversations that conform to the given dialog flow. Table 2 shows numbers for automatic coverage evaluation. We can see that the data-driven dialog flow has better coverage of the domain’s bot-bot transitions.

## 4 Conclusion and Future Work

We introduce a novel method for developing dialog flows that reflect the combined intrinsic knowledge of LLMs and existing domain-relevant dialogs. Our data-driven dialog flow approach achieves better domain coverage than the intrinsic flow approach across human and automatic evaluations. Our paper outlines a blueprint (in Figure 1) for integrating the generated dialog flows into LLM-based chatbots, with a primary focus on the methodologies for dialog flow generation. We believe these dialog flows can be a springboard for future interactive dialog systems that maintain a natural conversation flow within the domain.

### Limitations

In this study, our experimentation was confined to task-oriented dialogs, encompassing a relatively narrow spectrum of dialog flows. This specialization may limit the applicability of our findings to dialog domains characterized by a broader array of tasks and more open-ended dialogues. Additionally, our methodology relies solely on unsupervised clustering techniques, bypassing datasets that are annotated with slot values and user intents, which could potentially enhance dialog flow induction. Furthermore, we have not extended our research to test the performance of chatbots constrained by the dialog schemas we developed. Therefore, the efficacy of these schemas in practical chatbot applications remains an area for future investigation.

### Acknowledgment

We would like to thank the CS STARS program at UIUC for supporting Stuti and Nishi. We are grateful to members of the BlenderNLP group for their valuable comments and feedback. This research is based on work supported by U.S. DARPA KAIROS Program No. FA8750-19-2-1004 and U.S. DARPA INCAS Program No. HR001121C0165. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of DARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

## References

- 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*.
- Dan Bohus and Alexander I Rudnicky. 2009. The ravenclaw dialog management framework: Architecture and systems. *Computer Speech & Language*, 23(3):332–361.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gasic. 2018. Multiwoz-a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026.
- Minlie Huang, Xiaoyan Zhu, and Jianfeng Gao. 2020. Challenges in building intelligent open-domain dialog systems. *ACM Transactions on Information Systems (TOIS)*, 38(3):1–32.
- 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*.
- Shikib Mehri and Maxine Eskenazi. 2021. Schema-guided paradigm for zero-shot dialog. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 499–508.
- Johannes E. M. Mosig, Shikib Mehri, and Thomas Kober. 2020. [Star: A schema-guided dialog dataset for transfer learning](#). *ArXiv*, abs/2010.11853.
- OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madeleine 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, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondrasiuk, 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, Rajeef 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 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 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. 2023. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- 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.
- Igor Shalyminov, Sungjin Lee, Arash Eshghi, and Oliver Lemon. 2019. Few-shot dialogue generation without annotated data: A transfer learning approach. *arXiv preprint arXiv:1908.05854*.
- Zheng Zhang, Ryuichi Takanobu, Qi Zhu, MinLie Huang, and XiaoYan Zhu. 2020. Recent advances and challenges in task-oriented dialog systems. *Science China Technological Sciences*, 63(10):2011–2027.
- Jeffrey Zhao, Yuan Cao, Raghav Gupta, Harrison Lee, Abhinav Rastogi, Mingqiu Wang, Hagen Soltau, Izhak Shafran, and Yonghui Wu. 2023. [AnyTOD: A programmable task-oriented dialog system](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16189–16204, Singapore. Association for Computational Linguistics.

## A Appendix

<b>MetaLWoz</b>	<b>Intrinsic</b>	<b>Data-driven</b>
Alarm set	32.9	<b>42.2</b>
Apartment finder	30.9	<b>45.2</b>
Bank bot	<b>34.2</b>	30.8
Bus schedule	<b>37.2</b>	14.4
City info	29.4	<b>33.4</b>
Edit playlist	<b>44.2</b>	39.4
Event reserve	28.8	<b>30.5</b>
Library Request	<b>35.7</b>	30.1
Movie listings	30.7	<b>34.4</b>
Music suggester	<b>34.0</b>	25.3
Name suggester	<b>43.2</b>	16.7
Order pizza	31.6	<b>36.1</b>
Pet advice	<b>33.8</b>	31.7
Phone plan	31.6	<b>37.8</b>
Restaurant picker	<b>29.4</b>	29.2
Scam lookup	22.6	<b>31.2</b>
Shopping	17.0	<b>22.9</b>
Ski Bot	27.2	<b>32.2</b>
Sports info	36.6	<b>37.1</b>
Store details	<b>35.7</b>	32.4
Update calendar	<b>38.4</b>	28.8
Update contact	<b>32.5</b>	30.3
Weather check	<b>36.1</b>	29.5
Wedding planner	17.0	<b>24.2</b>
<b>Average</b>	31.6	<b>33.1</b>

<b>MultiWOZ</b>	<b>Intrinsic</b>	<b>Data-driven</b>
Restaurant	31.0	<b>43.9</b>
Hotel	43.2	43.2
Attractions	43.3	<b>53.3</b>
Taxi	<b>75.3</b>	50.5
Train	6.9	<b>24.1</b>
<b>Average</b>	39.9	<b>43.0</b>

Table 3: Bot-Bot transition coverage (in %) for the proposed variants of dialog flows when measured on various domains in the MetalWoz (Shalyminov et al., 2019) and MultiWOZ (Budzianowski et al., 2018) datasets.

<b>MetaLWoz</b>	<b>Train</b>	<b>Test</b>
Alarm set	1345	336
Apartment finder	399	100
Bank bot	294	73
Bus schedule	718	180
City info	772	193
Edit playlist	459	115
Event reserve	431	108
Library request	1071	268
Movie listings	486	121
Music suggester	356	89
Name suggester	399	100
Order pizza	462	115
Pet advice	341	85
Phone plan	397	99
Restaurant picker	428	107
Scam lookup	1326	332
Shopping	722	181
Ski bot	486	121
Sports info	449	112
Store details	590	147
Update calendar	1593	398
Update contact	522	131
Weather check	441	110
Wedding planner	408	102

Table 4: Statistics of dialogs in various domains in the MetalWoz (Shalyminov et al., 2019) dataset.

### Intrinsic Flow Initial Prompt

Given the context of [DOMAIN], design a directed acyclic dialog flow suitable for visualization with mermaid.js. This flow should depict the nuances and potential branches of interactions between a bot and a user. Please adhere to the following guidelines:

**Nodes Definition:** Use distinct nodes to represent the bot ("B") and the user ("U").

**High-Level Dialog Action:** Each node should encapsulate that segment's core sentiment or function in the conversation, relevant to [DOMAIN]. It should be a label for the node representing a high-level dialogue action and not just the dialogue.

**Flow & Directionality:** Create directed connections between nodes to represent the progression of the conversation. The dialogue should flow from one node to potentially multiple nodes, allowing for various conversational turns.

**Diverse Conversational Possibilities:** Ensure that bot nodes can lead to multiple user nodes and vice versa. This should account for various user responses or bot prompts, showcasing the range of interactions possible within [DOMAIN].

**Acyclic Structure:** The dialog flow must not have loops or cyclic pathways. If a similar action or sentiment arises later in the conversation, introduce a new node to represent it, rather than looping back to an earlier node.

**Mermaid.js Compatibility:** Ensure that the constructed flow is adherent to mermaid.js graph notation, guaranteeing its seamless rendering.

Considering the guidelines, craft a dialogue flow focused on [DOMAIN]. The bot always begins by greeting the user and asking for what they want. The graph should be connected. The bot and user nodes should be in different colors. A bot node is only followed by user nodes and user nodes are by bot nodes.

### Intrinsic Flow Feedback Prompt

Based on the below evaluation criteria, suggest some improvements and provide concise + actionable feedback on the flow just generated:

**Optimality:** Check for redundancy. Ensure that nodes aren't replicating the same or very similar dialog actions, even if they arise at different points in the conversation

**Clarity of High-Level Dialog Action:** For every node, evaluate if the high-level dialog action is clear and meaningful. Avoid nodes that are vague or overly complex. Can someone unfamiliar with the domain understand the flow and interactions by looking at the flow?

**Extensiveness:** Does the flow account for diverse conversational possibilities? Are all the nodes interconnected to the graph? Does the flow cover all major high level topics and interactions within the given domain?

**Representativeness of the Domain:** Bot Nodes (B): Do the bot nodes represent clear and unambiguous actions? Are they too broad or too specific? User Nodes (U): Do user nodes accurately capture an adequate range of potential user responses and inquiries relevant to the domain?

### Intrinsic Flow Update Prompt

Taking into consideration the feedback and the original design guidelines - keep it in directed acyclic graph structure and make sure all new components are labeled and connected to the graph correctly- revise the flow. Ensure your revised flow addresses the identified areas of improvement while still adhering to the primary instructions for flow construction. Make sure to account for all new nodes including merged nodes and their labels/colors. Make sure all user nodes connect with bot nodes and bot nodes are the end of the conversations. Give your updates in the below format:

```
'split_nodes':  
# 'NodeToSplit': ['NewNode1', 'NewNode2', ...],  
  
'add_nodes':  
# 'NodeToAdd': 'Label',  
  
'remove_nodes':  
# 'NodeToRemove1', 'NodeToRemove2', ...
```

```
'relabel_nodes':  
# 'NodeToRelabel': 'NewLabel',  
  
'add_edges':  
# ('Start Node', 'End Node'),  
  
'remove_edges':  
# ('Start Node', 'End Node'),
```

### Intrinsic Flow Finalization Prompt

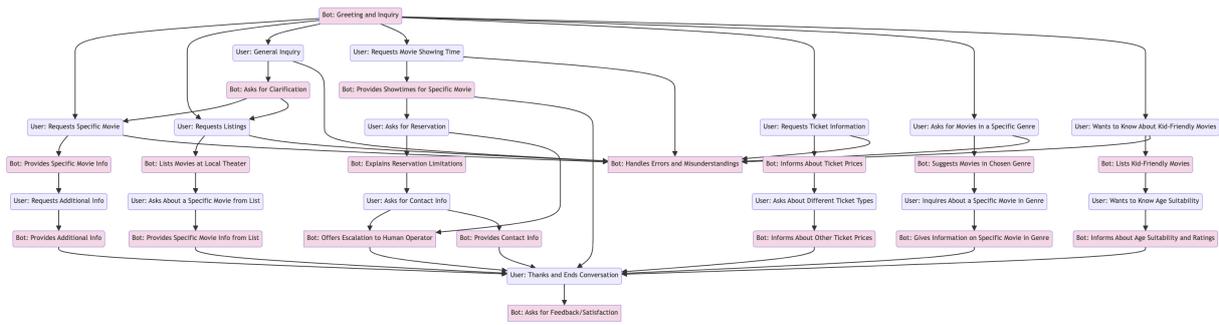
Clean up the flow to create a final flow. Ensure your revised flow addresses the identified areas of improvement while still adhering to the primary instructions for flow construction. Get rid of hanging/loose user nodes (user nodes with no output), have graph in directed acyclic structure, bot nodes shouldn't be connected to other bot nodes, and user nodes shouldn't be connected to other user nodes. All nodes should have input/output except begin and end nodes, one node shouldn't point to the another node more than once, and make sure all bot nodes are correctly colored.

### Intrinsic and Data-Guided Flows Merging Prompt

Given the two dialogue flows for the [DOMAIN] bot. One flow is LLM generated and the other is from data examples, merge all the unique elements of the two flows and do not duplicate similar elements. Merge the two flows based on the following design guidelines:

< Design guidelines >

Figure 4: Figure showing prompts for intrinsic and data-guided dialog flow generation.

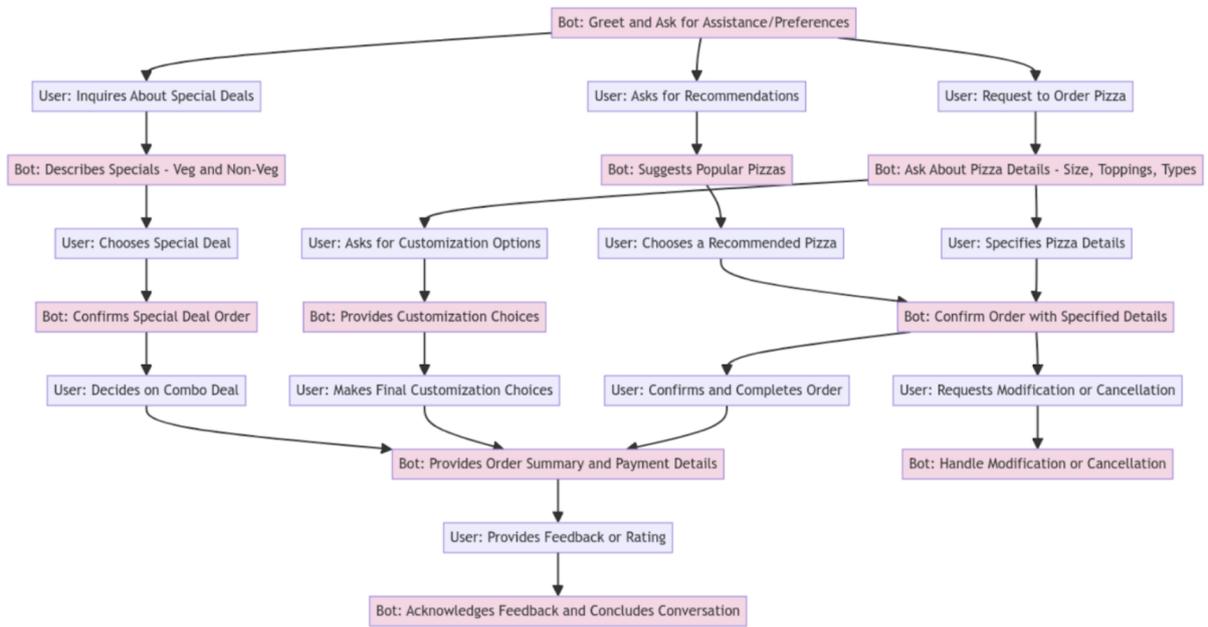


(a)

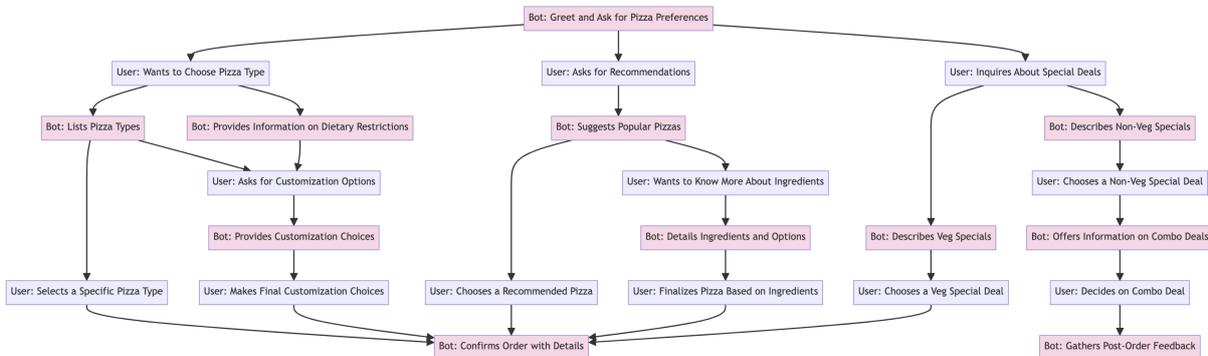


(b)

Figure 5: Data-driven (a) and Intrinsic (b) flows for the movie listings domain from MetalWoz.

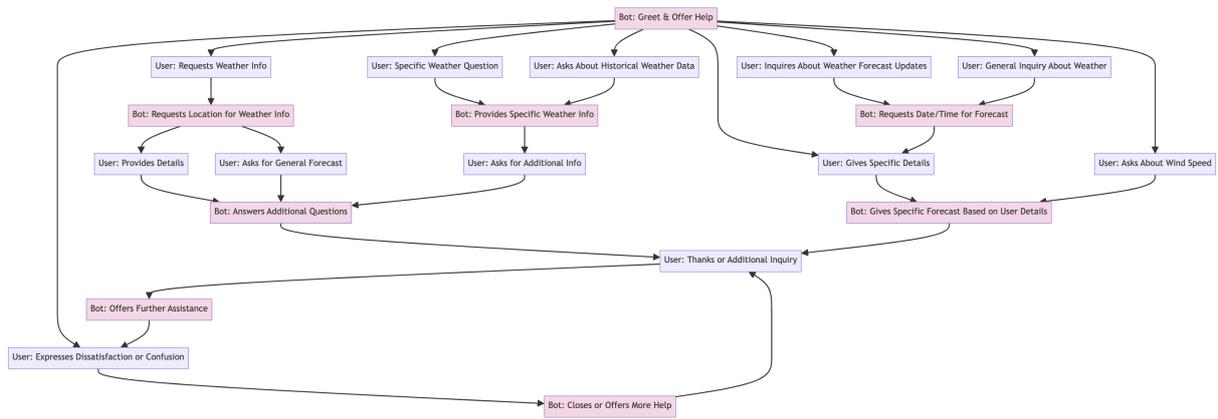


(a)

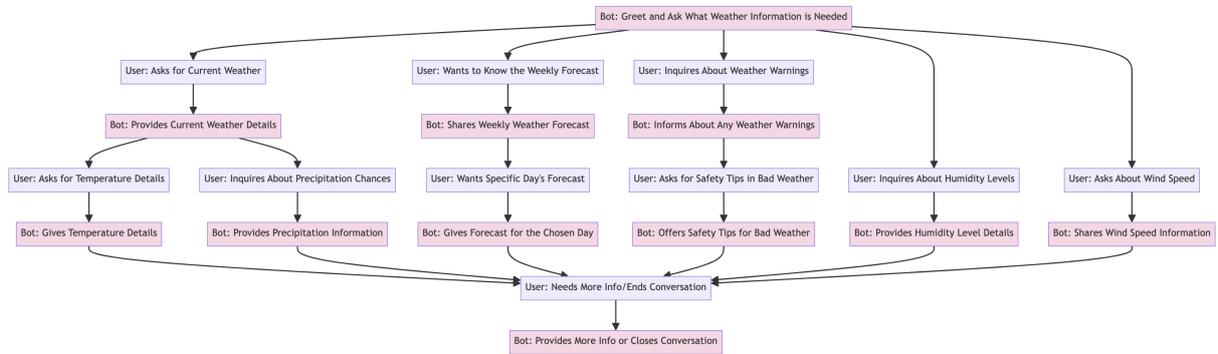


(b)

Figure 6: Data-driven (a) and Intrinsic (b) flows for the order pizza domain from MetaLWoz.



(a)



(b)

Figure 7: Data-driven (a) and Intrinsic (b) flows for the order weather domain from MetaLWoz.

### **What is a task-oriented dialog flow?**

A dialog flow is like a roadmap for conversations between a user and a chatbot, outlining all the possible exchanges they can have. It guides the chatbot on how to respond to different user inputs, ensuring the conversation flows smoothly and logically. Dialog flows are composed of bot nodes, which correspond to bot actions, and user nodes, corresponding to user actions. In the below dialog flows, all bot nodes are shown in pink and all user nodes are shown in blue. Arrows indicate the flow of the conversation and the potential action(s) a bot or user could take.

### **Directions**

You will be provided 2 variants of dialog flows for each conversation domain. Keeping in mind how a regular task-oriented chatbot might work, rate each flow on a scale of 1-5 on each of the four metrics: *domain (or topic) coverage*, *conclusiveness*, and *coherence*. You don't have to explain your answers.

**Examples:** Please refer to the document [here](#) for examples of a few flows along with some sample ratings.

**Note:** The annotator is recommended to come up with a rough working on what actions they believe a task-oriented dialog system for a given domain should “have”, even before looking at the provided flows. This will help judge better when evaluating the provided flows for domain coverage.

### **General Rubric**

#### **Domain (or Topic) Coverage (1-5):**

- Score 1: The flow is generic and barely covers any relevant aspects of the domain.
- Score 3: The flow covers key aspects of the domain but still misses some of them
- Score 5: The flow comprehensively covers all major and relevant aspects of the domain or topic, providing a thorough and detailed exploration.

#### **Conclusiveness (1-5):**

- Score 1: Conversations in the flow often end abruptly or leave the main query unresolved, leading to dissatisfaction.
- Score 3: Conversations tend to lead toward a resolution, but some paths may still end with questions or lack finality
- Score 5: Each conversation path leads to a clear and satisfactory conclusion or task completion, ensuring user queries are fully addressed.

#### **Coherence (1-5):**

- Score 1: The flow of conversation is disjointed or confusing, with many leaps and complex connections that disrupt understanding.
- Score 3: The conversation flow is natural for the most part but does have some non-logical paths or jumps
- Score 5: The conversation flows logically and naturally from one point to the next, with all parts making sense in the context and enhancing comprehension.

Figure 8: Evaluation Instructions for Human Annotators

# Knowledge-Grounded Dialogue Act Transfer using Prompt-Based Learning for Controllable Open-Domain NLG

Alain Vazquez Risco and Angela Ramirez and Neha Pullabhotla  
and Nan Qiang and Haoran Zhang and Marilyn Walker and M. Inés Torres  
University of the Basque Country and University of California Santa Cruz  
alain.vazquez@ehu.eus, aramir62@ucsc.edu, npullabh@ucsc.edu,  
nqiang@ucsc.edu, hzhan264@ucsc.edu, mawalker@ucsc.edu, manes.torres@ehu.eus

## Abstract

Open domain spoken dialogue systems need to controllably generate many different dialogue acts (DAs) to allow Natural Language Generation (NLG) to create interesting and engaging conversational interactions with users. We aim to create an NLG engine that can produce a variety of DAs that make substantive knowledge-grounded contributions to a conversation. Training such an NLG typically requires dialogue corpora that are labelled for DAs, which are expensive to produce and vulnerable to quality issues. Here, we present a prompt-based learning approach to transfer DAs from one domain, video games, to 7 new domains. For each novel domain, we first crawl WikiData to create Meaning Representations that systematically vary both the number of attributes and hops on the WikiData Knowledge Graph. The proposed method involves a self-training step to create prompt examples for each domain followed by an overgeneration and ranking step. The result is a novel, high-quality dataset, Wiki-Dialogue, of 71K knowledge-grounded utterances, covering 9 DAs and the Art, Movies, Music, Sports, TV, Animal, and Boardgames domains, whose combined DA and semantic accuracy is 89%. We assess the corpus quality using both automatic and human evaluations and find it high. The corpus is found to be safe, lexically rich, and large in vocabulary, when compared to similar datasets.

## 1 Introduction

Open domain spoken dialogue systems need to be able to controllably generate many different dialogue acts (DAs) in order to create interesting and engaging conversational interactions with users. For example, they should be able to ask questions of different types, inform the user of facts and express opinions, make recommendations and suggestions, and confirm what the user said. Moreover, using knowledge to ground DAs supports taking the initiative to drive the conversation forward, and

has been shown to help avoid hallucinations in generated outputs (Dziri et al., 2021; Gopalakrishnan et al., 2019; Chi et al., 2022).

We aim to create a Natural Language Generation (NLG) engine that can produce a variety of DAs and make substantive knowledge-grounded contributions to a conversation. A synthetic example of a conversation about music with a variety of DAs is shown in Table 1. This example is shown because it uses the nine different DAs from the ViGGO video games corpus (Juraska et al., 2019, 2021), with sets of triples from the Wikidata Knowledge Graph (KG) for knowledge grounding.

Dialogues like this require a semantically-controlled NLG that also controls the expression of DAs. In order to create such an NLG, training data consisting of dialogue corpora with utterances labeled with DAs, such as Multi-Woz (Budzianowski and Vulić, 2019), or a parallel corpus of domain-specific DA representations and reference utterances, such as ViGGO or RNNLG (Juraska et al., 2019; Wen et al., 2015), is typically needed. Such training data is typically collected via crowdsourcing, making it expensive to produce, and vulnerable to quality issues (Qian et al., 2021; Dušek et al., 2019). As a result, existing corpora primarily use the INFORM DA (Gardent et al., 2017; Novikova et al., 2017; Lebret et al., 2016; Nan et al., 2021; Parikh et al., 2020) or cover a limited number of topics (Wen et al., 2015; Budzianowski and Vulić, 2019; Juraska et al., 2019; Rastogi et al., 2020).

However, recent work using prompt-based learning (PBL) with LLMs shows that PBL can be used to create high quality dialogic datasets (Kim et al., 2023; Chen et al., 2023; Ramirez et al., 2023). For example, Ramirez et al. used PBL to enlarge the DA labelled ViGGO video games and RNN Laptop and TV corpora (Juraska et al., 2019; Wen et al., 2015). Here we extend their overgenerate-and-rank NLG method to transfer the DAs from the ViGGO corpus onto a set of 6 WikiData domains and a

Dialogue Act Representation	System Output
GREETING	System: Hi, it's so good to talk to you again! How are you doing today? User: I'm doing good, I want to talk about music
REQUEST-ATTRIBUTE (song genre.)	Sys: Sure, What kind of music do you like to listen to? User: I enjoy pop music
REQUEST (songs specifier catchy. songs genres pop.)	Sys: Do you know of any pop songs that are catchy? User: Yes, never gonna give you up
CONFIRM (Never Gonna Give You Up performer Rick Astley. Never Gonna Give You Up domain songs.)	Sys: Is it Rick Astley's song Never Gonna Give You Up that you're referring to? User: Yes
GIVE OPINION (Never Gonna Give You Up rating excellent. Never Gonna Give You Up performer Rick Astley. Never Gonna Give You Up publication year 1987. Never Gonna Give You Up domain songs.)	Sys: Rick Astley did a nice job with the single 'Never Gonna Give You Up' from 1987. It's quite good! User: That sounds great, I totally agree with you.
SUGGEST (Never Gonna Give You Up genre pop music. All I see genre pop music. All I see performer Kylie Minogue.)	Sys: Do you also like pop songs by Kylie Minogue, for example, All I See? User: No, I don't know about that song
INFORM (All I see rating excellent. All I see from album X. All I see performer Kylie Minogue. All I see domain songs.)	Sys: I must say, Kylie Minogue really nailed it with their performance of All I See from the album X. It's a fantastic track that I love listening to. User: I see. Are there any other songs that you can recommend
RECOMMEND (Kylie Minogue genre synth pop. Berlin genre synth pop. Berlin performer of Take My Breath Away.)	Sys: Speaking of synth-pop artists, have you heard the song Take My Breath Away by Berlin? User: Oh I love this artist and this song
REQUEST-EXPLANATION (songs rating excellent. songs performer Berlin.)	Sys: What do you think is it that makes the songs by Berlin so excellent? User: His songs always bring tears to my eyes

Table 1: Music conversation using a dialogue policy based on all the Dialogue Act (DA) types from ViGGO

database about Animals. Since there is no existing training data for these domains, we start with prompt examples from ViGGO, and add a self-training step. We show that self-training improves the results.

We first specify the constraints on different types of ViGGO DAs for each domain, then crawl Wiki-Data to create meaning representations (MRs) consisting of sets of KG triples for each DA type and domain. Example MRs can be seen on the left side of Table 1. We build on WebNLG’s approach for generating multihop KG Triples as the basis for some MRs, which then specify paths through the KG, e.g. talking about the spouse of a person’s sibling involves a path of two hops (Gardent et al., 2017). We systematically vary the complexity of the MRs in terms of number of attributes (up to 8) and number of hops through the KG (up to 3). To test generalizability, we then extend the method to create MRs for an Animals database whose relations and values are very different. To improve the quality of the corpus, we first overgenerate multiple outputs and then rank them, by automatically estimating DA and semantic accuracies and fluency at generation time. We then take the best rated outputs and repeat the procedure with a self-training step. The result is a novel, high-quality, synthetic dataset, Wiki-Dialogue, of knowledge-grounded DAs for the Art, Movies, Music, Sports, TV, Animal, and Boardgames domains, whose combined DA and semantic accuracy (PERFECT) is 89%. Our contributions include:

- Wiki-Dialogue: A new multi-domain dialogue

act and semantically-controlled corpus for the NLG community.

- A novel method that transfers DAs from one domain to another with prompt-based learning.
- A systematic analysis of methods for improving the quality of LLM generated corpora.

## 2 Related Work

Most knowledge-grounded dialogue corpora are based on crowdsourcing utterances matching different types of MRs. One motivation for generating corpora using LLMs is that crowdsourcing is expensive, and crowdsourced corpora can be noisy, requiring extensive filtering or additional annotation to ensure accuracy. For example, Dušek et al. states that up to 40% of the utterances in the E2E corpus either omitted information that was present in the MR or contained additional information (Dušek et al., 2019), while Qian et al. state that 70% of the dialogues in MultiWOZ contained dialogue state (semantic) annotation errors, which had to be corrected (Ye et al., 2022). Web-NLG (Gardent et al., 2017) also contained mismatches between the MRs and the crowdsourced utterances that have been corrected.<sup>1</sup> We show in Section 4 that there are fewer semantic errors in Wiki-Dialogue than in similar crowdsourced corpora.

The most similar corpus to Wiki-Dialogue is WebNLG, a multi-domain corpus that has been used for NLG challenges (Colin et al., 2016;

<sup>1</sup><https://gitlab.com/shimorina/webnlg-dataset>

Zhou and Lampouras, 2020; Ferreira et al., 2018). WebNLG is based on the DBPedia KG (Lehmann et al., 2015), and covers 19 domains, with utterances realizing the MRs collected via crowdsourcing. The English dataset contains about 17,000 triple sets and 45,000 crowdsourced texts. However, the only DA in WebNLG is INFORM, and dialogues consisting of only INFORM utterances are unnatural (See et al., 2019; Rastogi et al., 2020; Hedayatnia et al., 2020), while Wiki-Dialogue provides DA diversity with nine DAs (Juraska et al., 2019).

Previous work by Moon et al. (2019) created the OpenDialKG dataset by crowdsourcing human-human dialogues consisting of 15K utterances for the Music, Movies, Sports and Books domains. Each utterance corresponds to either a one or two-hop fact in the Freebase KG (Bast et al., 2014) and the KG paths are hand-annotated on each utterance. This corpus contains different types of DAs, but there are no DA annotations.

The Schema-Guided Dialogue (SGD) dataset consists of over 20k annotated multi-domain, task-oriented conversations between a human and a virtual assistant (Rastogi et al., 2020). These conversations target interactions with services, such as travel, spanning 20 domains. The dialogues are generated in two steps: first a simulator automatically generates a dialogue, given a task, in the form of a sequence of DAs whose semantic values are filled by queries to Freebase. The DAs used by the simulator are distinct for the system and the user, with 10 system DAs and 11 user DAs. Then crowdworkers are tasked with paraphrasing each dialogue act/MR combination in a dialogue to ensure naturalness of the utterance realizations.

The Topical Chat corpus consists of 235K utterances from 8 domains, but differs from other knowledge-grounded corpora in that the knowledge is represented by sentences which are automatically aligned with the dialogues. DAs were automatically labelled on this corpus using the 11 DAs from the ISO DA standard with an F1 of 0.54 (Hedayatnia et al., 2020; mez). The DA labels in Wiki-Dialogue are much less noisy: See Table 10.

Other knowledge-grounded NLG corpora with rich sets of DAs have focused on specific domains. The ViGGO Video Games corpus contains 9 DAs (Juraska et al., 2019), the RNNLG corpus encompasses 13 DAs for domains such as laptops, TVs, hotels, and restaurants (Wen et al., 2015), and MultiWOZ offers 34 task-oriented, domain-specific

DAs, for the restaurant, hotel, attraction, taxi, train, hospital and police domains (Eric et al., 2021).

Recent work by Wu et al. (2023) on controlling DAs in NLG for task-oriented dialogue introduces DiactTOD, a model based on learning latent DAs from pre-existing datasets, achieving state-of-the-art performance on MultiWOZ (Ye et al., 2022). However, the evaluation DiactTOD is based on a benchmark set and automatic evaluation metrics, which do not evaluate DA accuracy. We use a fine-tuned classifier that filters for DA accuracy as one step in the automatic ranking of possible responses. We also apply both automatic and human evaluation, showing that we achieve average DA accuracies of .98 for one-hop and .89 for multi-hop.

Other work on creating synthetic data for NLG has focused on creating whole dialogues or augmenting existing corpora (Kim et al., 2023; Xu et al., 2021; Chen et al., 2023). One of the main challenges with synthetic dialogue generation is producing high quality outputs without human evaluation. Here we show that the quality of synthetic dialogue data can be improved using a cycle of self-training, along with an overgenerate-and-rank step that uses a DA classifier and semantic accuracy estimates. We build on previous work by Ramirez et al. by using their definitional prompt style, ranking function, and DA tagger. We extend their approach by transferring the DAs used for controlled generation to 7 new domains and incorporating a self-training step required to bootstrap high quality generation outputs for completely novel domains.

### 3 Experimental Method

Figure 1 provides an overview of the experimental architecture. Our method consists of five steps:

- Specifying DA constraints;
- Creating DA specific MRs;
- Prompt Creation and LLM selection;
- Overgenerate and Rank from the MRs for one round using ViGGO examples in the prompts;
- Self-Training: Select new in-domain prompt examples and conduct a second round of overgenerate-and-rank.

#### 3.1 Specifying Dialogue Act Constraints

We use the DAs from the ViGGO corpus to enable more highly varied dialogue policies. One possible policy is illustrated by the conversation shown in Table 1. The utility of controlling DAs and being

DA	Num Rels	ViGGO		Athletes		Wild Animals	
		Mandatory Rels	Hops	Mandatory Rels	Hops	Mandatory Rels	Hops
INFORM	3-8	NAME, GENRES	1	NAME, SPORT	1-3	NAME, COMMON_CLASS	1
CONFIRM	2-3	NAME	1	NAME, SPORT	1	NAME, COMMON_CLASS	1
GIVE_OPINION	3-5	NAME, RATING	1	NAME, RATING, SPORT	1-3	NAME, RATING, COMMON_CLASS	1
RECOMMEND	2-3	NAME	1	NAME, SPORT	1-2	NAME, COMMON_CLASS	1
REQUEST	1-2	SPECIFIER	1	SPECIFIER	1	SPECIFIER	1
REQUEST_ATTRIBUTE	1	–	1	–	1	–	1
REQUEST_EXPLANATION	2-3	RATING	1	NAME, POPULARITY	1-2	POPULARITY	1
SUGGEST	2-3	NAME	1	NAME	1-2	NAME	1
VERIFY_ATTRIBUTE	3-4	NAME, RATING	1	NAME, RATING	1-2	NAME, POPULARITY	1

Table 2: Semantic Constraints on Dialogue Acts for VideoGames (ViGGO), Athletes and Wild Animals.

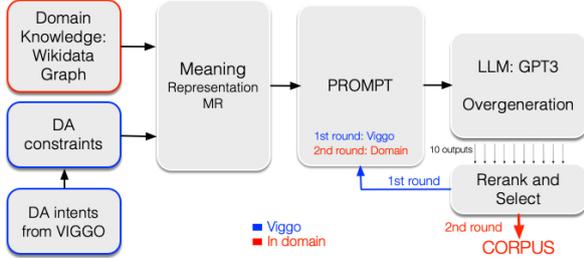


Figure 1: Experimental Architecture

able to produce different DAs is also illustrated by the utterances for the Art domain in Table 3, which demonstrate how the same MR can be realized very differently depending on the DA. A range of DAs allows a dialogue system to vary the dialogue policy in interesting ways (Juraska et al., 2021).

#### Attributes and Values

(NAME [Andromeda Chained to the Rocks], CREATOR [Rembrandt], GENRES [mythological painting, nude], INCEPTION [1630])

#### confirm

Is the painting you’re talking about **Andromeda Chained to the Rocks** by **Rembrandt**, a **mythological nude** piece from **1630**?

#### recommend

I recall you enjoy **mythological** paintings and **nudes**, so have you heard of **Rembrandt’s Andromeda Chained to the Rocks** from **1630**?

#### inform

**Andromeda Chained to the Rocks**, created by **Rembrandt** in **1630**, is a **mythological** painting featuring a **nude** figure.

Table 3: Sample dialogue acts (DAs) in the Art domain showing how the same attributes and values can be realized as different DAs.

The initial step of our method involves specifying semantic attribute constraints for each domain and entity and DA types. NUM RELS in Table 2 is the number of relations that can be included in a particular DA. As seen in Table 2, the CONFIRM DA (Row 2) should only include 2–3 relations, or it is difficult to understand, while the INFORM DA typically has more slots than other DAs.

We also specify mandatory and optional rela-

tions for each DA and entity type. MANDATORY RELS in Table 2 specifies relations that must be included for each DA type. The GIVE\_OPINION DA (Row 3) requires an ordinal attribute similar to RATING in order to orient the polarity of the opinion. For domains such as Athletes and Animals, that do not have directly such an attribute, so the number of Wikipedia page views is used to create an ordinal and equivalent POPULARITY attribute, with values ranging from LOW to HIGH. Potential values for the SPECIFIER attribute needed by the REQUEST DA must be provided for each entity type, e.g. the specifier CATCHY used in the REQUEST DA in the conversation in Table 1. The ViGGO columns in Table 2 show the mandatory relations that were based on ViGGO’s 14 video game attributes. Each DA also has optional relations that define the attributes that can be added to the mandatory ones when creating the MRs. They are shown for each entity type and domain in the corpus GitHub.<sup>2</sup>

For each domain, there are two or more entity types, e.g. the entity types for Movies are Actor and Movie; for Music they are Song, Album and Musician; for Sports they are Athlete and Team; and for Animals they are Wild Animals, Cats and Dogs. Each domain’s specific slots are provided on GitHub with the corpus description. For the entity types that are media, like Movies, Songs or TV Shows, the constraints are identical to those shown for ViGGO in Table 2.

However, other domains require different semantic constraints as illustrated in the Athletes and Wild Animals columns of Table 2, e.g. for Athletes, the attribute of genre doesn’t apply, but the sport that the Athlete plays serves a similar function. Similarly, Animals doesn’t have a genre, but common\_classes of animal, such as reptile, fish or mammal, are needed to specify the general type of the animal. For human entity types such as Actors, Musicians and Athletes the REQUEST\_EXPLANATION DA (Row 7 of Table 2) must include the name

<sup>2</sup><https://github.com/aramir62/Wiki-Dialogue>

slot: see the example conversation in Table 1. In addition, for both Athletes and Animals, the CONFIRM (Row 2), GIVE OPINION (Row 3), and RECOMMEND (Row 4) DAs require that the sport be mentioned for the athlete, and the common\_class mentioned for the animal.

### 3.2 Generating KG Triples from WikiData

After specifying the DA attributes, we generate KG triples adhering to DA constraints using WikiData queries for Art, Movies, Music, TV Shows, Boardgames, and Sports. WikiData offers detailed knowledge across many domains and shares canonical IDs with Wikipedia. To test generalization, we used API Ninjas<sup>3</sup> to create MRs for the Animal domain. For all entity types for each domain, we selected 60 entities: 30 popular and 30 lesser-known, based on Wikipedia page visits in the past 6 months. These entities become the starting nodes for all paths, and for multi-hop data, both incoming and outgoing relations are included.

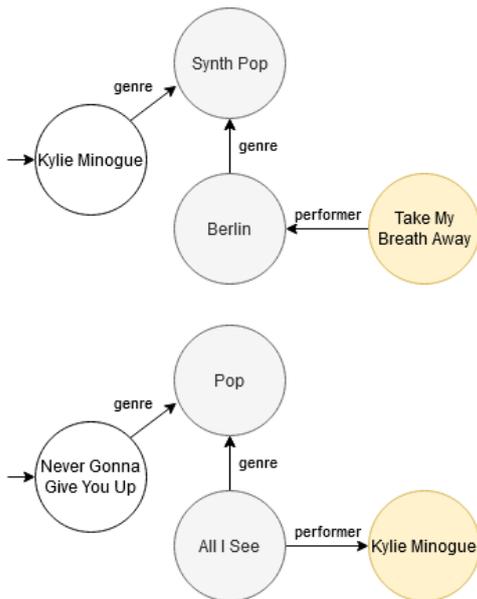


Figure 2: Subgraphs of Wikidata for the Recommend Dialogue Act and Suggest Dialogue Act in the music conversation in Table 1

For multi-hop paths, there are constraints on the maximum number of hops for each DA. This is necessary since some DAs such as CONFIRM should not realize multiple relations between entities, so it is constrained to being 1-hop. For other DAs, we performed a qualitative analysis that showed that hops larger than 3 often cause the model to fail to

<sup>3</sup><https://api-ninjas.com/api/animals>

properly realize either the values or the relations. We thus restrict the number of hops to 3. Figure 2 shows two subgraphs of WikiData that correspond to 3-hops (paths of length 3). A 3-hop path derived from the second subgraph is below:

```
[Never Gonna Give You Up] - genre → [Pop]
- genre ← [All I See] - performer → [Kylie
Minogue]
```

Before using this path as the MR for a prompt to an LLM, we convert it to a more text-like representation, namely the representation shown in the 6th Row of Table 1 for the SUGGEST DA.

We then utilize the MR generation package from ViGGO<sup>4</sup> to generate the MRs. The mandatory relations are used first, and then optional relations are randomly sampled to vary the total number of relations in the MRs. In total we create ~71K MRs across the 7 domains, as summarized in the bottom row (Total) of Table 8.

### 3.3 Prompt Creation and LLM Selection

Domain	Dialogue Act	Example
Movies	give_opinion	description of give_opinion: An expression of opinion along with its justification. The response may consist of 1 or 2 sentences, but it must contain both an opinion and its justification. The justification will also be based on the given attributes. No extra information should be added in Data to Text for give_opinion>: Data: Apollo 13 review score excellent. Apollo 13 producer Brian Grazer. Apollo 13 domain movies. Data to Text for give_opinion: I think Apollo 13 is an excellent movie. Brian Grazer is a great producer and he did an outstanding job with this one.
Music	recommend	description of recommend: A question asking if your friend is familiar with a <b>song</b> you would recommend. In the response, bring up the given song (in a recommending way) because it has certain attributes that were discussed earlier in the conversation. Make sure you ask about the <b>song</b> , not the attributes. No extra information should be added in Data to Text for recommend: Data: Littlest Things part of Alright, Still. Alfie part of Alright, Still. Data to Text for recommend: Have you heard the album Alright, Still? It has two great tracks, Littlest Things and Alfie.

Table 4: Example of the Definitional Prompt for Give-Opinion and Recommend for Movies and Music.

Recent work on data-to-text NLG suggests that even though LLMs will have rarely, if at all, seen examples of data-to-text NLG in their training data (Brown et al., 2020; Raffel et al., 2020; Devlin et al., 2019), many LLMs do well on NLG for dialogue (Soltan et al., 2022; Ramirez et al., 2023). However, since we are transferring the DAs from the ViGGO corpus onto new domains, we start off with only ViGGO examples in the prompts. We

<sup>4</sup><https://github.com/jjuraska/slug2slug>

piloted several prompt styles, and decided to use a Definitional prompt format (Gupta et al., 2022; Rastogi et al., 2020; Ramirez et al., 2023). Table 4 provides an example of a DA representation, used with a Definitional prompt, for two DAs across two domains. The section of the prompt labelled “Data” is a conversion of the KG path as described above.

After piloting our approach with ChatGPT-4, Vicuna 13B and GPT-3.5, we selected GPT-3.5 as the best performing LLM. For the model configuration, we employed gpt-3.5-turbo with a top-p value of 1, temperature set to 0.7, and a newline token as the stop token, following Ramirez et al. (2023). Then for all 71K MRs, we overgenerate 10 outputs, and then automatically rank them.

### 3.4 Overgenerate-and-Rank

To create a high-quality data set, we use an overgenerate-and-rank method. The basis for using ranking is a direct translation of the probability of a generated output  $y$ , conditioned on a DA  $d$ , and an MR  $a$ , as in Equation 1. This requires a ranking function that selects outputs that maximize DA accuracy, semantic accuracy, and fluency by assigning a score to each utterance.

$$p(y|d, a) = p(d|y, a) * p(a|y) * p(y) \quad (1)$$

The term  $p(d|y, a)$  requires a highly accurate DA classifier to use in automatic ranking. We utilize the ViGGO DA classifier, which achieves an average F1 over .97 for the ViGGO DAs.<sup>5</sup>

In order to estimate  $p(a|y)$ , semantic accuracy, at generation time in a domain-independent way, we use Beyond-BLEU (BBLEU) (Wieting et al., 2019), which was shown to perform better than other off-the-shelf measures of semantic accuracy such as BLEU, BERTScore and BLEURT (Papineni et al., 2002; Sellam et al., 2020; Zhang et al., 2019). Since these metrics require comparisons with reference utterances, which are not available at generation time, we define a referenceless version based on *pseudo-references*,  $S_{pseudo}$ , created from the input DAs (Juraska, 2022). For any set of KG triples, we create its  $S_{pseudo}$  by converting each triple to a simple sentence Ent1 relation Ent2, and then concatenating all the triples together. Because pseudo-references are available at generation time, we can use pseudo-Beyond-BLEU (pBBLEU) for ranking.

<sup>5</sup><https://github.com/aramir62/da-nlg>

The term  $p(y)$  requires an estimate of fluency. In general, NLG outputs from very large LLMs do not suffer from problems of fluency, but recent work suggests that the probability  $P(S)$  of a generated output  $S$  according to an LLM is a good automatic measure of fluency (Kann et al., 2018; Suzgun et al., 2022). We thus adopt  $P(S)$  to measure fluency, and use GPT-2 to calculate  $P(S)$ .

**RF<sub>DA</sub>: DAC | pBBLEU | P(S)**

Table 5: Ranking function. DAC = probability of the correct DA using a classifier. pBBLEU = pseudo-Beyond-BLEU to measure semantic accuracy.  $P(S)$  = LM probability to measure fluency. The | indicates stepwise evaluation.

For ranking, we adopt the ranking function  $RF_{DA}$  in Table 5. Ramirez et al. compared this ranking function with a ranking function that simply multiplies all the terms as in Equation 1.  $RF_{DA}$  filters first for DA correctness, and then for semantic accuracy, reflecting the importance of DA correctness. Interestingly, Ramirez et al. showed that the  $RF_{DA}$  ranking function also increased **semantic accuracy**, in addition to increasing DA accuracy.

### 3.5 Self Training

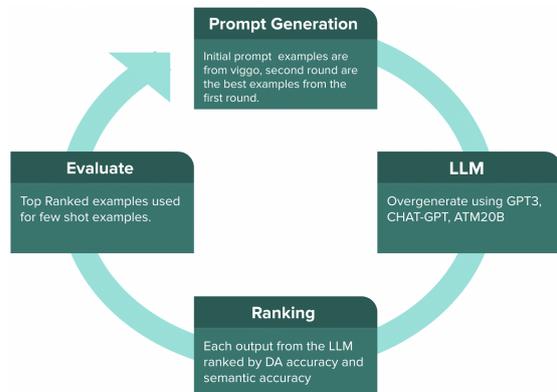


Figure 3: Self-Training Setup for In domain Prompts

One of the challenges with transferring the DA types to new domains is that we have no training data or prompt samples. We investigate a novel approach that uses ViGGO examples in the prompts followed by a round of self-training with silver-generated data. Figure 3 shows the self-training loop in more detail. We start with prompt generation using out of domain, ViGGO examples, then use an LLM to overgenerate examples using these prompts, and rank the outputs. We then select in-domain prompt examples from the top ranked

outputs for self-training. We show in Section 4 that self-training yields a significant improvement in the quality of the Wiki-Dialogue corpus.

The selection of the 10 examples for the in-domain prompts for each DA intent and domain is done manually. We select them following specific criteria in order to provide the LLM a wider knowledge of each DA intents’ realizations. For example, all the attributes of each DA intent are included in the prompt examples in a similar proportion. We also include examples with the different possible number of relations of each DA intent (Table 2). Finally, for scalar attributes with a reduced number of values like RATING, we ensure that all the values are present in a similar proportion.

Interestingly, this two-round self-training process worked successfully for every domain except for Animals, where we received error messages from the LLM complaining about being given data in the wrong domain when using the ViGGO examples. So for Animals, we constructed the 10 in-domain examples for the prompt of each DA intent by hand from a sample of MRs.

## 4 Results and Corpus Quality Evaluation

The Wiki-Dialogue corpus includes more than 71K utterances ( $\sim 50K$  for one-hop and  $\sim 21K$  for multi-hop). Table 8 presents a summary of the resulting corpus. Below, we show the benefits of our self-training and overgenerate-and-rank methods in Tables 6 and 7, respectively. We also carry out a comparison with other KG corpora (Table 9) and a human evaluation (Table 10).

### 4.1 With and Without Self-Training

Domain	BBLEU		DAC	
	N	ST	N	ST
Art	0.84	<b>0.85</b>	<b>1.00</b>	0.99
Music	0.81	0.82	<b>0.98</b>	0.97
Movies	0.77	<b>0.78</b>	0.96	<b>0.97</b>
Sports	0.84	<b>0.86</b>	<b>0.99</b>	0.97
TV	0.75	<b>0.82</b>	0.98	<b>0.98</b>
Boardgames	0.77	<b>0.80</b>	0.94	<b>0.99</b>

Table 6: No self-training = N (Out of Domain Prompts). After self-training = ST (In Domain Prompts). DAC is the DA accuracy using DA classifier. BBLEU is the Beyond BLEU score. Bolded text indicates a significant result (paired t-test,  $p < 0.05$ )

Table 6 compares the performance of generating outputs with the self-trained domain-specific examples and using ViGGO prompt examples for all the domains except animals (due to its problems with the generation with the ViGGO prompt examples).

We compare them using two metrics: BBLEU (Wieting et al., 2019), and DAC. After self-training (ST), the BBLEU scores show a significant improvement (paired t-test,  $p < 0.05$ ), except for the Music domain. Overall for DAC, self-training improves some domains but not others. However, the DAC is consistently high across all domains, with values ranging from 0.97 to 0.99 after self-training.

### 4.2 Before and After Overgenerate-and-Rank

Domain	BBLEU		DAC	
	B	A	B	A
Art	0.80	<b>0.85</b>	0.91	<b>0.99</b>
Music	0.78	<b>0.82</b>	0.88	<b>0.97</b>
Movies	0.72	<b>0.78</b>	<b>0.99</b>	0.95
Sports	0.82	<b>0.86</b>	0.87	<b>0.97</b>
TV	0.68	<b>0.82</b>	0.89	<b>0.98</b>
Animals	0.61	<b>0.75</b>	0.86	<b>0.90</b>
Boardgames	0.78	<b>0.80</b>	0.72	<b>0.99</b>

Table 7: Before (B) and After (A) Overgenerate-and-Rank. DAC is the DA accuracy of the DA classifier using self-training examples. BBLEU is the Beyond BLEU score. Bolded text indicates a significant result (paired t-test  $p < 0.05$ ).

Table 7 compares performance before and after applying the overgenerate-and-rank method. The results show that the method consistently improves the performance across all domains, e.g. in the Music domain, the BBLEU score increases from 0.78 to 0.82, and the DAC improves remarkably from 0.88 to 0.97, while in Sports, the BBLEU score increases from 0.82 to 0.86 and the DAC from 0.87 to 0.97 (paired t-test,  $p < 0.05$ ). We see similar results for the rest of the domains. One reason for an increase in the BBLEU scores is that the overgenerate-and-rank method produces a diverse number of outputs, which increases the chances of a candidate output capturing all the attributes from an MR, while outputs that perform worse are dropped after ranking.

### 4.3 Automatic Evaluation

Table 8 presents a summary of the Wiki-Dialogue corpus, with results for each domain (and also the total) split into one-hop and multi-hop generation. The results for all the domains for BBLEU ( $> 0.75$ ) and DAC ( $> 0.92$ ) are very good.

Table 8 also reports four more automatic metrics: vocabulary size, Canary% (Kim et al., 2022), MLTD (McCarthy and Jarvis, 2010) and Flesch-Kincaid (Kincaid et al., 1975). The vocabulary size is a common metric reported for NLG engines,

Domain	Counts		BBLEU		DAC		Vocab Size		Canary%		MLTD		Flesch-Kincaid	
	One	Multi	One	Multi	One	Multi	One	Multi	One	Multi	One	Multi	One	Multi
Art	6297	-	0.85	-	0.99	-	1953	-	0.11	-	56.15	-	9.16	-
Music	5342	3000	0.81	0.86	0.98	0.94	2047	2573	0.00	0.07	40.72	76.56	6.67	8.09
Sports	3473	3000	0.84	0.88	0.99	0.95	3025	3321	0.00	0.00	47.50	67.06	7.56	8.36
TV	7030	5956	0.80	0.85	0.97	1.00	2847	4640	0.00	0.02	45.99	38.10	8.07	9.07
Movies	7083	8295	0.78	0.78	0.97	0.99	3721	4053	0.00	0.00	48.68	36.40	7.37	7.02
Animals	19092	-	0.75	-	0.92	-	4248	-	0.74	-	45.91	-	7.31	-
Boardgames	1500	1500	0.75	0.85	1.00	0.99	504	913	0.40	0.00	54.89	77.23	6.64	7.10
<b>TOTAL</b>	49817	21751	0.78	0.83	0.96	0.98	12985	11051	0.31	0.01	47.36	45.43	7.60	7.94

Table 8: Automatic Evaluation Metrics. Counts are the number of unique MRs. DAC is the DA accuracy using a DA classifier. BBLEU is the Beyond BLEU score. Vocab size defines the number of unique tokens. Canary% is the percentage of sentences considered unsafe by the Canary model. MLTD is a measure of lexical richness. Flesch-Kincaid is a metric of readability. The data is split into one-hop (One) vs. multi-hop (Multi).

especially since neural training methods tend to reduce the size of the vocabulary from the original corpus (Juraska, 2022). The **Vocab Size** column shows that both one and multi-hop utterances use a large vocabulary.

The **Canary%** column is based on the use of the Canary model to analyze which utterances could be problematic in terms of ethical issues, rudeness, toxicity or bias, inspired by work on SODA (Kim et al., 2023). While Kim et al. filtered 5% of the SODA outputs based on Canary, we found that less than the 0.3% of the utterances are considered ethically inappropriate by the model. We did not filter these utterances because a manual check showed that Canary is very sensitive to certain entities, but the utterances are not actually ethically dangerous.

The MLTD and Flesch-Kincaid metrics estimate the lexical richness and readability of the corpus. For both metrics, the results for the Wiki-Dialogue corpus show no large differences across domains for one-hop and multi-hop. The Flesch-Kincaid values show that the Wiki-Dialogue outputs can be understood by the average American, so they are appropriate for a dialogue.

Table 9 shows a comparison of Wiki-Dialogue with other corpora based on KGs, namely WebNLG which is based on DBPedia and OpenDialKG, which is based on FreeBase (Han and Gardent, 2023; Moon et al., 2019). Wiki-Dialogue is larger than WebNLG but smaller than OpenDialKG (column **N**), even though WebNLG covers more domains, and OpenDialKG covers fewer domains. Wiki-Dialogue uses nine different DAs (column **DAs**), while WebNLG only has the **INFORM** DA, and OpenDialKG is not labelled for DAs. Wiki-Dialogue covers 7 domains (column **Dom**), while WebNLG covers 19 domains, and OpenDialKG covers 4 domains. The Music, Sports, and Movies domains are represented in all three datasets. This

suggests that future work could possibly benefit from using a combination of these corpora.

Corpus	Wiki-Dial	WebNLG	OpenDialKG
<b>N</b>	71568	47915	91829
<b>DAs</b>	9	1	?
<b>Dom</b>	7	19	4
<b>Can%</b>	0.22	0.15	0.03
<b>Vocab</b>	18359	6646	20574
<b>MLTD</b>	46.75	27.27	66.23
<b>FK</b>	7.69	8.93	3.71

Table 9: Comparison of Wiki-Dialogue with other corpora based on a Knowledge-Graph. **N** is the number of unique MRs. **DAs** is the number of Dialogue Act types. **Dom** is the number of domains. **FK** is Flesch-Kincaid.

The **Can%** column shows that all of the corpora are very safe (Kim et al., 2022), perhaps because they are all knowledge grounded. The **Vocab** and **MLTD** columns show that Wiki-Dialogue has a larger vocabulary and is more lexically diverse than WebNLG despite the fact that WebNLG covers more domains. Compared to OpenDialKG, Wiki-Dialogue has lower lexical diversity (column **MLTD**), which may be due to the fact that OpenDialKG is human-human. Both Wiki-Dialogue and WebNLG have a higher Flesch-Kincaid (column **FK**) reading level than OpenDialKG, probably because OpenDialKG restricts MRs to 1 and 2 hops, making utterances shorter on average.

#### 4.4 Human Evaluation

Table 11 and Table 12 in the Appendix provide example realizations of every DA for all 7 domains. These examples show that the quality of the corpus is high: the realizations are natural and the DAs are correctly realized with high accuracy.

For human evaluation, we selected 100 examples from each domain for both one-hop and multi-hop yielding 1200 examples with 700 one-hop examples and 500 multi-hop examples. Five expert anno-

Domain	↓ HAL		↑ PERF		↑ DAC		↑ SAC	
	One	Multi	One	Multi	One	Multi	One	Multi
Art	0.01	-	0.98	-	1.00	-	0.98	-
Animals	0.11	-	0.82	-	0.89	-	0.93	-
BoardGames	0.00	0.18	0.97	0.76	1.00	0.81	0.97	0.95
Movies	0.01	0.06	0.91	0.82	0.98	0.82	0.94	1.00
Music	0.01	0.18	0.97	0.89	1.00	1.00	0.97	0.89
Sports	0.02	0.09	0.97	0.95	1.00	1.00	0.97	0.95
TV	0.00	0.10	0.93	0.73	0.98	0.80	0.95	0.92
TOTAL	0.02	0.10	0.94	0.83	0.98	0.89	0.96	0.94

Table 10: Human annotation results for HAL (Hallucinations), DAC (DA accuracy), and SAC (Semantic Accuracy). PERF (Perfect Semantic and DA accuracy) is calculated based on DAC and SAC.

tators were given a manual for DAs, and provided examples of hallucinations and utterances that were both perfect and not perfect. Each set of utterances were annotated for hallucinations (HAL), DA accuracy (DAC), and semantic accuracy (SAC). The 100 example-set for each domain and hop type was annotated by one annotator, and then 30 of these were re-annotated by a second expert to estimate inter-annotator agreement using Cohen’s Kappa. The average DAC Kappa is 0.94, and the SAC Kappa is 0.89 showing a very high level of agreement between the annotators.

We then calculated the percentage of utterances that had both perfect DAC and SAC (PERF). The results are in Table 10. Overall, Table 10 shows that the quality of the Wiki-Dialogue corpus is high, with perfect outputs that correctly realize both the specified DA and the set of KG triples in the MR ranging from 73% for TV multi-hop to 98% for Art, with an overall average over both hop types of 89% PERFECT outputs. We see that one-hop datasets have fewer hallucinations, better DA accuracy and more perfect utterances. While some values for hallucinations seem high, e.g. 18% for Music multi-hop, these values compare favorably to crowdsourced corpora such as MultiWOZ or E2E, as discussed in Section 2.

## 5 Conclusion and Future Work

This paper describes and provides a novel 71K utterance corpus called Wiki-Dialogue, covering 9 DAs and 7 KG based domains that are useful for both social conversation and task-oriented dialogue. The corpus includes both one and multi-hop sets of KG triples, and the MRs vary from a single triple for some REQUEST dialogue acts up to 8 triples for INFORM DAs.

Our novel method involves a self-training step to create prompt examples for novel domains, fol-

lowed by an overgenerate-and-rank step, and we show that these two steps combined drastically improve the quality of the corpus. We assess the quality of the corpus with both automatic and human evaluation and find that the quality is high. We hand annotate for hallucinations and semantic errors and find their frequency to be lower than reported values for crowdsourced corpora such as E2E and MultiWOZ (Dušek et al., 2019, 2020; Qian et al., 2021). We also check the corpus for safety using Canary and find that only 0.22% of the utterances are flagged as needs\_intervention, while a manual inspection of these suggests that there are no safety issues with them.

In a comparison with similar corpora such as WebNLG and OpenDialog, we observe that although WebNLG cover more domains Wiki-Dialogue is lexically richer and has a larger vocabulary. While WebNLG has only INFORM DAs, Wiki-Dialogue covers 9 DAs, providing a complementary and unique resource to the dialogue community.

Future work should explore how the Wiki-Dialogue corpus can be used to train an NLG for dialogue. In preliminary experiments, a subset of Wiki-Dialogue was used to fine-tune a 3B multi-domain NLG engine, which was tested in Athena, a real-time Amazon Alexa Prize system, with human users of Amazon “Let’s Chat” (Yue Fan and Wang, 2023). The deployment used a cross-domain universal dialogue policy based on the nine Wiki-Dialogue DAs. An example of this policy is shown in the music conversation in Figure 1.

**Ethical Considerations.** There are several potential risks with such an NLG. LLMs introduce the possibility of disinformation, often called hallucinations, whose control is an active area of research. One of the challenges is that it is very difficult to automatically identify them. Here we experiment with ranking functions for better control of hallucinations, hand-label hallucinations and characterize them. Another potential risk is that some of the DAs, like recommend and suggest, could be used in an application to persuade a user to buy something. **Acknowledgments.** This work has been partially funded by Spanish MCIU by the BEWORD project (grant number PID2021-126061OB-C42) and by the Basque Government under grant PRE 2020 1 0274.



## References

- Hannah Bast, Florian Baurle, Björn Buchhold, and Elmar Haubmann. 2014. Easy access to the Freebase dataset. In *Proceedings of the 23rd International Conference on World Wide Web*, pages 95–98.
- 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.
- Paweł Budzianowski and Ivan Vulić. 2019. Hello, It’s GPT-2 - How Can I Help You? Towards the Use of Pretrained Language Models for Task-Oriented Dialogue Systems. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 15–22.
- Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Seokhwan Kim, Andy Rosenbaum, Yang Liu, Zhou Yu, and Dilek Hakkani-Tur. 2023. **PLACES: Prompting language models for social conversation synthesis**. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 844–868, Dubrovnik, Croatia. Association for Computational Linguistics.
- Ethan A Chi, Ashwin Paranjape, Abigail See, Caleb Chiam, Trenton Chang, Kathleen Kenealy, Swee Kiat Lim, Amelia Hardy, Chetanya Rastogi, Haojun Li, et al. 2022. Neural generation meets real people: Building a social, informative open-domain dialogue agent. In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 376–395.
- Emilie Colin, Claire Gardent, Yassine M’rabet, Shashi Narayan, and Laura Perez-Beltrachini. 2016. The WEBNLG challenge: Generating Text from DBpedia data. In *Proceedings of the 9th international natural language generation conference*, pages 163–167.
- 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.
- Ondřej Dušek, David M Howcroft, and Verena Rieser. 2019. Semantic noise matters for neural natural language generation. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 421–426.
- Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2020. Evaluating the state-of-the-art of end-to-end natural language generation: The E2E NLG challenge. *Computer Speech & Language*, 59:123–156.
- Nouha Dziri, Andrea Madotto, Osmar R Zaiane, and Avishek Joey Bose. 2021. Neural path hunter: Reducing hallucination in dialogue systems via path grounding. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2197–2214.
- Mihail Eric, Nicole Chartier, Behnam Hedayatnia, Karthik Gopalakrishnan, Pankaj Rajan, Yang Liu, and Dilek Hakkani-Tur. 2021. **Multi-sentence knowledge selection in open-domain dialogue**. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 76–86, Aberdeen, Scotland, UK. Association for Computational Linguistics.
- Thiago Castro Ferreira, Diego Moussallem, Emiel Kraemer, and Sander Wubben. 2018. Enriching the WebNLG corpus. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 171–176.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. Creating training corpora for nlg micro-planners. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 179–188.
- Karthik Gopalakrishnan, Behnam Hedayatnia, Qinlang Chen, Anna Gottardi, Sanjeev Kwatra, Anu Venkatesh, Raefer Gabriel, Dilek Hakkani-Tür, and Amazon Alexa AI. 2019. Topical-chat: Towards knowledge-grounded open-domain conversations. *Proc. Interspeech 2019*, pages 1891–1895.
- Raghav Gupta, Harrison Lee, Jeffrey Zhao, Yuan Cao, Abhinav Rastogi, and Yonghui Wu. 2022. **Show, don’t tell: Demonstrations outperform descriptions for schema-guided task-oriented dialogue**. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4541–4549, Seattle, United States. Association for Computational Linguistics.
- Kelvin Han and Claire Gardent. 2023. **Generating and answering simple and complex questions from text and from knowledge graphs**. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 285–304, Nusa Dua, Bali. Association for Computational Linguistics.
- Behnam Hedayatnia, Karthik Gopalakrishnan, Seokhwan Kim, Yang Liu, Mihail Eric, and Dilek Hakkani-Tur. 2020. Policy-driven neural response generation for knowledge-grounded dialog systems. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 412–421.
- Juraj Juraska. 2022. *Diversifying Language Generated by Deep Learning Models in Dialogue Systems*. Ph.D. thesis, UC Santa Cruz.

- Juraj Juraska, Kevin Bowden, Lena Reed, Vrindavan Harrison, Wen Cui, Omkar Patil, Rishi Rajasekaran, Angela Ramirez, Cecilia Li, Eduardo Zamora, et al. 2021. Athena 2.0: Contextualized Dialogue Management for an Alexa Prize SocialBot. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 124–133.
- Juraj Juraska, Kevin K Bowden, and Marilyn Walker. 2019. ViGGO: A video game corpus for data-to-text generation in open-domain conversation. In *Proceedings of the 12th International Conference on Natural Language Generation*.
- Katharina Kann, Sascha Rothe, and Katja Filippova. 2018. Sentence-level fluency evaluation: References help, but can be spared! In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 313–323.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Bras, Malihe Alikhani, Gunhee Kim, et al. 2023. SODA: Million-scale Dialogue Distillation with Social Commonsense Contextualization. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12930–12949.
- Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. 2022. ProsocialDialog: A prosocial backbone for conversational agents. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4005–4029, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. 1975. Derivation of new readability formulas (automated reliability index, fog count and flesch reading ease formula) for navy enlisted personnel. Technical Report Research Branch Report 8-75, Naval Technical Training, US Naval Air Station, Millington, TN.
- Rémi Lebret, David Grangier, and Michael Auli. 2016. Neural text generation from structured data with application to the biography domain. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1203–1213.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia—a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic web*, 6(2):167–195.
- Philip M McCarthy and Scott Jarvis. 2010. MTL D, VOCD-D, and HD-D: A validation study of sophisticated approaches to lexical diversity assessment. *Behavior Research Methods*, 42(2):381–392.
- Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. OpenDialKG: Explainable conversational reasoning with attention-based walks over knowledge graphs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 845–854, Florence, Italy. Association for Computational Linguistics.
- Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xiangru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, et al. 2021. DART: Open-domain structured data record to text generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 432–447.
- Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017. The E2E Dataset: New Challenges For End-to-End Generation. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 201–206.
- 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.
- Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A controlled table-to-text generation dataset. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1173–1186.
- Kun Qian, Ahmad Beirami, Zhouhan Lin, Ankita De, Alborz Geramifard, Zhou Yu, and Chinnadhurai Sankar. 2021. Annotation inconsistency and entity bias in MultiWOZ. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 326–337, Singapore and 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. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67.
- Angela Ramirez, Kartik Aggarwal, Juraj Juraska, Utkarsh Garg, and Marilyn Walker. 2023. Controllable generation of dialogue acts for dialogue systems via few-shot response generation and ranking. In *Meeting of the Special Interest Group on Discourse and Dialogue, SIGDIAL*.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8689–8696.

- Abigail See, Stephen Roller, Douwe Kiela, and Jason Weston. 2019. What makes a good conversation? how controllable attributes affect human judgments. 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 1702–1723.
- 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.
- Salah Soltan, Shankar Ananthakrishnan, Jack FitzGerald, Rahul Gupta, Wael Hamza, Haidar Khan, Charith Peris, Stephen Rawls, Andy Rosenbaum, Anna Rumshisky, et al. 2022. AlexaTM 20b: Few-shot learning using a large-scale multilingual seq2seq model. arXiv preprint arXiv:2208.01448.
- Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2022. Prompt-and-rerank: A method for zero-shot and few-shot arbitrary textual style transfer with small language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2195–2222.
- Tsung-Hsien Wen, Milica Gasic, Nikola Mrkšić, Pei-Hao Su, David Vandyke, and Steve Young. 2015. Semantically conditioned LSTM-based natural language generation for spoken dialogue systems. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1711–1721.
- John Wieting, Taylor Berg-Kirkpatrick, Kevin Gimpel, and Graham Neubig. 2019. [Beyond BLEU: training neural machine translation with semantic similarity](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4344–4355, Florence, Italy. Association for Computational Linguistics.
- Qingyang Wu, James Gung, Raphael Shu, and Yi Zhang. 2023. DiactTOD: Learning generalizable latent dialogue acts for controllable task-oriented dialogue systems. arXiv preprint arXiv:2308.00878.
- Xinnuo Xu, Guoyin Wang, Young-Bum Kim, and Sungjin Lee. 2021. [AugNLG: Few-shot natural language generation using self-trained data augmentation](#). 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 1183–1195, Online. Association for Computational Linguistics.
- Fanghua Ye, Jarana Manotumruksa, and Emine Yilmaz. 2022. MultiWOZ 2.4: A multi-domain task-oriented dialogue dataset with essential annotation corrections to improve state tracking evaluation. In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 351–360.
- Wen Cui Winson Chen Vrindavan Harrison Angela Ramirez Saaket Agashe Xinyue Gabby Liu Neha Pulabhotla Nan Qiang Jeshwanth Bheemanpally Sugam Garg Marilyn Walker Yue Fan, Kevin K. Bowden and Xin Eric Wang. 2023. Athena 3.0: Personalized multimodal chatbot with neuro-symbolic dialogue generators. *Proceedings of the Alexa Prize*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTScore: Evaluating text generation with BERT. In *International Conference on Learning Representations*.
- Giulio Zhou and Gerasimos Lampouras. 2020. WebNLG Challenge 2020: Language agnostic delexicalisation for multilingual RDF-to-Text generation. In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*, pages 186–191.

# A Appendix

domain	MR	text
<b>verify_attribute</b>		
sports	Venus Williams rating okay. Venus Williams start time 1994. Venus Williams domain athlete.	You said you think Venus Williams is okay. Do you think she is an average athlete compared to other players who started in 1994?
movies	Avengers: Endgame review score good. Avengers: Endgame cast member Robert Redford, Taika Waititi. Avengers: Endgame domain movies.	I remember you thought Avengers: Endgame was good. Do you usually enjoy movies with Robert Redford and Taika Waititi?
music	Kings Of The Wild Frontier Good. Kings Of The Wild Frontier	I know you like Kings of the Wild Frontier. Do you like other songs from the Kings of the Wild Frontier album?
tv	The Crown review score average. The Crown genre biographical film, drama television series. The Crown domain tv.	I know you said The Crown had an average review score. Do you feel the same way about other biographical films and drama television series?
art	Raphael rating excellent. Raphael genres portrait, portrait painting. Raphael movements Italian Renaissance, High Renaissance.	You rated the portrait painting by Raphael as excellent. Would you say you generally have a strong liking for Italian Renaissance and High Renaissance portrait paintings?
animals	Carolina Dog good with other dogs friendly to other dogs. Carolina Dog good with children kid-friendly. Carolina Dog domain animals.	I know you said that Carolina Dogs are friendly to other dogs and kid-friendly. Are you generally interested in animals that are friendly to other dogs and kid-friendly?
board games	A Game of Thrones average rating good. A Game of Thrones genre fantasy. A Game of Thrones domain board game.	You said A Game of Thrones was good. Do you usually like fantasy board games?
<b>recommend</b>		
sports	Jake Paul sport boxing. Jake Paul height 1.83. Jake Paul domain athlete.	Have you heard of Jake Paul? He's a professional boxer who stands at 1.83 meters tall.
movies	Elvis review score good. Elvis producer Gail Berman, Catherine Martin, Baz Luhrmann. Elvis domain movies.	Have you seen Elvis, produced by Gail Berman, Catherine Martin, and Baz Luhrmann? It has a good review score.
music	Money For Nothing Rock Music Littlest Things Mark Ronson	Have you heard of the rock song Money for Nothing produced by Mark Ronson?
tv	Peaky Blinders review score good. Peaky Blinders creator David Caffrey. Peaky Blinders domain tv.	Have you seen Peaky Blinders from David Caffrey? It's got a really good review score.
art	The Haywain Triptych creator Hieronymus Bosch. The Haywain Triptych location Museo del Prado.	Have you seen The Haywain Triptych by Hieronymus Bosch at the Museo del Prado?
animals	Staffordshire Bull Terrier trainability very easy to train. Staffordshire Bull Terrier domain animals.	I think you'd really be interested in the Staffordshire Bull Terrier, which is very easy to train. Have you heard of it before?
board games	HeroQuest publisher Games Workshop. HeroQuest genre fantasy. HeroQuest domain board game.	Have you ever heard of HeroQuest, the fantasy board game from Games Workshop?
<b>give_opinion</b>		
sports	Ronaldo sport association football. Ronaldo rating good. Ronaldo domain athlete.	I think Ronaldo is a great footballer. He has consistently put in strong performances for his teams and is known for his skill and passion for the game.
movies	Scooby-Doo review score average. Scooby-Doo director Raja Gosnell. Scooby-Doo MPA film rating PG. Scooby-Doo domain movies.	Scooby-Doo is an average movie, but I think it's worth watching for the great direction of Raja Gosnell and the PG rating.
music	Y.M.C.A. rating excellent. Y.M.C.A. producer Jacques Morali. Y.M.C.A. record label Casablanca Records. Y.M.C.A. domain songs.	Y.M.C.A. by the Village People, produced by Jacques Morali and released by Casablanca Records, is an excellent song.
tv	Gunpowder review score average. Gunpowder screenwriter Ronan Bennett. Gunpowder domain tv.	Gunpowder was an average show, but I appreciate the work of the screenwriter, Ronan Bennett, for making the story interesting.
art	Henri de Toulouse-Lautrec rating good. Henri de Toulouse-Lautrec genres animal painting, portrait.	I think Henri de Toulouse-Lautrec is a good painting. The combination of animal painting and portraiture in this piece is quite engaging.
animals	Fin Whale popularity very popular. Fin Whale weight 1,800kg - 70,000kg 4,000lbs - 150,000lbs. Fin Whale domain animals.	I'm really impressed by the Fin Whale with its huge weight range of 1,800kg - 70,000kg (4,000lbs - 150,000lbs).
board games	Twilight Struggle average rating excellent. Twilight Struggle maximum number of players 2. Twilight Struggle domain board game.	I highly recommend Twilight Struggle. It's an excellent board game with a maximum number of players of 2, making it perfect for a competitive game night.
<b>request_attribute</b>		
sports	athletes sports teams.	What sports team does your favorite athlete play for?
movies	movies review score None. movies domain movies.	What rating would you give to movies that you have watched recently?
music	songs performer.	Who is your go-to artist for music?
tv	show screenwriter None. show domain tv.	What is your favorite TV show and its screenwriter?
art	paintings genres.	What types of painting genres resonate with you the most?
animals	animal good with other dogs. animal domain animals.	Is your favorite animal good with other dogs?
board games	board game genre. board game domain board game.	What board game genre do you find the most enjoyable?

Table 11: Examples of transferred DAs for verify\_attribute, recommend, give\_opinion, and request\_attribute to all 7 domains.

domain	MR	text
<b>suggest</b>		
sports	Wilt Chamberlain Nba Most Valuable Player Award, Nba All-Star Game Kobe Bryant Most Valuable Player Award	Have you ever heard of Wilt Chamberlain, who won the NBA Most Valuable Player Award and the NBA All-Star Game Kobe Bryant Most Valuable Player Award?
movies	Ant-Man and the Wasp cast member Dax Griffin, Hannah John-Kamen. Ant-Man and the Wasp publication date 2018. Ant-Man and the Wasp domain movies.	Did you ever watch Ant-Man and the Wasp, starring Dax Griffin and Hannah John-Kamen, which was released in 2018?
music	Waka Waka performer Freshlyground, Shakira. Waka Waka domain songs.	Do you like the music of Shakira and Freshlyground, like Waka Waka?
tv	Breaking Bad character Tuco Salamanca, Mike Ehrmantraut. Breaking Bad domain tv.	Do you know any characters from Breaking Bad, such as Tuco Salamanca and Mike Ehrmantraut?
art	The Bridge at Narni creator Jean-Baptiste Camille Corot. The Bridge at Narni genres landscape art.	Have you explored landscape art, such as The Bridge at Narni by Jean-Baptiste Camille Corot?
animals	Pygora Goat weight 75-95 lbs. for males, 64-75 lbs. for females.. Pygora Goat lifespan 12-15 years. Pygora Goat domain animals.	Have you heard about the Pygora Goat? It typically weighs 75-95 lbs. for males and 64-75 lbs. for females, and has a lifespan of 12-15 years.
board games	Enchanted Forest maximum number of players 6. Enchanted Forest instance of board game. Enchanted Forest domain board game.	Have you ever played a board game for up to 6 players, like Enchanted Forest?
<b>request</b>		
sports	athletes specifier boring.	Who do you think is the most boring athlete to watch?
movies	movies specifier realistic. movies director Cary Joji Fukunaga. movies domain movies.	Have you seen any realistic movies directed by Cary Joji Fukunaga?
music	songs specifier beautiful.	What's the most beautiful song you can think of?'
tv	show specifier overrated. show screenwriter Arif Aliyev. show domain tv.	Do you think any of Arif Aliyev's TV shows are overrated?
art	paintings specifier bizarre. paintings location Museo Nacional Centro de Arte Reina Sofia.	Can you think of a bizarre painting from the Museo Nacional Centro de Arte Reina Sofia?
animals	animals specifier cheerful. animals habitat Tropical bushland, grassland and savannas.	Do you know of any cheerful animals that live in Tropical bushland, grassland and savannas?
board games	board game specifier frustrating. board game maximum number of players 5. board game domain board game.	What's the most frustrating board game you can play with up to 5 players?
<b>request_explanation</b>		
sports	athletes rating good. athletes team positions center.	What is it about athletes who play the center position that you find so good?
movies	movies review score good. movies screenwriter Steven Knight. movies domain movies.	What do you think makes the movies written by Steven Knight so good?
music	songs rating good. songs from album Katharine McPhee.	What do you like about the songs from Katharine McPhee's album that make them good?
tv	show review score average. show awards received None. show creator Jonathan Frakes, Hanelle Culpepper. show domain tv.	Given the average review score and lack of awards, what do you think Jonathan Frakes and Hanelle Culpepper did right in their TV shows?
art	painters rating good. painters movements Impressionism, post-impressionism.	What is it about Impressionism and post-impressionism paintings that you find so good?
animals	animal coat length short coat. animal protectiveness very alert. animal name Airedale Terrier. animal good with other dogs friendly to other dogs. animal domain animals.	What do you like about the Airedale Terrier that has a short coat, is very alert, and is friendly to other dogs?
board games	board game genre science fiction. board game average rating excellent. board game domain board game.	What makes science fiction board games your favorite type of game?
<b>confirm</b>		
sports	Dwayne Johnson sport professional wrestling. Dwayne Johnson sports teams Calgary Stampeders, Miami Hurricanes football. Dwayne Johnson domain athletes.	Are you talking about Dwayne Johnson, the professional wrestler who previously played for the Calgary Stampeders and the Miami Hurricanes football team?
movies	Thor: Love and Thunder cast member Jacob Yee, Gemma Dart. Thor: Love and Thunder domain movies.	Are you referring to the movie Thor: Love and Thunder, starring Jacob Yee and Gemma Dart?
music	Break My Soul genres house music. Break My Soul domain songs.	Are you referring to the house music song Break My Soul?
tv	Obi-Wan Kenobi character Darth Vader, Obi-Wan Kenobi. Obi-Wan Kenobi screenwriter Hossein Amini. Obi-Wan Kenobi domain tv.	Are you referring to the Obi-Wan Kenobi TV show, written by Hossein Amini, featuring Darth Vader and Obi-Wan Kenobi?
art	The Wayfarer creator Hieronymus Bosch. The Wayfarer genres genre art. The Wayfarer movements Early Netherlandish painting.	Are you referring to The Wayfarer, the Early Netherlandish genre art painting by Hieronymus Bosch?
animals	Carolina Parakeet common class Bird. Carolina Parakeet predators Birds of prey, raccoons, skunks, squirrels, and snakes. Carolina Parakeet domain animals.	Do you mean the Bird known as the Carolina Parakeet, which is preyed upon by Birds of prey, raccoons, skunks, squirrels, and snakes?
board games	Monopoly publisher Hasbro. Monopoly country of origin United States of America. Monopoly domain board game.	Are you talking about the Hasbro board game Monopoly, from the United States of America?
<b>inform</b>		
sports	Magic Johnson sport basketball. Magic Johnson stats leader 1992 NBA All-Star Game. Magic Johnson domain athletes.	Magic Johnson is a basketball athlete who was the leading stats leader for the 1992 NBA All-Star Game.
movies	Captain Phillips 2013 'Film Based On Books', 'Survival Film' 'Faysal Ahmed As Nour Najee' Excellent 'Louis Mahoney', 'Catherine Keener'	Captain Phillips is a 2013 survival film based on books, starring Faysal Ahmed as Nour Najee, Louis Mahoney, and Catherine Keener. It has an excellent review score.
music	As It Was record label Columbia Records. Mireille Mathieu record label Columbia Records. Mireille Mathieu genre pop music. As It Was genre synth-pop. As It Was performer Harry Styles.	Harry Styles' As It Was is a synth-pop album released on Columbia Records. Mireille Mathieu's record label is also Columbia Records and her genre is pop music. (3 hop)
tv	Peaky Blinders genre drama television series, costume drama. Peaky Blinders creator David Caffrey. Peaky Blinders screenwriter Steven Knight. Peaky Blinders review score good. Peaky Blinders character Arthur Shelby, Jessie Eden. Peaky Blinders domain tv.	Peaky Blinders is a good costume drama television series created by David Caffrey and written by Steven Knight. It stars Arthur Shelby and Jessie Eden.
art	The Coronation of Napoleon creator Jacques-Louis David. The Coronation of Napoleon materials canvas, oil paint. The Coronation of Napoleon genres history painting.	The Coronation of Napoleon, a history painting by Jacques-Louis David, is made of oil paint on canvas.
animals	Australian Stumpy Tail Cattle Dog shedding average shedding. Australian Stumpy Tail Cattle Dog domain animals.	The Australian Stumpy Tail Cattle Dog has average shedding.
board games	Enchanted Forest genre fantasy. Enchanted Forest maximum number of players 6. Enchanted Forest domain board game.	Enchanted Forest is a fantasy board game for up to 6 players.

Table 12: Examples of transferred DAs for suggest, request, request\_explanation, confirm and inform to all 7 domains

# Incremental Learning for Knowledge-Grounded Dialogue Systems in Industrial Scenarios

Izaskun Fernandez<sup>1</sup>, Cristina Aceta<sup>1</sup>, Cristina Fernandez<sup>1</sup>, María Inés Torres<sup>2</sup>,  
Aitor Etxalar<sup>2</sup>, Ariane Mendez<sup>3</sup>, Maia Agirre<sup>3</sup>, Manuel Torralbo<sup>3</sup>,  
Arantza del Pozo<sup>3</sup>, Joseba Andoni Agirre<sup>4</sup>, Egoitz Artetxe<sup>4</sup>, Iker Altuna<sup>4</sup>

<sup>1</sup>TEKNIKER - Basque Research and Technology Alliance (BRTA),

<sup>2</sup>Speech Interactive Research Group (SPIN) - University of the Basque Country (UPV/EHU),

<sup>3</sup>Vicomtech Foundation - Basque Research and Technology Alliance (BRTA), <sup>4</sup>IMH Campus

Correspondence: [izaskun.fernandez@tekniker.es](mailto:izaskun.fernandez@tekniker.es)

## Abstract

In today’s industrial landscape, seamless collaboration between humans and machines is essential and requires a shared knowledge of the operational domain. In this framework, the technical knowledge for operator assistance has traditionally been derived from static sources such as technical documents. However, experienced operators hold invaluable know-how that can significantly contribute to support other operators. This work focuses on enhancing the operator assistance tasks in the manufacturing industry by leveraging spoken natural language interaction. More specifically, a Human-in-the-Loop (HIL) incremental learning approach is proposed to integrate this expertise into a domain knowledge graph (KG) dynamically, along with the use of in-context learning for Large Language Models (LLMs) to benefit other capabilities of the system. Preliminary results of the experimentation carried out in an industrial scenario, where the graph size was increased in a 25%, demonstrate that the incremental enhancing of the KG benefits the dialogue system’s performance.

## 1 Introduction

Human-Machine Interaction (HMI) is revolutionizing traditional industrial processes. Smart manufacturing relies on the collaboration between highly advanced machinery and the knowledge and decision-making abilities of human operators. The industry of the near future requires qualified personnel specialized in technologies such as robotics and artificial intelligence (AI), capable of making informed decisions based on these factors. In this context, a human-centered approach positions operators as a crucial element in new industrial plants. Thanks to the latest technological advances, voice interaction between operators and industrial manufacturing systems or machines is now feasible. Moreover, these technologies are *hands-free* and *eyes-free*,

enabling operators to perform physical tasks, *support natural language communication* that requires minimal training, and are *highly flexible*, allowing communication at various levels of detail. Consequently, there has been an increase in the number of prototypes and systems exploring the use of voice as a natural interaction interface between operators and machines in industrial environments in recent years. Additionally, dialogue modeling and management have drastically changed due to the recent success of large language models (LLMs). However, any application based on LLMs needs reliable and up-to-date knowledge sources. In particular, industrial scenarios require robust models capable of handling very technical and precise knowledge, which is necessary for tasks shared by humans and machines.

Traditionally, HMI has relied on rule-based systems to represent knowledge and actions, ensuring everything remains under control. As a result, these interaction systems are static, failing to capture the expert human knowledge of the factory that is not documented or included in the system’s knowledge base. This limitation can be addressed through the concept of Human in the Loop (HIL), also known as Operator in the Loop (OIL) in industrial contexts. These AI systems facilitate collaboration between humans and machines to enhance results and accelerate the learning process. The HIL paradigm involves continuous interaction throughout all post-deployment stages of AI models. As illustrated in Figure 1, in the industrial sector, the OIL paradigm enables the integration of expert knowledge into HMI interfaces by providing feedback using natural language. This approach allows voice interaction systems to evolve over time, adapting to the unique dynamics of each factory and incorporating the expertise that operators develop.

In this work, an OIL incremental learning approach to manage knowledge-grounded, task-oriented dialogue (TOD) systems in industrial set-

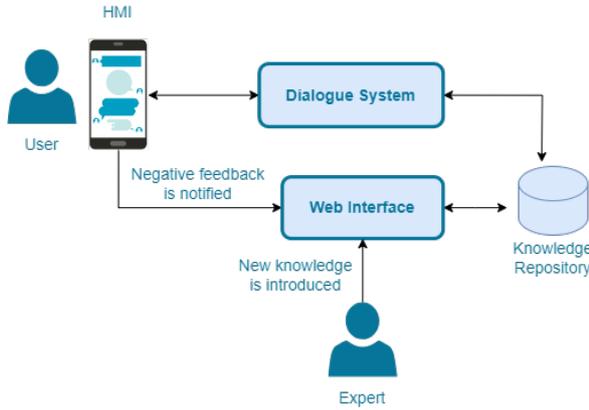


Figure 1: Operator in the Loop paradigm.

tings is proposed, being its main contributions (1) extending a previously defined ontology to support the management and storage of new knowledge provided by experts; (2) developing online learning capabilities to collect user feedback, thereby updating and expanding a knowledge graph and (3) developing an LLM-based natural language understanding (NLU) system that queries a KG to constrain it within the task. Preliminary evaluations show promising results in NLU performance and KG grounding.

The rest of the paper is structured as follows: Section 2 presents the related work. The proposed knowledge graph-based incremental dialogue system with the detailed description of each of its modules is described in Section 3, and the initial validation results in Section 4. Finally, conclusions of this work are shown in Section 5.

## 2 Related work

The current state of the art considers knowledge graphs as a useful asset in industrial settings and, more specifically in human-centric approaches (Abonyi et al., 2024), such as robot interaction and collaborative manufacturing (Nagy et al., 2024). In this line, approaches such as the one proposed by Nagy et al. (2024) are observed, in which knowledge graphs are used to model factors related to the operator and their conditions, such as movements or collaboration with machines. Moreover, knowledge graphs have been used in this scenario for task-oriented dialogue, which enable operators to communicate to industrial systems in a more natural way. In this context, knowledge graphs have been traditionally used to model the domain of the use case, providing a detailed representation

of the scenario and reducing ambiguity between the agents involved (Sidi Yakoub et al., 2015). However, more modern approaches also make use of knowledge graphs for dialogue management (Teixeira et al., 2021; Aceta et al., 2022)

Of course, this process also has an impact on dialogue management, since one of the most widespread techniques is to obtain this information from users. To do this, dialogues are generated dynamically to be able to obtain the necessary information for the system to learn, as well as the appropriate moment for it, based on a strategy (Liu and Mazumder, 2021). Some approaches also base these interactions on the feedback obtained from the user taking into account, for example, evaluations such as “it’s not what I wanted” or “you didn’t understand me well” (Veron et al., 2021).

In the field of Natural Language Processing (NLP), there’s a clear surge in leveraging state-of-the-art strategies across multiple applications, particularly through the deployment of pre-trained Large Language Models (LLMs) in dialogue systems. Ozdemir (2023) describes these models as AI models that often, though not exclusively, stem from the Transformer architecture. They are crafted to understand and generate human language, code, and beyond. Also, they are trained on immense troves of text, and they can tackle a vast array of language-related tasks, from simple text classification to elaborate text generation. As highlighted by this author, the LLMs available in the market (like various versions of GPT, Gemini, Llama, among others) have been pre-trained on extensive datasets from diverse sources using distinct methodologies. Thus, not all LLMs perform equally, and their training processes significantly influence their performance in specific applications.

Therefore, to optimize pre-trained language models for task-oriented dialogue systems, different works employ models like Alpaca, GPT-Neo, BART, T5, Llama 2 and GPT-3.5 (Hudeček and Dušek, 2023; Andreas et al., 2022; Li et al., 2022; Hu et al., 2024); among others. Also, different authors adopt various approaches for constructing these dialogue systems. Prominent among these is fine-tuning pre-trained language models using methods like LoRA (Low Rank Adaptation) (Andreas et al., 2022; Li et al., 2022), prompt tuning (Cao, 2023; Hudeček and Dušek, 2023) and Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022), among others. This indicates a clear trend in using LLMs in TOD sys-

tems and various optimization strategies. However, many of these methods demand substantial specific data for training as they are data-driven, which may not be available for certain industrial use cases.

To address this issue and avoid the need for hand-crafted rules, in-context learning approaches are becoming increasingly popular. These approaches involve designing prompts using snippets of example dialogues, the user’s goal, and the dialogue history (Sekulić et al., 2024). This optimization method, known as prompt tuning, allows adapting the model to task requirements without requiring a corpus or extensive training, just relying on natural language instructions to guide the model’s behaviour.

### 3 Knowledge-Grounded Incremental Task-oriented Dialogue System

Two of the most common applications of TOD systems in industrial scenarios are to provide assistance through processes and to deliver tasks to a certain industrial intelligent system. Therefore, the expected interactions from the user can mainly be classified as navigation instructions through processes and action requests to industrial systems, respectively. The TOD system’s responses, on the other side, must be in the form of steps of the processes on which the user will request assistance for the former, and the corresponding machine-readable action for the latter.

So, in this type of scenarios, towards an incremental approach, feedback may be useful in these two situations, mainly: (1) the content presented does not meet the needs of the user or (2) the interpretation of the interaction indicates that what the user wants to do next or deliver to the system is not appropriate.

This work presents the extension and adaptation of KIDE4I, presented in Aceta et al. (2022) and based on the TODO Ontology (Aceta et al., 2021), to provide it with feedback-capturing and management capabilities. The aim of such task is to achieve a system that is capable of learning from interactions with users over time and, thus, improve its interpretation and dialogue capacities, as well as adapting to the users’ needs. To this end, the following aspects have been addressed:

1. Extension of the TODO ontology to support the management and storage of new knowledge based on feedback (described in Section 3.1).

2. New functionalities to generate dialogues aimed at collecting feedback and to update knowledge extracted from it (described in Section 3.2).

Likewise, and towards assessing the benefits when updating KIDE4I with the most recent technologies, in-context learning (ICL) of LLMs through prompt-tuning has been explored and implemented in the natural language understanding (NLU) module, as detailed in Section 3.3. This task has allowed to compare more traditional strategies, such as rule-based ones, with the most disruptive one nowadays: the use of LLMs in scenarios with limited resources (in terms of training corpus), such as industrial ones.

#### 3.1 Industrial-Assistance-Oriented Incremental Knowledge Graph

By definition, a knowledge graph focuses on representing relationships and capturing real-world connections, ideally based on an ontology that provides the formal framework for defining the terms and concepts used in that representation.

As described previously, the focus of this work is developing a knowledge-graph-based TOD system for industrial scenarios which is based on technical documentation and expert knowledge that can be extended over time through, for instance, feedback gathering. To achieve such a system it is necessary to construct a knowledge graph that formally represents all this information, relying on a agreed ontology that allows an incremental learning approach.

The core ontology for developing the knowledge graph in the context of this work is TODO (Aceta et al., 2021), the main modules of which can be seen in Figure 2. This modular ontology is designed to enable task-oriented dialogue systems to interact naturally with users at both understanding and communication levels by distinguishing two main areas of knowledge: domain (TODODom) and dialogue (TODODial), respectively. It can be readily adapted to various industrial settings, thus minimizing the time and cost of adaptation. Additionally, it supports the storage and reproduction of the dialogue process, allowing for learning from new interactions. However, this tracing capability, although being a good starting point for supporting an incremental learning approach, does not support the generation and management of user feedback. In order to solve this gap, the TODO ontology has

been extended.

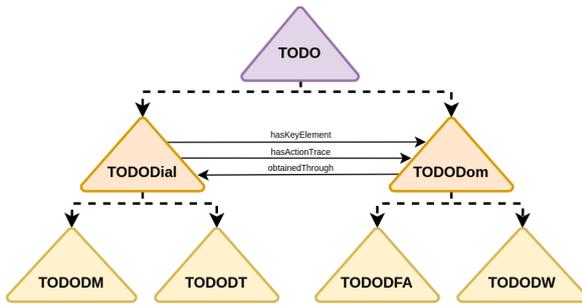


Figure 2: TODO ontology (Aceta et al., 2021)

In that extension and adaptation task, new classes and relationships that allow representing the key concepts aimed at collecting feedback have been added. More precisely, 2 classes (C) and 4 object properties (OP) have been created in the TODODom module and 2 classes in TODODM, which are listed below, by module.

#### *TODODom (domain)*

- DefinitiveLexicalUnit (C), to depict lexical units (i.e., variants) that have not been added through feedback (i.e. manually or in a supervised way) or lexical units that have been added through feedback several times.
- ProvisionalLexicalUnit (C), to depict lexical units that have been added through feedback but the confidence to consider them as definitive is still low.
- hasDefinitiveLexicalUnit (OP), to relate frame heads (i.e., generic terms to agglutinate different variants) to their corresponding definitive lexical units.
- hasProvisionalLexicalUnit (OP), to relate frame heads to their corresponding provisional lexical units.
- isDefinitiveLexicalUnitOf (OP). Inverse property of hasDefinitiveLexicalUnit.
- isProvisionalLexicalUnitOf (OP). Inverse property of hasProvisionalLexicalUnit.

#### *TODODM (dialogue management)*

- NewLexicalUnitConfirmationRequest (C), to request the user for confirmation to relate a lexical unit to a specific frame head.

- ActionDetectedResponse (C), to inform the user that it has detected an action (for which the command includes a new reference to potentially be added to the graph).

These classes and properties have been added by following the LOT methodology (Poveda-Villalón et al., 2019), which makes sure that knowledge is modelled into the ontology ensuring its quality. Therefore, the quality of the ontology (compared to the results obtained in Aceta et al. (2021)), has not been affected.

With the ontology ready, a manual instantiation of the newly-modelled, dialogue-related classes has been carried out, in order to offer the dialogue manager variations to interact with the user and direct the dialogue to capture feedback, such as “Can you confirm that {item} is a related word?”. The rest of the dialogue-related instantiations have been reused from the generic instantiation of TODODial (Aceta et al., 2022).

As for the domain section of the knowledge graph, it is instantiated automatically. First of all, the relevant procedures have been defined by the experts by using an interface designed to simplify the instantiation process. In a nutshell, this interface, once a procedure is defined, generates, first of all, a JSON file. This JSON file, by following an *Extract, transform and load* (ETL) process, is transformed into RDF and uploaded to the RDF store, which, in this case, has been Virtuoso 8.3. An example snippet of an instantiated procedure can be found in Appendix A.

This first graph version enables the system to be ready to be used and its knowledge to be extended through user feedback in subsequent interactions.

### 3.2 Dialogue Management Supporting Incremental Approach

Once the ontology is extended and the dialogue instances for collecting feedback are ready, as reported in Section 3.1, it is necessary to add to the dialogue manager the capability to extract the new knowledge to be included in the system.

As described previously, the two situations that may require feedback gathering would be when the system is not capable of correctly interpreting a user request and when the information provided by the system is not accurate.

To respond to the first situation, the dialogue manager has been extended so that, instead of asking the user to reformulate the request because they

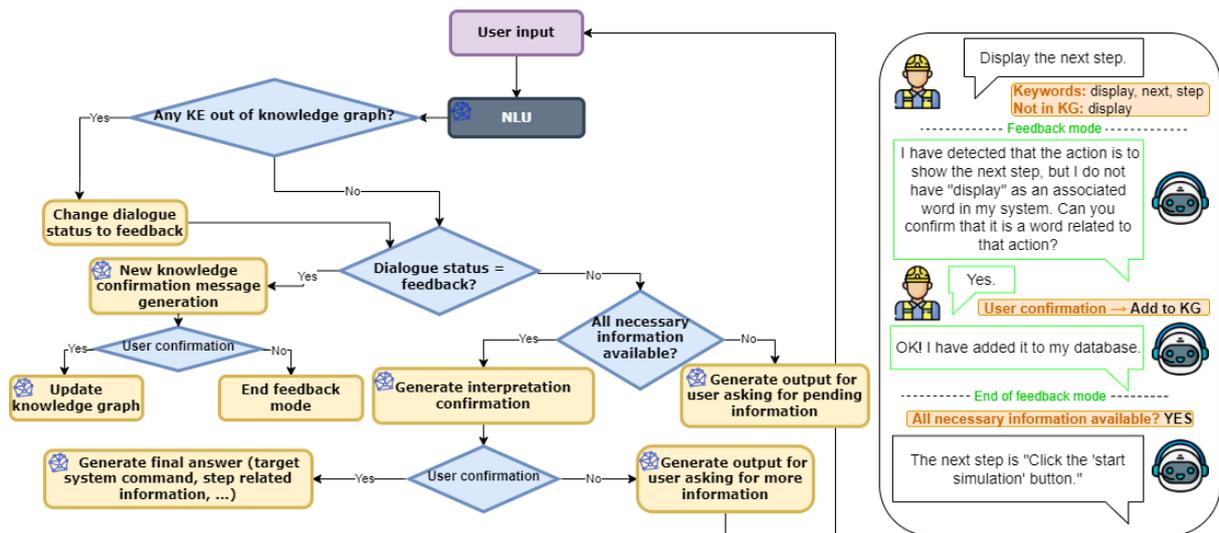


Figure 3: System workflow, including feedback management (left), along with an interaction example that requires feedback and follows the aforementioned flow (right).

are not able to understand it, the system enters *feedback mode*. The main goal of this mode is that the system is able to link new key elements to an action or action slot in subsequent interactions. For that, a clarifying request for the user, as a question, is triggered, in an intent to link the key element(s) extracted from the interpretation module with some of the classes/instances of the semantic repository. When the user responds to said system request, it is interpreted and the dialogue status is checked. If the status is feedback mode, once the user confirms the interpretation, the system launches a request to update the knowledge graph. This update, which has been automated by developing a REST API service, represents the extension of the base knowledge of the system. However, since it is an automated process, and to achieve controlled growth, this new knowledge is marked as obtained from feedback in the base (*provisional*, in accordance with what has been established in the adaptation of the TODO ontology, depicted in Section 3.1). Figure 3 visually summarizes the system’s dialogue flow, with the new feedback management capabilities to learn based on interactions with the user and update knowledge dynamically.

When it comes to the second case, in which the user’s disapproval of a system response is due to the fact that the content does not cover their needs, this feedback must trigger an action by an expert to review the system’s knowledge and update it if appropriate. For this case, a graphical interface has been developed so that it enables the user to indicate their disagreement with the content and the

expert to edit the content of the processes described in the repository when necessary. By the time this edition occurs, a functionality has been developed in the dialogue manager, which allows updating the knowledge graph with the new content. This new revised and improved data is what the system will use onwards as part of the extended knowledge graph.

### 3.3 Natural Language Understanding

The functions of the Natural Language Understanding (NLU) component are, first, to determine if a transcribed user voice command is classified as a polar interaction (e.g. “yes”, “no”). If it is, it is in charge of determining whether it is positive or negative. If no polarity is detected, a key element extraction (KEE) component is raised to extract the relevant information from the command, as shown in Figure 4.

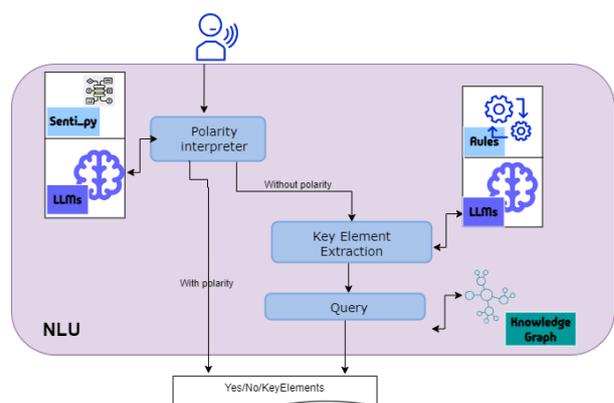


Figure 4: NLU pipeline

For the present work, two different approaches have been tested for the polarity interpreter. The first relies on the KIDE4I implementation, using a sentiment analysis algorithm. As for the second, an LLM-based approach has been implemented. More specifically, the GPT-3.5 Turbo model has been adapted through the in-context learning prompting strategy.

For the rest of interactions, namely non-polarized interactions, the KEE module intervenes fulfilling a slot-filling task. Similarly to the polarity interpreter, two different implementations have been used. For the first, again, the KIDE4I rule-based approach has been followed, whereas, for the second, the same LLM approach as above has been used. In this case, within the prompt, its function and the slots to be filled are indicated.

After detecting the slot values, it is essential to verify that these values correspond to recognizable world elements within the target system. This process involves querying the knowledge graph and comparing the detected slot values with those stored in the graph. If the key elements identified by the NLU are found, they are retained in the component's final result. Otherwise, the unrecognized values are discarded, asking for the missing information later on.

## 4 Initial Results

The extended TOD system with feedback capabilities has been tested and validated in an experimentation task. In this particular case, the system's function was to provide support through the different phases for manufacturing a piece using CNC programming on a milling machine. When given instructions, the users that were not satisfied with the answer given by the assistant would mark the response as inadequate. Some initial results related to this experimentation, mainly regarding the impact of those iterations in the knowledge graph, are presented in Section 4.1.

In terms of incrementing knowledge through feedback, the assistance scenario is suitable for evaluation. However, the variability in the interactions (e.g. "Show me the next step", "I need more information") is limited. Due to this, and in order to provide more insights, the LLM-based and rule-based NLU components have been tested in a collaborative bin-picking scenario, which is richer in terms of references to key elements in user commands. The obtained results for both have

been compared in order to determine if the LLM-based approach, which makes the task of adapting the dialogue system easier, is able to maintain or improve the accuracy of the rule-based approach. The experimental setup, as well as the benchmark results, are presented in Section 4.2.

### 4.1 Incremental TOD System in Use

The proposed validation scenario, as noted previously, is about CNC programming to manufacture a piece on an IKASMAK 5.1 milling machine. To do so, 15 users were requested to be assisted by the assistant described in Section 3 and to give insights about different procedures along the different phases of the programming process upon a given user request. The language used in this scenario was Spanish.

Furthermore, while the user is interacting with the system, if it does not present the desired information or does not perform as expected, the user can vote negatively the answer. This vote triggers a review alert for the expert, who will review the dialogue flow and, if necessary, update the knowledge graph to try to solve the gap, as shown in Figure 5.

This expert review, then, improves and, some times, even increases the information in the knowledge graph, and so, the accuracy of the system in further uses. The following subsection shows the evolution of the graph after the experimentation where, at some point within the dialogues, 22 of the total of 551 turns were marked as the response from the system was not valid, which triggered an expert review. Although it could sound like a high number, the time required for the review and update of the knowledge graph has not exceeded 5 hours, a process that would have taken much longer if done through other methods (e.g. manual instantiation) and would have required an ontology expert to perform it.

#### 4.1.1 Impact on the Knowledge Graph and Initial Analysis

So as to show the evolution of the graph before and after the expert review process, Tables 1 and 2, respectively, are depicted below.

In the case of Table 1, the average and total number of the different instances for each relevant class in procedure definition can be seen, out of a total of 341 total instances.

Thanks to the user feedback and expert review process, and as it can be seen in Table 2, the number of total instances has increased. More specifically,



Figure 5: User discontent triggering expert review requirement

Proc	AP	P	Met	Task	Step	PSI
10	14.9	1	1.3	1.3	6.9	8.7
<b>Total</b>	149	10	13	13	69	87

Table 1: Number of procedures (“Proc”) and average and total number of activation phrases (“AP”), procedures (“P”), methods (“Met”), tasks, steps and additional information (“PSI”), in the moment of the experimentation.

Proc	AP	P	Met	Task	Step	PSI
12	14.6	1	1.6	1.6	8	9.2
<b>Total</b>	175	12	19	19	96	110

Table 2: Number of procedures (“Proc”) and average and total number of activation phrases (“AP”), procedures (“P”), methods (“Met”), tasks, steps and additional information (“PSI”), after the user feedback and expert review.

90 more instances have been added, making a total of 431 (that is, a 25% more knowledge). Among these instances, new activation phrases have been added for the existing procedures and, furthermore, two new procedures have been included: “Detener un programa” (“Stop a program”) and “Configurar el avance” (“Configure the advance”). These two procedures have been added following the same format as the rest of procedures, an example of which can be found in Appendix A.

## 4.2 Natural Language Understanding: Benchmarking

The scenario used for the NLU component validation is a classification task, in which a bin-picking collaborative robot is able to classify cartridges by depositing them in different boxes, according to

user commands in Spanish. More specifically, the robot can pick up different ink cartridges from a table, identify their color and brand, and sort them into two separate containers, based on the operator’s instructions. The operator must use natural communication to inform the robot about the type of cartridge and the designated box. This communication involves not only voice commands but also gestures to indicate the destination. Consequently, the key element extraction module must identify actions and targets related to brands, colors (of the cartridges), and containers. Additionally, it must detect references to gestures indicated by phrases like “here” or “this”, which enhance the verbal instructions and provide supplementary information.

### 4.2.1 In-context Learning LLM vs Rule-based

In order to evaluate the behaviour of NLU in the different systems (rule-based and LLM-based), similar dialogues have been established with the same start of dialogue and the same end goal. In this way, they can be compared in number of turns and the performance of the NLU can be analysed. Therefore, based on these dialogues, the results of the KEE have been analysed for each turn.

A total of 74 dialogues were established with a total of 12 different users. However, the number of dialogue turns (159 for the rule-based and 176 for the LLM-based) and total KEE module intervention (130 and 108, respectively) varies due to the structure of the dialogue –which is slightly different for each system– and the performance of the different modules. The performance of both approaches can be observed in Table 3. In order to have a better approximation of the results, they have been classified between “fully detected”, “partially detected” and “not detected” to refer to when

	Fully detected		Partially detected		Not detected		Out-of-scope errors	
	%	#	%	#	%	#	%	#
Rule-based	64.61	84	17.69	23	13.07	17	4.61	6
LLM-based	98.14	106	0.92	1	0	0	0.92	1

Table 3: KEE results. Results are represented in percentages (%) and absolute numbers (#).

all, some, or none of the elements to be identified have been detected, respectively. Finally, it is worth mentioning that, due to out-of-scope causes, in both systems there have been elements that have been erroneously sent as input to the KEE, also presented in Table 3. These interactions, despite having had an output from the KEE, have not been taken into account in this analysis as NLU as they are caused by external errors.

All in all, we can observe a better performance of the LLM-based approach for key element extraction. More specifically, the LLM-based method outperforms the rule-based approach by a 33.5% in terms of fully detected key elements. Furthermore, the rule-based approach is more prone to partially detected and not detected elements, a situation with is rare in the LLM-based approach, with an only case of the former and no cases in the latter, which emphasises the capacity of these methods in this type of tasks.

As for the polarity component, the results have reported a 100% accuracy in both approaches.

## 5 Conclusions

This work introduces a knowledge graph-based method for managing the knowledge base of a task-oriented dialogue system for industrial settings, in which the knowledge graph is in charge of storing both domain and dialogue-management-related knowledge. This dialogue system features incremental learning capabilities that, by using the HIL/OIL paradigms, allows, on the one hand, for users to give feedback regarding the output of the system and, on the other hand, for experts to improve the knowledge included in the knowledge graph according to operators’ feedback. For this, the ontology used in the knowledge graph, which originates from an existing ontology for task-oriented dialogue systems, has been extended to cover the addition of knowledge and the generation of additional dialogues for that end.

Furthermore, for the natural language understanding (NLU) module, which originally was de-

signed by following a rule-based approach, has been implemented by using Large Language Models (LLMs) to improve both the system’s maintenance and the quality of the interpretations obtained by it.

The system has been evaluated in two real-world industrial settings: a bin-picking scenario, in which the NLU component was implemented by using LLMs, and a manufacturing scenario, in which the incremental learning capabilities of the system have been tested. For the first scenario, the results show that the performance of LLM-based NLU is higher than the rule-based approach by 16%, which is a significant improvement, especially for the fact that LLMs are easier to adapt to other scenarios than rules. For the second scenario, the addition of feedback interfaces has allowed to improve the existing knowledge graph of the system. The result of this is the addition of more explanations to existing procedures and even two new procedures; all in all, this translates into 25% more knowledge than at the time of the experimentation. This is expected to impact positively in the system’s performance from now on, which will be evaluated in a new experimentation task as part of future work.

These results show that the use of knowledge graphs for managing the knowledge base of task-oriented dialogue systems in industrial settings is a promising approach, especially when combined with incremental learning capabilities, and that the use of LLMs for other modules of the system leads to systems that are easy to maintain over time and to adapt to new scenarios.

## Acknowledgments

This project has received funding from the Department of Economic Development and Infrastructure of the Basque Government under grant number KK-2022/00102 (BERREKIN), KK-2023/00012 (BEREZ-IA) and KK-2024/00064 (IKUN).

## References

- János Abonyi, László Nagy, and Tamás Ruppert. 2024. *Knowledge Graph-Based Framework to Support the Human-Centric Approach*, pages 127–156. Springer Nature Switzerland, Cham.
- Cristina Aceta, Izaskun Fernández, and Aitor Soroa. 2021. TODO: A Core Ontology for Task-Oriented Dialogue Systems in Industry 4.0. In *Further with Knowledge Graphs*, pages 1–15. IOS Press.
- Cristina Aceta, Izaskun Fernández, and Aitor Soroa. 2022. KIDE4I: A generic semantics-based task-oriented dialogue system for human-machine interaction in industry 5.0. *Applied Sciences*, 12(3):1192.
- Vinsen Marselino Andreas, Genta Indra Winata, and Ayu Purwarianti. 2022. *A comparative study on language models for task-oriented dialogue systems*. CoRR, arXiv abs/2201.08687.
- Lang Cao. 2023. Diaggpt: An llm-based chatbot with automatic topic management for task-oriented dialogue. *arXiv preprint arXiv:2308.08043*.
- Songbo Hu, Xiaobin Wang, Zhangdie Yuan, Anna Korhonen, and Ivan Vulić. 2024. DIALIGHT: Lightweight Multilingual Development and Evaluation of Task-Oriented Dialogue Systems with Large Language Models. *arXiv preprint arXiv:2401.02208*.
- Vojtěch Hudeček and Ondřej Dušek. 2023. Are large language models all you need for task-oriented dialogue? In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 216–228.
- Chen Li, Xiaochun Zhang, Dimitrios Chrysostomou, and Hongji Yang. 2022. Tod4ir: A humanised task-oriented dialogue system for industrial robots. *IEEE Access*, 10:91631–91649.
- Bing Liu and Sahisnu Mazumder. 2021. Lifelong and continual learning dialogue systems: learning during conversation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 15058–15063.
- László Nagy, János Abonyi, and Tamás Ruppert. 2024. *Knowledge Graph-Based Framework to Support Human-Centered Collaborative Manufacturing in Industry 5.0*. *Applied Sciences*, 14(8).
- 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.
- Sinan Ozdemir. 2023. *Quick Start Guide to Large Language Models: Strategies and Best Practices for Using ChatGPT and Other LLMs*. Addison-Wesley Professional.
- María Poveda-Villalón, Alba Fernández-Izquierdo, and Raúl García-Castro. 2019. *Linked Open Terms (LOT) Methodology*. <https://doi.org/10.5281/zenodo.2539305>.
- Ivan Sekulić, Silvia Terragni, Victor Guimarães, Nghia Khau, Bruna Guedes, Modestas Filipavicius, André Ferreira Manso, and Roland Mathis. 2024. Reliable LLM-based User Simulator for Task-Oriented Dialogue Systems. *arXiv preprint arXiv:2402.13374*.
- Mohammed Sidi Yakoub, Sid-Ahmed Selouani, and Roger Nkambou. 2015. Mobile spoken dialogue system using parser dependencies and ontology. *International Journal of Speech Technology*, 18:449–457.
- Milene Santos Teixeira, Vinícius Maran, and Mauro Dragoni. 2021. The interplay of a conversational ontology and AI planning for health dialogue management. In *Proceedings of the 36th annual ACM symposium on applied computing*, pages 611–619.
- Mathilde Veron, Sahar Ghannay, Anne-Laure Ligozat, and Sophie Rosset. 2021. Lifelong learning and task-oriented dialogue system: what does it mean? In *Increasing Naturalness and Flexibility in Spoken Dialogue Interaction: 10th International Workshop on Spoken Dialogue Systems*, pages 347–356. Springer.

## A Example procedure

Listing 1: Snippet of the instances in the “Editar o modificar un programa” (“Edit or modify a procedure”) procedure. This example is presented in TTL format for readability.

```
1 [...]
2
3 ### https://w3id.org/todo/tododw-ekin-inst#
4   ↪ Method0_Procedure_65113d9f5d9c3075571719b5
5 :Method0_Procedure_65113d9f5d9c3075571719b5 rdf:type owl:NamedIndividual ,
6   tododwHowto:Method ;
7   var:hasFirstSegment :Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
8   var:isMadeOf :Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
9   var:processSegmentId "0" .
10
11 ### https://w3id.org/todo/tododw-ekin-inst#
12   ↪ PSIO_Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5
13 :PSIO_Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 rdf:type owl:
14   ↪ NamedIndividual ,
15   <http://www.mesa.org/xml/B2MML-V0600#ProcessSegmentInformation> ;
16   var:relatedToProcessSegment :
17     ↪ Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
18   var:taskImage "https://server/editar-programa/Metodo1-Paso2.png" ;
19   tododwHowto:index "0" .
20
21 ### https://w3id.org/todo/tododw-ekin-inst#Procedure_65113d9f5d9c3075571719b5
22 :Procedure_65113d9f5d9c3075571719b5 rdf:type owl:NamedIndividual ,
23   tododwHowto:Procedure ;
24   var:description "Editar o modificar un programa" ;
25   var:hasFirstSegment :Method0_Procedure_65113d9f5d9c3075571719b5 ;
26   var:isMadeOf :Method0_Procedure_65113d9f5d9c3075571719b5 ;
27   var:processSegmentId "65113d9f5d9c3075571719b5" .
28
29 ### https://w3id.org/todo/tododw-ekin-inst#
30   ↪ Step0_Task0_Method0_Procedure_65113d9f5d9c3075571719b5
31 :Step0_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 rdf:type owl:NamedIndividual
32   ↪ ,
33   tododwHowto:Step ;
34   var:description "Abrir el programa deseado. " ;
35   var:isPrevious :Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
36   var:processSegmentId "0" ;
37   tododwHowto:hasAssociatedProcedure :Procedure_65113a0c5d9c3075571719b2 .
38
39 ### https://w3id.org/todo/tododw-ekin-inst#
40   ↪ Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5
41 :Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 rdf:type owl:NamedIndividual
42   ↪ ,
43   tododwHowto:Step ;
44   var:description "Con el programa en pantalla , tal y como se muestra en la
45     ↪ siguiente imagen, se podrá comenzar a modificar o extender el código G
46     ↪ para programar la pieza." ;
47   var:hasRelatedInformation :
48     ↪ PSIO_Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
49   var:isNext :Step0_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
50   var:isPrevious :Step2_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
51   var:processSegmentId "1" .
52
53 ### https://w3id.org/todo/tododw-ekin-inst#
54   ↪ Step2_Task0_Method0_Procedure_65113d9f5d9c3075571719b5
55 :Step2_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 rdf:type owl:NamedIndividual
56   ↪ ,
57   tododwHowto:Step ;
58   var:description "Para guardar los cambios , no es necesaria ninguna acción especí
59     ↪ fica: se guarda automáticamente." ;
60   var:isNext :Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
61   var:processSegmentId "2" .
```

```
52
53 ### https://w3id.org/todo/tododw-ekin-inst#
54 ↪ Task0_Method0_Procedure_65113d9f5d9c3075571719b5
55 :Task0_Method0_Procedure_65113d9f5d9c3075571719b5 rdf:type owl:NamedIndividual ,
56   tododwHowto:Task ;
57   var:hasFirstSegment :Step0_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
58   var:isMadeOf :Step0_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ,
59               :Step1_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ,
60               :Step2_Task0_Method0_Procedure_65113d9f5d9c3075571719b5 ;
61   var:processSegmentId "0" .
62 [...]
```

# Anticipating Follow-Up Questions in Exploratory Information Search

Graham Wilcock

CDM Interact and University of Helsinki

Helsinki, Finland

graham.wilcock@cdminteract.com

## Abstract

The paper describes methods for anticipating follow-up questions in exploratory information search. There are two main cases: information stored in knowledge graphs, and information in unstructured texts such as Wikipedia. In the first case, follow-up questions are anticipated by extracting subgraphs relevant to user queries, passing the subgraphs to an LLM to generate responses. In the second case, entities and their relationships are extracted from the texts and added to short-term knowledge graphs relevant to initial queries. Follow-up questions are then anticipated by extracting subgraphs relevant to subsequent queries and passing the subgraphs to the LLM, as in the first case. The short-term graphs in dialogue memory are often sufficient to answer follow-up questions. If they are not, the described steps are repeated as required.

## 1 Introduction

Dialogue systems that support users in exploratory information search typically need to handle many follow-up questions. The paper describes methods for anticipating follow-up questions in dialogues for exploratory information search. There are two cases: exploring information stored in knowledge graphs, and exploring information in unstructured texts such as Wikipedia.

The dialogues are exploratory because the users do not yet know where the information is located, or even if it exists. They may not know the structure of the knowledge graphs, or what taxonomy has been used to classify the information into different categories. As a result, users need to keep asking questions as they learn to navigate around different information spaces.

The proposed approach aims to anticipate likely follow-up questions by constructing subgraphs of entities and relationships relevant to current and recent user queries. This can be done while the user is thinking what question to ask next.

If a user is searching existing knowledge graphs, likely follow-up questions can be anticipated by extracting subgraphs relevant to the current user query. The subgraphs are included in prompts to LLMs to generate responses to the user.

If a user is searching unstructured texts such as Wikipedia, there is no knowledge graph from which subgraphs can be extracted. In this case an LLM is prompted to extract entities from the user query, and to extract relevant entities and relationships from the texts, and finally to construct a small short-term knowledge graph from them.

The paper is structured as follows. Section 2 discusses related work. Section 3 summarizes existing methods for generating natural language responses from Wikipedia texts and from knowledge graphs. Section 4 describes new methods for generating subgraphs from existing knowledge graphs and for generating new knowledge graphs from texts. In Section 5 the new methods are used to anticipate follow-up questions in a hybrid retrieval approach combining structured and unstructured retrieval.

## 2 Related Work

Hogan et al. (2022) is a comprehensive guide to knowledge graphs. Schneider et al. (2022) survey the increasing use of knowledge graphs in NLP.

Sarkar et al. (2020) study methods for extracting subgraphs from DBpedia for use in conversational recommender systems. This is similar to subgraph extraction from knowledge graphs stored in Neo4j graph databases, described in Section 4.1.

A system combining conversational agents with knowledge graphs in Neo4j databases is described by Wilcock and Jokinen (2022). A similar system from Schneider et al. (2023b) aims for synergy between knowledge graphs and conversational agents by bridging the gap between structured and unstructured information retrieval, a topic also addressed here in Section 5 on hybrid retrieval.

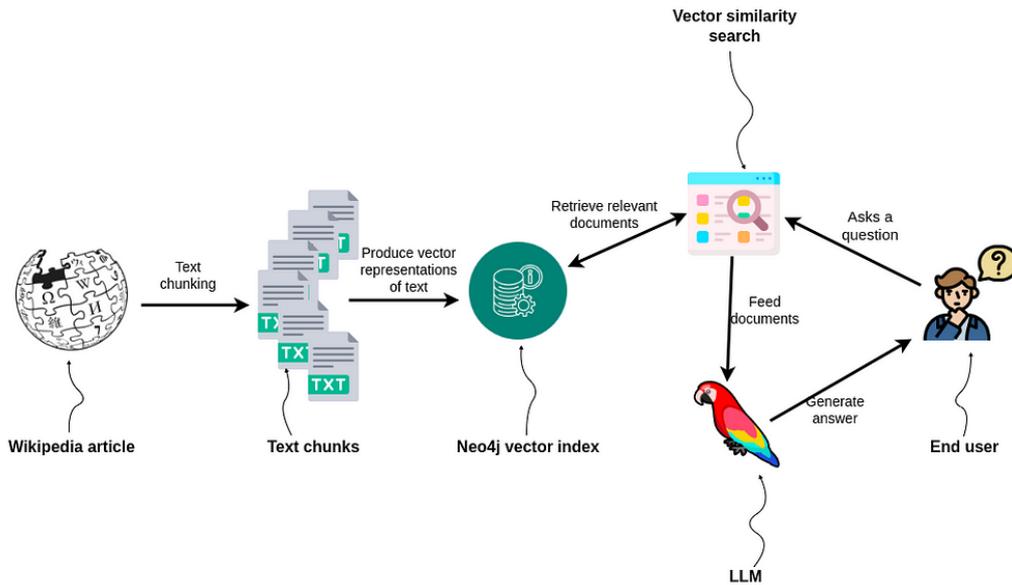


Figure 1: Simple RAG from Wikipedia texts. Image by Tomaz Bratanic, from (Bratanic, 2023a).

Concerning methods for anticipating follow-up questions in exploratory search, Schneider et al. (2023b) mention WikiTalk (Wilcock, 2012), an early robot dialogue system for exploratory search in Wikipedia. Using no knowledge graphs, WikiTalk extracted sets of hyperlinks from Wikipedia articles to transition smoothly between topics by anticipating what the user will ask about next.

Jokinen and Wilcock (2016) proposed a method for anticipation of follow-up topics in Wikipedia search based on hyperlinks and keywords extracted from the current article. This enables anticipating follow-up topics that have no explicit link, and also works for documents without hyperlinks.

The WikiTalk approach of extracting small sets of Wikipedia hyperlinks from the current topic to related topics was motivated by the need at that time to restrict speech recognition vocabulary to a finite list of predicted phrases (Wilcock, 2012). However, the basic idea is similar to retrieving a subgraph or neighborhood of relevant nodes from a knowledge graph, as described in Section 4.1.

RAG (Lewis et al., 2020) is often described as a way of *grounding* LLM responses in the retrieved information, but *conversational* grounding has a long history in dialogue systems research (Traum, 1995; Jokinen, 1996). Grounding is especially important in open-ended conversational exploratory search for navigation in unknown information landscapes (Schneider et al., 2023a).

Theory of Mind errors often arise from failure to build shared knowledge during the dialogue (Wilcock and Jokinen, 2023). Jokinen et al. (2024) investigate the capacity of LLMs to build shared knowledge by classifying grounding-related dialogue acts and by extracting mutually grounded information.

### 3 LLMs that Generate Responses

RAG enables LLMs to generate natural language responses from retrieved information that is not in their training corpora. This section compares existing methods for RAG from Wikipedia texts and RAG from knowledge graphs.

#### 3.1 Simple RAG from Wikipedia texts

Figure 1 shows a simple RAG application described by Bratanic (2023a) that answers questions based on information from Wikipedia. For a given topic, Wikipedia articles are downloaded and split into text chunks using LangChain. Vector embeddings of the chunks are generated and stored in a Neo4j database with the texts.

When users ask questions, embeddings of the questions are generated and the most relevant chunks are found by semantic similarity using a Neo4j vector index. The questions and the most relevant chunks are passed to an LLM to generate the answers. Follow-up questions are enabled by using LangChain memory components.

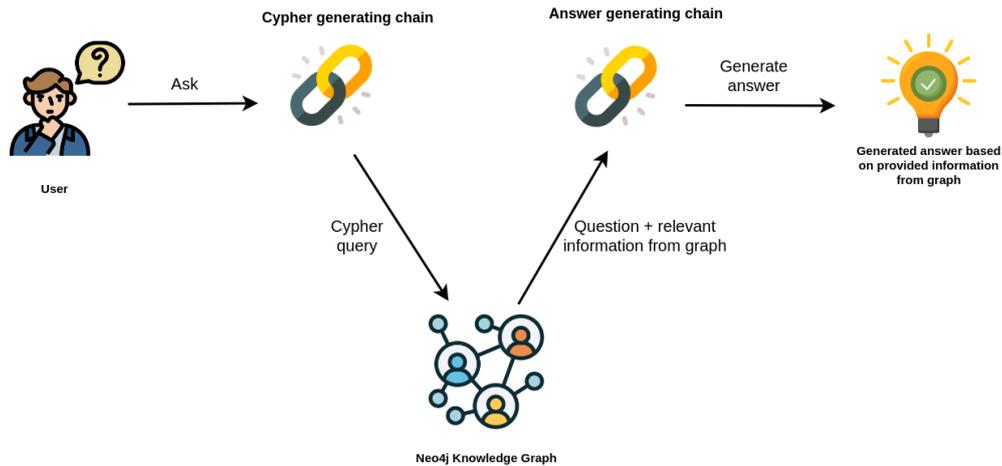


Figure 2: RAG from Knowledge Graphs. Image by Tomaz Bratanic, from (Bratanic, 2023b).

### 3.2 When simple RAG goes wrong

In order to reduce hallucinations, LLMs can be prompted to avoid making up false facts by using only the information given in the context. However, this can sometimes cause LLMs to avoid telling true facts, by answering as if the facts given in the context are the only true facts in the world.

An example is given by Wilcock (2024), from a *Chat with Wikipedia* application that was given the topic *William Shakespeare*. When asked the question *Did he have any children, grandchildren or other descendants?* the set of most relevant chunks retrieved by RAG did not mention his child Susanna. This caused a conflict between the LLM’s own knowledge of Shakespeare and the instructions to use only the information given in the context.

To resolve this conflict, the LLM gave a correct but misleading reply *Yes, William Shakespeare had at least two known children*. The absence of his child Susanna from the context caused the LLM to invent a false justification *There is no direct evidence that he had any other children*.

The follow-up question *Who was Susanna Shakespeare?* caused a new set of chunks to be retrieved and the LLM replied *Susanna Shakespeare was the daughter of William Shakespeare and his wife Anne Hathaway*. It then contradicted its previous reply by adding *Susanna is one of three children known to have been born to Shakespeare and his wife*.

### 3.3 RAG from knowledge graphs

Recently Neo4j graph databases have been widely used to manage knowledge graphs (Barrasa and

Webber, 2023). RAG applications can retrieve information from Neo4j knowledge graphs using Cypher database queries.

Figure 2 from (Bratanic, 2023b) shows RAG from knowledge graphs using two LLMs. The first LLM generates database query code based on the user question. The query retrieves relevant information from the knowledge graph. The second LLM uses the question and the retrieved information to generate the response to the user.

An advantage of RAG from knowledge graphs is that semantic metadata such as taxonomies can be added to the graphs and used to generate more intelligent responses. An example of using knowledge graph metadata in a dialogue system is given by Wilcock (2024). When a user asks for restaurants that serve European cuisine, the graph query finds restaurants serving Italian cuisine. As a taxonomy of cuisines from WikiData was added to the graph, the RAG retrieves the Italian restaurants because Italian cuisine is a subclass of European cuisine in the taxonomy. The LLM gives an intelligent response, explaining that the restaurants serve Italian cuisine which is a type of European cuisine.

## 4 LLMs that Generate Graphs

We now describe methods for generating subgraphs from existing knowledge graphs and for generating new knowledge graphs from texts.

### 4.1 Generating subgraphs from graphs

A graph retriever function (Bratanic, 2024) that extracts subgraphs from knowledge graphs in Neo4j graph databases is shown in Figure 3.

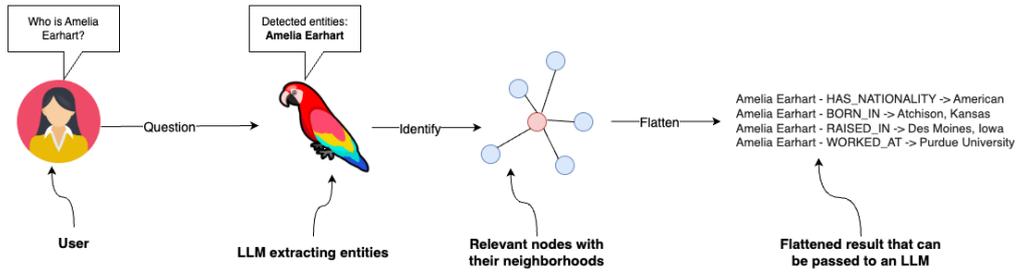


Figure 3: A graph retriever function. Image by Tomaz Bratanić, from (Bratanić, 2024)

The function first extracts entities from the user query. Next, it iterates over the detected entities and uses a Cypher template to retrieve a neighborhood of relevant nodes. The subgraph of relationships between these nodes is converted to a flattened text format that can be passed to an LLM to generate a natural language response to the user.

#### 4.2 Generating knowledge graphs from texts

LLMs can help with knowledge graph construction by analyzing unstructured texts and generating new structured data from them. LLMs must identify the entities mentioned in the texts and identify the relationships between them. They must generate code to create entities and relationships as nodes and relationships in the knowledge graph.

LLMGraphTransformer (Bratanić, 2024) helps to construct a knowledge graph by using an LLM to convert texts into graph documents, which can then be imported into Neo4j graph databases. Links to the sources of the texts can be included in the graph documents for provenance checking.

Bratanić (2024) introduces a hybrid approach to retrieval that aims to enhance RAG accuracy by combining vector-based search of unstructured text with structured retrieval of knowledge graph data. The new approach is shown in Figure 4.

To demonstrate the hybrid approach, Bratanić uses LLMGraphTransformer to extract entities and relationships from Wikipedia texts about Elizabeth I, convert the texts to graph documents, and import them into a knowledge graph in a Neo4j database.

Elizabeth I - RULED -> England
Elizabeth I - RULED -> Ireland
Elizabeth I - BELONGED_TO -> House Of Tudor
Elizabeth I - PARENT -> Henry Viii
Elizabeth I - PARENT -> Anne Boleyn

Table 1: Generated relationships about Elizabeth I.

Table 1 shows some examples of relationships generated when a user asks about Elizabeth I. The

graph retriever has converted the relationships into a flattened text format that is passed in a prompt to an LLM to generate responses to the user.

## 5 Anticipating Follow-up Questions in Hybrid Retrieval for RAG

The hybrid retrieval approach has been applied to anticipating follow-up questions in exploratory information search from Wikipedia. This short paper does not include evaluation and analysis sections. However, the hybrid retrieval approach can be compared with the simple RAG from Wikipedia texts described in Section 3.1.

In a small trial, LLMGraphTransformer was used to construct a knowledge graph about William Shakespeare from Wikipedia texts. Appendix A lists 50 of the generated relationships. Table 2 shows some examples.

William Shakespeare - SPOUSE -> Anne Hathaway
William Shakespeare - PARENT -> Susanna
William Shakespeare - PARENT -> Hamnet
William Shakespeare - PARENT -> Judith
William Shakespeare - PART_OWNER -> Lord Chamberlain's Men

Table 2: Generated relationships about Shakespeare.

When asked the follow-up questions in Table 3, the hybrid retrieval gives the correct information about Shakespeare's children, unlike the error in Section 3.2 with simple RAG from texts.

Question: When was he born?
Assistant: <i>William Shakespeare was born on 23 April 1564.</i>
Question: Did he have any children, grandchildren, or descendants?
Assistant: <i>Yes, William Shakespeare had three children: Susanna, and twins Hamnet and Judith.</i>

Table 3: Follow-up questions about Shakespeare.

There are some technical problems in generating knowledge graphs using LLMs. One problem is

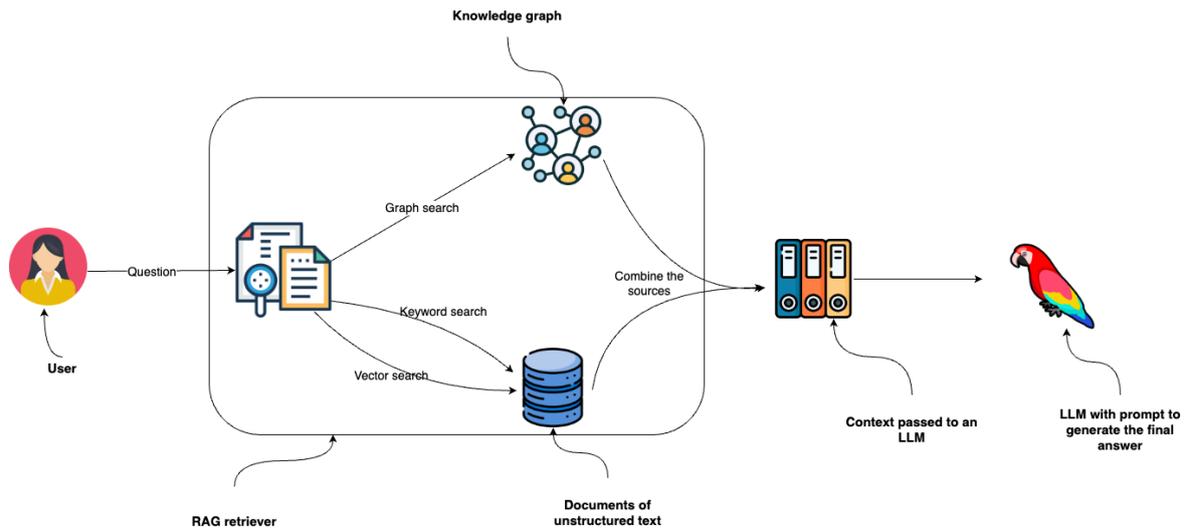


Figure 4: Hybrid Retrieval for RAG. Image by Tomaz Bratanić, from (Bratanić, 2024).

getting the direction of relationships correct. For example in Table 1, PARENT relationships go from Elizabeth I to her parents Henry VIII and Anne Boleyn, but in Table 2, PARENT relationships go from Shakespeare to his children Susanna, Hamnet and Judith. Work to resolve this problem is ongoing.

## 6 Conclusion

After summarizing existing methods for generating natural language responses from Wikipedia texts and from knowledge graphs, the paper described new methods for anticipating follow-up questions in dialogues for exploratory information search, considering two cases. When exploring information already stored in knowledge graphs, follow-up questions are anticipated by extracting subgraphs that are likely to be relevant to subsequent user queries, and passing the subgraphs to an LLM to generate responses.

When exploring information in unstructured texts such as Wikipedia, entities and relationships are extracted from the texts and used to construct new short-term knowledge graphs relevant to initial user queries. Follow-up questions are anticipated by extracting subgraphs likely to be relevant to subsequent user queries, and continuing as in the first case.

Although there are some problems to be solved in automatic construction of knowledge graphs by LLMs, this kind of approach is attractive. Ongoing work will aim to explore its potential benefits both for anticipating follow-up questions in exploratory

information search, and more widely in other areas of spoken dialogue systems.

## Acknowledgements

The author thanks Kristiina Jokinen for valuable discussions and fruitful collaboration on dialogue-related topics.

## References

- Jesús Barrasa and Jim Webber. 2023. *Building Knowledge Graphs: A Practitioner’s Guide*. O’Reilly Media.
- Tomaz Bratanić. 2023a. LangChain library adds full support for Neo4j vector index. <https://neo4j.com/developer-blog/langchain-library-full-support-neo4j-vector-index/>.
- Tomaz Bratanić. 2023b. neo4j\_cypher. <https://github.com/langchain-ai/langchain/tree/master/templates/neo4j-cypher>.
- Tomaz Bratanić. 2024. Enhancing RAG-based applications accuracy by constructing and leveraging knowledge graphs. [https://github.com/tomasonjo/blob/master/llm/enhancing\\_rag\\_with\\_graph.ipynb](https://github.com/tomasonjo/blob/master/llm/enhancing_rag_with_graph.ipynb).
- Aidan Hogan, Eva Blomqvist, Michael Cochez, Cláudia d’Amato, Gerard de Melo, Claudio Gutiérrez, Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian Neumaier, Axel-Cyrille Ngonga Ngomo, Axel Polleres, Sabbir M. Rashid, Anisa Rula, Lukas Schmelzeisen, Juan F. Sequeda, Steffen Staab, and Antoine Zimmermann. 2022. *Knowledge Graphs*. Morgan & Claypool.
- Kristiina Jokinen. 1996. Cooperative Response Planning in CDM: Reasoning about Communicative Strategies. In Anton Nijholt, editor, *Twente Workshop Series in Language Technology*.
- Kristiina Jokinen, Phillip Schneider, and Taiga Mori. 2024. Towards harnessing large language models for comprehension of conversational grounding. In *14th International Workshop on Spoken Dialogue Systems Technology (IWSDS 2024)*, Sapporo, Japan.
- Kristiina Jokinen and Graham Wilcock. 2016. Double topic shifts in open domain conversations: Natural language interface for a Wikipedia-based robot application. In *Proceedings of the Open Knowledge Base and Question Answering Workshop (OKBQA 2016)*, pages 59–66, Osaka, Japan. The COLING 2016 Organizing Committee.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS 2020)*, pages 9459–9474, Vancouver, Canada.
- Rajdeep Sarkar, Koustava Goswami, Mihael Arcan, and John P. McCrae. 2020. Suggest me a movie for tonight: Leveraging knowledge graphs for conversational recommendation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4179–4189, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Phillip Schneider, Anum Afzal, Juraj Vladika, Daniel Braun, and Florian Matthes. 2023a. Investigating conversational search behavior for domain exploration. In *Advances in Information Retrieval: 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2–6, 2023, Proceedings, Part II*, page 608–616, Berlin, Heidelberg. Springer-Verlag.
- Phillip Schneider, Nils Rehtanz, Kristiina Jokinen, and Florian Matthes. 2023b. From data to dialogue: Leveraging the structure of knowledge graphs for conversational exploratory search. In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, pages 609–619, Hong Kong, China. Association for Computational Linguistics.
- Phillip Schneider, Tim Schopf, Juraj Vladika, Mikhail Galkin, Elena Simperl, and Florian Matthes. 2022. A decade of knowledge graphs in natural language processing: A survey. 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 601–614, Online only. Association for Computational Linguistics.
- David R. Traum. 1995. *A computational theory of grounding in natural language conversation*. Ph.D. thesis, University of Rochester, USA.
- Graham Wilcock. 2012. WikiTalk: A spoken Wikipedia-based open-domain knowledge access system. In *Proceedings of the COLING 2012 Workshop on Question Answering for Complex Domains*, pages 57–69, Mumbai, India.
- Graham Wilcock. 2024. New technologies for spoken dialogue systems: LLMs, RAG and the GenAI Stack. In *14th International Workshop on Spoken Dialogue Systems Technology (IWSDS 2024)*, Sapporo, Japan.
- Graham Wilcock and Kristiina Jokinen. 2022. Conversational AI and knowledge graphs for social robot interaction. In *ACM/IEEE International Conference on Human-Robot Interaction (HRI 2022)*, pages 1090–1094, Sapporo, Japan. Association for Computing Machinery.
- Graham Wilcock and Kristiina Jokinen. 2023. To Err Is Robotic; to Earn Trust, Divine: Comparing ChatGPT and Knowledge Graphs for HRI. In *32nd IEEE International Conference on Robot and Human Interactive Communication (RO-MAN 2023)*, pages 1396–1401, Busan, Korea.

## A Appendix A

Relationships relevant to William Shakespeare extracted by LLMGraphTransformer from Wikipedia texts and imported into a Neo4j knowledge graph. They are shown in a flattened text format that can be passed in prompts to LLMs.

Only 50 relationships are listed here.

William Shakespeare - SPOUSE -> Anne Hathaway  
William Shakespeare - PARENT -> Susanna  
William Shakespeare - PARENT -> Hamnet  
William Shakespeare - PARENT -> Judith  
William Shakespeare - PART\_OWNER -> Lord Chamberlain'S Men  
Lord Chamberlain'S Men - NAME\_CHANGE -> King'S Men  
King James Vi Of Scotland - ASCENSION -> King'S Men  
William Shakespeare - FRIEND -> John Heminges  
William Shakespeare - FRIEND -> Henry Condell  
Shakespeare - FAMILY -> John Shakespeare  
Shakespeare - FAMILY -> Mary Arden  
Shakespeare - MARRIAGE -> Anne Hathaway  
Shakespeare - ACQUAINTANCE -> Ben Jonson  
Shakespeare - ACQUAINTANCE -> William Oldys  
Shakespeare - ACQUAINTANCE -> George Steevens  
Shakespeare - AUTHOR -> Plays  
William Shakespeare - AUTHOR -> Plays  
Plays - CLASSIFICATION -> Tragedy  
Plays - CLASSIFICATION -> History  
Plays - CLASSIFICATION -> Comedy  
Plays - CLASSIFICATION -> Problem Plays  
Plays - CLASSIFICATION -> Romances  
Shakespeare - ARRIVAL -> London  
Shakespeare - INVOLVEMENT -> The Curtain  
Tudor Morality Plays - INFLUENCE -> Shakespeare  
Classical Aesthetic Theory - INFLUENCE -> Shakespeare  
Classical Aesthetic Theory - DERIVED\_FROM -> Aristotle  
Classical Aesthetic Theory - DERIVED\_FROM -> Plautus  
Classical Aesthetic Theory - DERIVED\_FROM -> Terence  
Rose - SIMILARITY -> Globe  
Public Theatres - HAS\_FEATURE -> Three Stories High  
Public Theatres - HAS\_FEATURE -> Open Space At The Center  
Public Theatres - HAS\_FEATURE -> Polygonal In Plan  
Public Theatres - HAS\_FEATURE -> Inward-Facing Galleries  
Public Theatres - HAS\_FEATURE -> Stage  
Stage - SURROUNDED\_BY -> Platform  
Platform - SURROUNDS -> Audience  
Stage - HAS\_FEATURE -> Rear  
Rear - HAS\_FEATURE -> Entrances And Exits  
Entrances And Exits - USED\_BY -> Actors  
Entrances And Exits - USED\_BY -> Musicians  
Public Theatres - HAS\_FEATURE -> Upper Level  
Upper Level - CAN\_BE\_USED\_AS -> Balcony  
Public Theatres - MADE\_OF -> Timber  
Public Theatres - MADE\_OF -> Lath And Plaster  
Public Theatres - HAS\_FEATURE -> Thatched Roofs  
Public Theatres - VULNERABLE\_TO -> Fire  
Public Theatres - REPLACED\_BY -> Globe  
Globe - REPLACED\_WITH -> Tile Roof  
Blackfriars Theatre - ASSOCIATED\_WITH -> Shakespeare

# Bridging Information Gaps in Dialogues With Grounded Exchanges Using Knowledge Graphs

Phillip Schneider<sup>1</sup>, Nektarios Machner<sup>1</sup>, Kristiina Jokinen<sup>2</sup>, and Florian Matthes<sup>1</sup>

<sup>1</sup>Technical University of Munich, Department of Computer Science, Germany

<sup>2</sup>National Institute of Advanced Industrial Science and Technology, AI Research Center, Japan  
{phillip.schneider, nektarios.machner, matthes}@tum.de  
kristiina.jokinen@aist.go.jp

## Abstract

Knowledge models are fundamental to dialogue systems for enabling conversational interactions, which require handling domain-specific knowledge. Ensuring effective communication in information-providing conversations entails aligning user understanding with the knowledge available to the system. However, dialogue systems often face challenges arising from semantic inconsistencies in how information is expressed in natural language compared to how it is represented within the system’s internal knowledge. To address this problem, we study the potential of large language models for conversational grounding, a mechanism to bridge information gaps by establishing shared knowledge between dialogue participants. Our approach involves annotating human conversations across five knowledge domains to create a new dialogue corpus called *BridgeKG*. Through a series of experiments on this dataset, we empirically evaluate the capabilities of large language models in classifying grounding acts and identifying grounded information items within a knowledge graph structure. Our findings offer insights into how these models use in-context learning for conversational grounding tasks and common prediction errors, which we illustrate with examples from challenging dialogues. We discuss how the models handle knowledge graphs as a semantic layer between unstructured dialogue utterances and structured information items.

## 1 Introduction

Conversational grounding is an integral aspect of dialogues where interlocutors share information and build up a common understanding. This mutually established knowledge serves as context for subsequent interactions. For building effective dialogue systems, the natural language processing (NLP) community has long focused on conversational grounding, which involves inferential reasoning, dynamic feedback, and repair strategies (Udagawa

and Aizawa, 2021). Despite extensive research, challenges remain in adapting to different conversation domains, addressing semantic vocabulary mismatches, overcoming information gaps between user knowledge and the system’s internal knowledge model, as well as the lack of appropriate training data (Lemon, 2022). Owing to rapid technical advances regarding large language models (LLMs), novel opportunities arise to comprehend contextual intricacies within dialogues and reconcile information expressed in natural language with that stored in machine-readable data structures.

Recognizing the limited research on LLM-based conversational grounding, we investigated the capabilities of LLMs on knowledge grounding tasks. This involved annotating an existing corpus containing dialogues about different domain-specific tabular datasets. In addition to labeling grounding acts, we annotated grounded knowledge items in a knowledge graph structure, a powerful representation of complex relationships between entities and their attributes. Knowledge graphs have proven valuable in various NLP tasks, such as disambiguating ambiguous utterances by providing contextual information (Hogan et al., 2021; Schneider et al., 2022). For example, in dialogue systems, knowledge graphs can help identify the correct meaning of a word with multiple senses or resolve references to specific entities, enhancing the overall understanding and coherence of conversations. We opted for the JSON-LD format due to its simplicity and acceptance as a web standard, allowing interoperability by reusing existing namespaces with shared vocabularies to model knowledge from different sources and domains.

While JSON-LD primarily uses a tree-like structure, it can represent more complex graph structures by linking nodes using identifiers like *@id* and *@type*. As a serialization format for Resource Description Framework (RDF) data, JSON-LD can be transformed into other formats, such as

N-Triples, RDF/XML, or Turtle. This flexibility allows JSON-LD to be integrated with graph databases and other RDF tools, enhancing its utility in various applications. Table 1 shows an example annotation of grounded knowledge in JSON-LD format from a conversation about nature parks.

Our contributions include (1) creating a novel dialogue corpus called *BridgeKG* with over 250 conversational grounding annotations across five knowledge domains, (2) conducting a range of zero- and few-shot experiments by evaluating four LLMs on two grounding tasks, and (3) summarizing common prediction errors and prompting techniques for improving model performance. To ensure the reproducibility of our experiments, we provide the *BridgeKG* dataset, source code, and evaluation outputs in a public GitHub repository.<sup>1</sup>

## 2 Related Work

In regard to the literature on grounding in NLP, it is essential to first define the broadly used term. Grounding can be categorized into three main types. Conversational grounding ensures a common understanding of shared knowledge within a conversation (Traum, 1994). Perceptual grounding links language to sensory experiences of the real world like visual information (Cangelosi, 2010). Knowledge grounding incorporates external information sources to support NLP systems, such as providing factual knowledge to generative language models (Lewis et al., 2020).

Our study focuses solely on conversational grounding by employing LLMs, a topic addressed in only a few recent studies. One related work by Shaikh et al. (2024) examines whether LLM generations contain grounding acts, simulating turn-taking from various conversation datasets. They found that LLMs generate language with less conversational grounding than humans, often producing text that appears to assume common ground. Both their study and ours focus on the three grounding acts: explicit grounding, implicit grounding, and clarification, as proposed by Clark and Schaefer (1989). Two other closely related studies, conducted by Jokinen et al. (2024) and Mohapatra et al. (2024), involve annotating dialogue corpora and employing language models to classify grounding acts and extract grounded knowledge items. While the former conducts preliminary experiments on two conversations with GPT-3.5-Turbo, the lat-

ter presents two annotated dialogue corpora with grounding acts, grounding units, a measure of their degree of grounding, and a baseline evaluation with the open-source T5 model (Raffel et al., 2020).

Unlike the mentioned related work, we are the first to conduct a series of LLM experiments aimed at knowledge identification in information-seeking conversations utilizing an in-context knowledge graph structure for identifying referenced and grounded knowledge items in dialogues.

## 3 Method

**Dataset Annotation** The source dialogue corpus we reuse was collected in a study on exploratory information-seeking conversations from Schneider et al. (2023). It comprises 26 conversations about tabular datasets on real-world knowledge spanning the domains of geography, history, media, nutrition, and sports. Every conversation involved a pair where one person was the information seeker and the other was the information provider, using a text-based chatroom for communication. The information seekers were instructed to discover and gather new information about their partner’s previously unknown dataset. Two researchers annotated each written dialogue with labels for grounding acts (explicit, implicit, and clarification). Explicit grounding involves a response that clearly confirms understanding or acceptance of received information (e.g., “okay, thanks”), whereas implicit grounding moves the conversation forward without explicitly acknowledging or questioning the recently shared information (implicit acceptance). Clarification occurs when a conversation partner seeks more information about thus far presented knowledge, which does not result in grounded knowledge since mutual acceptance has not yet been reached.

### Example Annotation of Grounded Knowledge

```
{["@context": "http://www.w3.org/ns/csvw", "@schema": "http://schema.org"}, {"@id": "http://example.org/nature-parks", "url": "nature-parks.csv", "schema:description": "The table contains information about nature parks in Germany", "tableSchema": {"columns": [{"name": "name", "datatype": "string"}, {"name": "state", "datatype": "string"}, {"name": "year", "datatype": "integer"}, {"name": "area_in_km2", "datatype": "integer"}, {"name": "summary", "datatype": "string"}], "primaryKey": "name"}}, {"@type": "schema:Place", "name": "Barnim", "state": "Brandenburg Berlin", "year": 1999, "area_in_km2": 749, "summary": "The park includes the Barnim heath habitats dating back to the ice age. It lies between the glacial valleys of Eberswalde in the north and Berlin in the south, and is more than half forested. The region is shaped by many individual lakes and meltwater gullies."}]
```

Table 1: Example JSON-LD annotation of grounded knowledge from the *BridgeKG* dataset, representing the system’s knowledge concerning a dialogue about nature parks. Properties are displayed in blue color.

<sup>1</sup>[github.com/philotron/Bridge-KG](https://github.com/philotron/Bridge-KG)

Model	Zero-Shot Prompt				Few-Shot Prompt			
	Accuracy	Precision	Recall	F1-Score	Accuracy	Precision	Recall	F1-Score
GPT-3.5-Turbo (n=1)	0.64	0.50	0.46	0.43	0.55	0.50	0.51	0.50
GPT-3.5-Turbo (n=3)	0.66	<b>0.81</b>	0.50	0.50	0.69	0.59	0.54	0.54
GPT-3.5-Turbo (n=all)	0.59	0.39	0.44	0.41	0.57	0.51	0.45	0.45
GPT-4o (n=1)	0.39	0.55	0.54	0.42	0.64	0.66	0.64	0.61
GPT-4o (n=3)	0.59	0.66	<b>0.67</b>	0.59	0.73	<b>0.74</b>	0.69	<b>0.70</b>
GPT-4o (n=all)	0.64	0.68	0.66	0.62	0.71	0.73	0.67	0.67
Llama-3-8B (n=1)	0.61	0.54	0.53	0.54	0.59	0.65	0.69	0.59
Llama-3-8B (n=3)	0.65	0.60	0.60	0.60	0.57	0.60	0.61	0.55
Llama-3-8B (n=all)	0.44	0.55	0.39	0.38	0.55	0.54	0.51	0.51
Llama-3-70B (n=1)	0.41	0.54	0.56	0.43	0.51	0.61	0.63	0.53
Llama-3-70B (n=3)	0.59	0.66	<b>0.67</b>	0.59	0.65	0.68	0.69	0.64
Llama-3-70B (n=all)	<b>0.71</b>	0.66	0.64	<b>0.64</b>	<b>0.76</b>	0.70	<b>0.70</b>	<b>0.70</b>

Table 2: Zero-shot and few-shot performance metrics for grounding act classification evaluated by macro-averaged accuracy, precision, recall, and F1-score. The variable n denotes the number of preceding input utterances. Bold values highlight the best value for each metric.

For explicit and implicit labels, the grounded knowledge items that have been shared until this point in the dialogue were annotated as a knowledge graph structure in JSON-LD format (Sporny et al., 2020). Annotation disagreements were collaboratively resolved to reach a consensus. Knowledge is incorporated into the grounding annotation only if it is a subset of the underlying tabular dataset and can be represented within the modeled internal system knowledge, which we defined using vocabulary from the namespaces *Schema.org* and *CSVW* (W3C, 2017, 2024). An example conversation illustrating labeled grounding acts and grounded knowledge items for individual dialogue utterances is provided in Table 4 in Appendix A.

**Experimental Setup** Based on the annotated dataset with conversational grounding labels, we conducted several experiments using four state-of-the-art LLMs: the open-source Llama-3-8B-Instruct as well as Llama-3-70B-Instruct (Meta AI, 2024) from the Llama 3 model family, and the closed-source models GPT-3.5-Turbo (version: 0125) and GPT-4o (version: 2024-05-13) (OpenAI, 2022, 2024). We defined two model prompts: one for classifying grounding acts and another for identifying grounded knowledge. For the knowledge identification prompt, which tasked the LLM to predict the grounded knowledge subset in the conversation thus far, we provided both the input dialogue and the complete system knowledge (i.e., the annotated grounded knowledge for the entire conversation). All models were prompted using a chat completion format, which included a system instruction and, in the few-shot setting, three in-context examples presented as user and assistant turns. Both model prompts are provided in

the Appendix in full length (Tables 5 and 6). To promote deterministic generation, we set the generation seed to 1 and the temperature parameter to 0. The maximum token limit was set to 128 for classification and 4096 for grounded knowledge identification. All generated outputs with extra text were preprocessed using a regular expression to match and extract the first occurrence of either the grounding act or JSON-LD array.

## 4 Results and Discussion

**Classification of Grounding Acts** Table 2 shows the performance for classifying grounding acts, using macro-averages to ensure equal class importance. Nearly all tested LLMs benefited from the added context of few-shot examples, with F1-scores generally improving; however, this improvement diminishes as the number of input dialogue turns (n) increases, suggesting potential redundancy when in-context examples are already provided. The results indicate that n=3 often optimizes

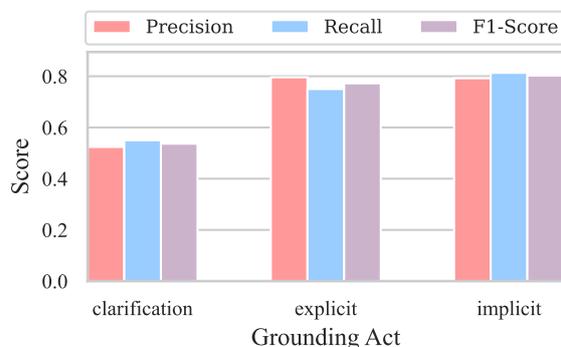


Figure 1: Performance comparison of precision, recall, and F1-score by grounding act for the Llama-3-70B model with all input utterances (n=all).

Issue Type	GPT-3.5-Turbo	GPT-4o	Llama-3-8B	Llama-3-70B
	Relative Frequency: Zero-Shot / Few-Shot			
Invalid JSON-LD	0.00 / 0.01	0.00 / 0.00	0.02 / 0.09	0.20 / 0.00
Property Hallucination	0.01 / 0.00	0.00 / 0.02	0.08 / 0.22	0.38 / 0.26
Value Hallucination	0.02 / 0.00	0.01 / 0.03	0.22 / 0.05	0.46 / 0.07
Property Excess	0.49 / 0.48	0.29 / 0.24	0.50 / 0.38	0.61 / 0.51
Property Deficit	0.37 / 0.22	0.31 / 0.09	0.50 / 0.36	0.39 / 0.20
Value Excess	0.68 / 0.63	0.40 / 0.31	0.66 / 0.32	0.76 / 0.47
Value Deficit	0.22 / 0.22	0.29 / 0.28	0.34 / 0.62	0.24 / 0.34

Table 3: Relative frequency of issues in zero- and few-shot predictions for grounded knowledge identification.

performance in both zero- and few-shot settings by balancing context retention, noise reduction, and efficient usage of tokens. While Llama-8B’s performance drops from 0.54 F1-score at n=1 to 0.38 at n=all, larger LLMs like Llama-70B and GPT-4o handle longer input better, probably due to a higher parameter count and superior noise handling.

Another significant finding is the competitive performance of open-source LLMs against proprietary ones: Llama-8B surpasses GPT-3.5 in the zero-shot run, and Llama-70B matches GPT-4o in the few-shot run. The breakdown of Llama-70B’s performance by grounding act, illustrated in Figure 1, reveals clarification as the most challenging act to classify, consistent with our observation of the other LLMs. For instance, the models often struggled when users tried to clarify a previously introduced concept. Instead of recognizing the clarification (e.g., “And category describes whether it is a movie, tv show, or work of literature?”), the models often misinterpreted it as introducing a new topic, falsely assuming that the previous concept is already implicitly grounded. Contrary to clarification acts, the F1-scores for explicit and implicit classification are comparable. Despite achieving the same overall F1-score, GPT-4o tends to overpredict implicit labels in contrast to the more balanced Llama-70B, as revealed by the confusion matrices in Figure 3 in Appendix A. The latter shows that GPT-4o excels at predicting explicit grounding accurately, avoiding false positives altogether, but it tends to overpredict the implicit class, particularly in cases where participants acknowledge information explicitly before asking a new question (e.g., “Ok very interesting! What is the highest level of protein in the chart?”).

**Identification of Grounded Knowledge** The second series of experiments aimed at identifying grounded knowledge for a suitable dialogue context, which is a significantly more complex task than classifying grounding acts (Wu et al., 2021;

Oh et al., 2023). Knowledge identification required the LLMs to uniquely pinpoint specific knowledge items from a set of possibilities within the system knowledge model, bridging between vague conversation utterances and structured JSON-LD arrays.

Figure 2 depicts the count of JSON-LD generations accurately matching our 127 annotations with valid properties, values, or completely identical content. The open-source models notably struggle more compared to the proprietary LLMs. While both open-source Llama models produce multiple valid outputs for properties and values with few-shot prompting, they fail to generate any valid predictions in the zero-shot setting. Therefore, these model runs are not displayed in the chart. Remarkably, GPT-4o outperforms GPT-3.5 by almost double, even in the zero-shot experiment, surpassing all other models by a great margin. In the few-shot cases, every third prediction from GPT-4o is identical to our annotated groundings, totaling 42 out of 127 instances. In some cases, The GPT-4o model even succeeded in precisely matching the annotated JSON-LD in a given conversation across a number of subsequent turns.

Table 3 provides a detailed analysis of the most common prediction issues and their relative fre-

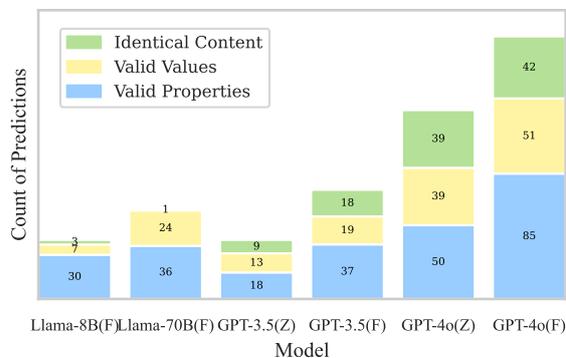


Figure 2: Count of predictions in JSON-LD format with valid properties, valid values, or identical content for evaluated models in zero- (Z) and few-shot (F) settings.

quencies for each model-prompt experiment. Examples for each issue type are listed in Table 7 in Appendix A. Open-source models generally produce more invalid JSON-LD arrays and hallucinate properties and values that are not part of the system knowledge. All tested LLMs tend to overpredict properties and values in zero-shot settings, even though these are grounded later in the conversation. Few-shot prompting can reduce excess properties and values, as well as counteract property deficits. However, in few-shot prompting, open-source models, particularly Llama-3-8B, tend to increase value deficits, becoming too hesitant to identify knowledge. This often results in empty JSON-LD arrays with generated statements such as “The conversation does not mention any specific knowledge items from the system knowledge.”

Our findings corroborate existing benchmarks, highlighting the sophisticated reasoning abilities of state-of-the-art proprietary LLMs such as GPT-4o in highly complex tasks. A similar task complexity-based LLM performance gap is also observable in the direct comparison of the MMLU and HumanEval benchmark scores between GPT-4o and Llama-3 (Hendrycks et al., 2020; Chen et al., 2021; OpenAI, 2024). While Llama-70B performs competitively in the language-focused grounding act classification task, the superiority of GPT-4o becomes apparent in identifying knowledge when handling structured JSON-LD data and fragmented information from dialogue utterances.

In short, when designing dialogue systems augmented with LLMs to handle conversational grounding, smaller open-source models like Llama-3-8B, especially fine-tuned versions, seem to be generally sufficient for basic NLP tasks such as detecting and classifying grounding-related dialogue acts. However, more complex tasks, such as identifying and integrating grounded knowledge from dialogue utterances with structured knowledge representations, require the use of more advanced and larger models like GPT-4o, which possess superior reasoning capabilities and proficiency in processing structured data formats.

## 5 Conclusion and Future Work

Our study examined LLMs for handling grounding-related knowledge in information-sharing dialogues. We found that classifying grounding acts was feasible for both open- and closed-source LLMs, with open-source LLMs performing on par

compared with leading proprietary ones. However, identifying grounded knowledge proved to be a distinctly more complex task. For the latter, the proprietary LLMs had a competitive edge, and the open-source models underperformed due to their higher predisposition to generate erroneous output. The experiment results from our newly created dataset highlight common prediction issues and demonstrate how few-shot prompting can enhance model outputs, offering valuable insights to advance research on conversational grounding.

Future work should concentrate on developing LLM-based dialogue systems that handle conversational grounding through a multi-component pipeline approach for recognizing grounding-specific dialogue acts as well as grounded knowledge (Jokinen et al., 2024). In previous studies, we have shown that LLMs can augment dialogue systems by performing semantic parsing for conversational question answering over knowledge graphs (Schneider et al., 2024a) and by verbalizing retrieved semantic triples into text responses (Schneider et al., 2024b). We believe conversational grounding is essential as it links the processes of semantic parsing of dialogue utterances, knowledge identification, and response generation, aligning the user’s prior knowledge with the system’s available knowledge base while maintaining the relevance and coherence of conversations.

## 6 Limitations

Our study has certain limitations that should be acknowledged. First, the experiments are based on a relatively small dataset, consisting of only 26 information-seeking conversations and 669 dialogue turns collected in a controlled laboratory setting. While these conversations span five distinct domains, the findings should be interpreted with caution, as they may not generalize to larger or more diverse dialogue corpora.

Additionally, the grounded knowledge annotations in our study are represented using the JSON-LD syntax. We chose the JSON-LD format because it is widely used, and many LLMs are trained to process JSON sequences effectively. However, it is important to recognize that other encoding formats, such as Turtle, RDF/XML, and N-Triples, may produce different performance results. Further, our experiments were restricted to the open-source Llama (Meta AI, 2024) and closed-source GPT (OpenAI, 2022, 2024) model families. It is

advisable for future work to explore an even bigger variety of LLMs, particularly those that are specifically trained on code and structured data like Codestral or Code Llama.

Lastly, conversational grounding in dialogue systems entails both the classification of grounding acts and the identification of grounded knowledge. While we have introduced and evaluated these tasks separately, incorporating our approach into an end-to-end evaluation could offer a more holistic understanding of end-to-end performance in more realistic dialogue scenarios.

## 7 Ethical Considerations

In our experiments, we used a publicly available dialogue dataset from [Schneider et al. \(2023\)](#) while ensuring that no personal identifying information of the participants was processed or disclosed. The information-seeking conversations from the dataset discuss only domain-specific knowledge from publicly accessible websites, such as Wikipedia. Moreover, to ensure optimal computing efficiency, evaluations of the Llama and GPT models were conducted on cloud computing platforms, with each inference run taking less than an hour.

## Acknowledgements

We would like to thank the anonymous reviewers for their helpful suggestions. Kristiina Jokinen acknowledges the support of Project JPNP20006 commissioned by the New Energy and Industrial Technology Development Organization (NEDO), Japan.

## References

Angelo Cangelosi. 2010. Grounding language in action and perception: From cognitive agents to humanoid robots. *Physics of life reviews*, 7(2):139–151.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. [Evaluating large language models trained on code](#). *arXiv preprint arXiv:2107.03374*.

Herbert H Clark and Edward F Schaefer. 1989. Contributing to discourse. *Cognitive science*, 13(2):259–294.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.

Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia D’amato, Gerard De Melo, Claudio Gutierrez, Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian Neumaier, Axel-Cyrille Ngonga Ngomo, Axel Polleres, Sabbir M. Rashid, Anisa Rula, Lukas Schmelzeisen, Juan Sequeda, Steffen Staab, and Antoine Zimmermann. 2021. [Knowledge graphs](#). *ACM Comput. Surv.*, 54(4).

Kristiina Jokinen, Phillip Schneider, and Taiga Mori. 2024. [Towards harnessing large language models for comprehension of conversational grounding](#). In *In 14th International Workshop on Spoken Dialogue Systems Technology (IWSDS 2024)*, Sapporo, Japan.

Oliver Lemon. 2022. [Conversational grounding in emergent communication—data and divergence](#). In *Emergent Communication Workshop at ICLR 2022*.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.

Meta AI. 2024. [Introducing Meta Llama 3: The most capable openly available LLM to date](#). *Meta AI Blog*.

Biswesh Mohapatra, Seemab Hassan, Laurent Romary, and Justine Cassell. 2024. [Conversational grounding: Annotation and analysis of grounding acts and grounding units](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3967–3977, Torino, Italia. ELRA and ICCL.

Minsik Oh, Joosung Lee, Jiwei Li, and Guoyin Wang. 2023. [PK-ICR: Persona-knowledge interactive multi-context retrieval for grounded dialogue](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16383–16395, Singapore. Association for Computational Linguistics.

OpenAI. 2022. [ChatGPT: Optimizing language models for dialogue](#). *OpenAI Blog*.

OpenAI. 2024. [Hello GPT-4o](#). *OpenAI Blog*.

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.

Phillip Schneider, Anum Afzal, Juraj Vladika, Daniel Braun, and Florian Matthes. 2023. [Investigating conversational search behavior for domain exploration](#). In *European Conference on Information Retrieval*, pages 608–616. Springer.

- Phillip Schneider, Manuel Klettner, Kristiina Jokinen, Elena Simperl, and Florian Matthes. 2024a. [Evaluating large language models in semantic parsing for conversational question answering over knowledge graphs](#). In *International Conference on Agents and Artificial Intelligence*.
- Phillip Schneider, Manuel Klettner, Elena Simperl, and Florian Matthes. 2024b. [A comparative analysis of conversational large language models in knowledge-based text generation](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 358–367, St. Julian’s, Malta. Association for Computational Linguistics.
- Phillip Schneider, Tim Schopf, Juraj Vladika, Mikhail Galkin, Elena Simperl, and Florian Matthes. 2022. [A decade of knowledge graphs in natural language processing: A survey](#). 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 601–614, Online only. Association for Computational Linguistics.
- Omar Shaikh, Kristina Gligorić, Ashna Khetan, Matthias Gerstgrasser, Diyi Yang, and Dan Jurafsky. 2024. [Grounding gaps in language model generations](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.
- Manu Sporny, Dave Longley, Gregg Kellogg, Markus Lanthaler, Pierre-Antoine Champin, and Niklas Lindström. 2020. [JSON-LD 1.1. W3C Recommendation](#).
- David Traum. 1994. [A computational theory of grounding in natural language conversation](#). *PhD thesis, Univ. Rochester*.
- Takuma Udagawa and Akiko Aizawa. 2021. [Maintaining common ground in dynamic environments](#). *Transactions of the Association for Computational Linguistics*, 9:995–1011.
- W3C. 2017. [CSVW Namespace Vocabulary Terms. W3C Document](#).
- W3C. 2024. [Schema.org. Schema.org](#).
- Zejiu Wu, Bo-Ru Lu, Hannaneh Hajishirzi, and Mari Ostendorf. 2021. [DIALKI: Knowledge identification in conversational systems through dialogue-document contextualization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1852–1863, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

## A Appendix

The Appendix provides one annotated conversation example (Table 4), the model prompts in full length (Tables 5 and 6), an overview of common issue types identified in the predictions (Table 7), and two confusion matrices of the classification results of the two best-performing model inference runs (Figure 3).

Dialogue Utterances	Dialogue Grounded Knowledge Act	
S: What is your dataset about?	-	-
P: it contains information about 11341 historical figures, including their full name, sex, birth year, city, country, continent, occupation, historical popularity index (HPI). The HPI represents the degree of this person’s online popularity	-	-
S: Who is the most popular?	implicit	<pre>[{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/historical-figures", "url": "historical-figures.csv", "schema:description": "The table contains information about historical figures", "tableSchema": {"columns": [{"name": "full_name", "datatype": "string"}, {"name": "sex", "datatype": "string"}, {"name": "birth_year", "datatype": "integer"}, {"name": "city", "datatype": "string"}, {"name": "country", "datatype": "string"}, {"name": "continent", "datatype": "string"}, {"name": "occupation", "datatype": "string"}, {"name": "historical_popularity_index", "datatype": "float"}], "primaryKey": "full_name"}]}</pre>
P: Aristotle, who is from Greece and has a largest HPI value: 31.9938.	-	-
S: I see, is there Socrate in the dataset?	explicit	<pre>[{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/historical-figures", "url": "historical-figures.csv", "schema:description": "The table contains information about historical figures", "tableSchema": {"columns": [{"name": "full_name", "datatype": "string"}, {"name": "sex", "datatype": "string"}, {"name": "birth_year", "datatype": "integer"}, {"name": "city", "datatype": "string"}, {"name": "country", "datatype": "string"}, {"name": "continent", "datatype": "string"}, {"name": "occupation", "datatype": "string"}, {"name": "historical_popularity_index", "datatype": "float"}, {"maximum": 31.9938}], "primaryKey": "full_name"}], {"@type": "schema:Person", "full_name": "Aristotle", "country": "Greece", "historical_popularity_index": 31.9938}]</pre>
P: Yes, Socrate is in the dataset.	-	-
S: What is its popularity index?	implicit	<pre>[{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/historical-figures", "url": "historical-figures.csv", "schema:description": "The table contains information about historical figures", "tableSchema": {"columns": [{"name": "full_name", "datatype": "string"}, {"name": "sex", "datatype": "string"}, {"name": "birth_year", "datatype": "integer"}, {"name": "city", "datatype": "string"}, {"name": "country", "datatype": "string"}, {"name": "continent", "datatype": "string"}, {"name": "occupation", "datatype": "string"}, {"name": "historical_popularity_index", "datatype": "float"}, {"maximum": 31.9938}], "primaryKey": "full_name"}], {"@type": "schema:Person", "full_name": "Aristotle", "country": "Greece", "historical_popularity_index": 31.9938}, {"@type": "schema:Person", "full_name": "Socrates"}]}</pre>
P: Historical popularity index (HPI) is metric that aggregates information on a biography’s online popularity. It aggregates information on the age and attention received by biographies in multiple language editions of Wikipedia to provide a summary statistic of their global popularity.	-	-

Table 4: Example of dialogue excerpt from the history domain with annotated grounding dialogue acts and grounded knowledge in JSON-LD format. Seeker (S) and provider (P) roles are abbreviated for each turn. Utterances are taken from the dialogue logs and may contain spelling errors. Newly grounded knowledge is displayed in blue color.

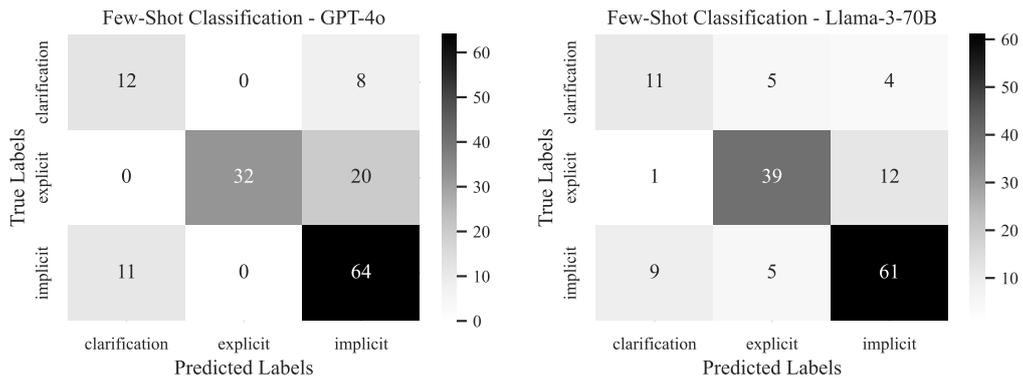


Figure 3: Confusion matrices for few-shot classification results of GPT-4o with three input utterances and Llama-3-70B with all input utterances.

<b>Grounding Act Classification Prompt</b>	
<b>Zero-Shot</b>	
SYSTEM: Predict the grounding label for the last response in the 'Input Dialogue:'. The label indicates whether the knowledge in the dialogue was accepted. Choose one of the following labels:	
explicit: The response confirms understanding or acceptance (e.g., 'okay', 'thanks', 'alright', 'nice') without seeking clarification.	
clarification: The response seeks clarification about a previous dialogue snippet.	
implicit: The response moves the conversation forward without explicitly confirming or seeking clarification.	
<b>Few-Shot</b>	
SYSTEM: Predict the grounding label for the last response in the 'Input Dialogue:'. The label indicates whether the knowledge in the dialogue was accepted. Choose one of the following labels:	
explicit: The response confirms understanding or acceptance (e.g., 'okay', 'thanks', 'alright', 'nice') without seeking clarification.	
clarification: The response seeks clarification about a previous dialogue snippet.	
implicit: The response moves the conversation forward without explicitly confirming or seeking clarification.	
USER: Input Dialogue:	seeker: Can you give me some information about your dataset?
provider:	My dataset includes information on buildings of Gothic architecture.
seeker:	How tall is the Cologne Cathedral?
ASSISTANT:	Output Label: implicit
USER: Input Dialogue:	provider: Monitors have different attributes like size or panel technology.
provider:	There are some with an aspect ratio of 21:9.
seeker:	What is aspect ratio?
ASSISTANT:	Output Label: clarification
USER: Input Dialogue:	provider: An elephant's average lifespan is around 65 years.
seeker:	I see, good to know.
ASSISTANT:	Output Label: explicit

Table 5: Overview of applied zero-shot and few-shot prompts for classification.

---

**Grounded Knowledge Identification Prompt**

---

**Zero-Shot**

---

SYSTEM: Your task is to identify the knowledge items that have been grounded by the conversation partners in the 'Input Dialogue'. The items of mutually grounded knowledge must be explicitly mentioned in the dialogue. Based on the complete set of 'System Knowledge', your task is to generate the subset of knowledge items that have been grounded so far. Ensure that the output is a valid JSON-LD structure (an array of JSON objects) and only include knowledge items from the formatted 'System Knowledge'.

---

**Few-Shot**

---

SYSTEM: Your task is to identify the knowledge items that have been grounded by the conversation partners in the 'Input Dialogue'. The items of mutually grounded knowledge must be explicitly mentioned in the dialogue. Based on the complete set of 'System Knowledge', your task is to generate the subset of knowledge items that have been grounded so far. Ensure that the output is a valid JSON-LD structure (an array of JSON objects) and only include knowledge items from the formatted 'System Knowledge'.

USER: System Knowledge: [{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/american-presidents", "url": "american-presidents.csv", "schema:description": "The table contains information about American presidents", "tableSchema": {"columns": [{"name": "name", "datatype": "string"}, {"name": "term", "datatype": "string"}, {"name": "party", "datatype": "string"}, {"name": "election\_year", "datatype": "integer"}]}, "primaryKey": "name", {"@type": "schema:Person", "name": "Barack Obama", "party": "Democratic"}]  
Input Dialogue:

seeker: Can you give me an example entry from your dataset?

provider: One of the presidents in the list is Barack Obama.

seeker: Thanks. What party does he belong to?

ASSISTANT: Output JSON-LD: [{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/american-presidents", "url": "american-presidents.csv", "schema:description": "The table contains information about American presidents", "tableSchema": {"columns": [{"name": "name", "datatype": "string"}]}, "primaryKey": "name", {"@type": "schema:Person", "name": "Barack Obama"}]

USER: System Knowledge: [{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/greek-islands", "url": "greek-islands.csv", "schema:description": "The table contains information about islands in Greece", "tableSchema": {"columns": [{"name": "island", "datatype": "string"}, {"name": "area\_in\_km2", "datatype": "integer", "minimum": 64, "maximum": 8336}, {"name": "cluster", "datatype": "string"}]}, "primaryKey": "island", {"@type": "schema:Place", "island": "Crete", "area\_in\_km2": 8336, "cluster": "Cretan"}, {"@type": "schema:Place", "island": "Alonnisos", "area\_in\_km2": 64, "cluster": "Sporades"}, {"@type": "schema:Place", "island": "Lesbos", "area\_in\_km2": 1633, "cluster": "North Aegean Islands"}]  
Input Dialogue:

provider: My dataset contains information on Greek islands. For example, there is Crete with an area of 8336 square kilometers.

provider: That makes it the largest island in Greece.

seeker: Which one is the smallest and what is its area?

ASSISTANT: Output JSON-LD: [{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/greek-islands", "url": "greek-islands.csv", "schema:description": "The table contains information about islands in Greece", "tableSchema": {"columns": [{"name": "island", "datatype": "string"}, {"name": "area\_in\_km2", "datatype": "integer", "maximum": 8336}], "primaryKey": "island", {"@type": "schema:Place", "island": "Crete", "area\_in\_km2": 8336}}]

USER: System Knowledge: [{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/android-smartphones", "url": "android-smartphones.csv", "schema:description": "The table contains information about Android smartphones", "tableSchema": {"columns": [{"name": "model", "datatype": "string"}, {"name": "developer", "datatype": "string"}, {"name": "release\_year", "datatype": "integer", "minimum": 2008, "maximum": 2024}, {"name": "android\_version", "datatype": "string"}]}, "primaryKey": "model", {"@type": "schema:Product", "model": "HTC Dream", "developer": "HTC", "release\_year": "2008"}, {"@type": "schema:Product", "model": "LG Wing", "developer": "LG", "release\_year": "2020", "android\_version": "Android 10"}, {"@type": "schema:Product", "release\_year": "2024"}]  
Input Dialogue:

provider: I can provide technical information about Android smartphones.

provider: One column contains data about the model and another specifies its release year.

seeker: I see, good to know.

ASSISTANT: Output JSON-LD: [{"@context": ["http://www.w3.org/ns/csvw", {"schema": "http://schema.org"}], "@id": "http://example.org/android-smartphones", "url": "android-smartphones.csv", "schema:description": "The table contains information about Android smartphones", "tableSchema": {"columns": [{"name": "model", "datatype": "string"}, {"name": "release\_year", "datatype": "integer"}]}, "primaryKey": "model"}]

---

Table 6: Overview of applied zero-shot and few-shot prompts for grounding knowledge.

Issue Type	Definition	Example
Invalid JSON-LD	Prediction is overall not a valid JSON-LD structure.	PRED: [{"@type": "schema:CreativeWork", "category": ["films", "series", "literature"]}] → <b>JSON syntax elements are missing.</b>
Property Hallucination	Predicted property is not part of system knowledge.	PRED: {"@type": "schema:Place", "name": "English Garden", "state": "Bavaria", "city": "Munich", "summary": "The English Garden is a large public park in Munich, Bavaria. It is one of the world's largest urban public parks, larger than New York's Central Park. The park is known for its beer gardens, streams, lakes, and Japanese teahouse."} → <b>English Garden was mentioned in the dialogue but is not part of the system knowledge. The "city" property does not exist.</b>
Value Hallucination	Predicted value is not part of system knowledge.	PRED: [{"@context": ["http://www.w3.org/ns/csvw", "http://schema.org"], "@id": "http://example.org/time-travel-media", "url": "time-travel-media.csv", "schema:description": "The table contains information about time travel media", "tableSchema": {"columns": [{"name": "year", "datatype": "string"}, {"name": "title", "datatype": "string"}, {"name": "author", "datatype": "string"}, {"name": "short text description", "datatype": "string"}, {"name": "category", "datatype": "string"}], "primaryKey": ["year", "title", "category"]}] → <b>System knowledge contains the columns "creators" and "description" but "author" and "short text description" are hallucinated.</b>
Property Excess	Predicted property is part of system knowledge but grounded too early.	PRED: [{"@context": ["http://www.w3.org/ns/csvw", "http://schema.org"], "@id": "http://example.org/football-records", "url": "football-records.csv", "schema:description": "The table contains information about football records", "tableSchema": {"columns": [{"name": "record", "datatype": "string"}], "primaryKey": "record"}] → <b>Properties should have been grounded at a later stage.</b>
Property Deficit	Property from system knowledge was not identified.	PRED: [{"@context": ["http://www.w3.org/ns/csvw", "http://schema.org"], "@id": "http://example.org/nature-parks", "url": "nature-parks.csv", "schema:description": "The table contains information about nature parks in Germany", "tableSchema": {"columns": [{"name": "name", "datatype": "string"}, {"name": "state", "datatype": "string"}, {"name": "year", "datatype": "integer"}, {"name": "area_in_km2", "datatype": "integer"}, {"name": "summary", "datatype": "string"}], "primaryKey": "name"}], {"@type": "schema:Place", "name": "Altmühl Valley Nature Park"}] → <b>Highlighted properties were not grounded.</b>
Value Excess	Predicted value is part of system knowledge but grounded too early.	PRED: [{"@context": ["http://www.w3.org/ns/csvw", "http://schema.org"], "@id": "http://example.org/nature-parks", "url": "nature-parks.csv", "schema:description": "The table contains information about nature parks in Germany", "tableSchema": {"columns": [{"name": "name", "datatype": "string"}, {"name": "state", "datatype": "string"}, {"name": "year", "datatype": "integer"}, {"name": "area_in_km2", "datatype": "integer"}, {"name": "summary", "datatype": "string"}], "primaryKey": "name"}] → <b>Values should have been grounded at a later stage.</b>
Value Deficit	Value from system knowledge was not identified.	PRED: [{"@context": ["http://www.w3.org/ns/csvw", "http://schema.org"], "@id": "http://example.org/historical-figures", "url": "historical-figures.csv", "schema:description": "The table contains information about historical figures", "tableSchema": {"columns": [{"name": "full_name", "datatype": "string"}, {"name": "birth_year", "datatype": "integer", "minimum": -3500, "maximum": 2005}], "primaryKey": "full_name"}], {"@type": "schema:Person", "birth_year": -3500}, {"@type": "schema:Person", "birth_year": 2005}] → <b>Highlighted values were not grounded.</b>

Table 7: Overview of six identified issue types with examples from generated model predictions (PRED). The manifestation of issues are highlighted in red color.

# "Keep up the good work!": Using Constraints in Zero-Shot Prompting to Generate Supportive Teacher Responses

E. Margaret Perkoff<sup>1</sup>, Angela Ramirez<sup>2</sup>, Sean Von Bayern<sup>1</sup>,  
Marilyn Walker<sup>2</sup>, James Martin<sup>1</sup>

<sup>1</sup>University of Colorado Boulder, <sup>2</sup>University of California Santa Cruz

Correspondence: [elpe8468@colorado.edu](mailto:elpe8468@colorado.edu)

## Abstract

Educational dialogue systems have been used to support students and teachers for decades. Such systems rely on explicit pedagogically-motivated dialogue rules. With the ease of integrating large language models (LLMs) into dialogue systems, applications have been arising that directly use model responses without the use of human-written rules, raising concerns about their use in classroom settings. Here, we explore how to constrain LLM outputs to generate appropriate and supportive teacher-like responses. We present results comparing the effectiveness of different constraint variations in a zero-shot prompting setting on a large mathematics classroom corpus. Generated outputs are evaluated with human annotation for Fluency, Relevance, Helpfulness, and Adherence to the provided constraints. Including all constraints in the prompt led to the highest values for Fluency and Helpfulness, and the second highest value for Relevance. The annotation results also demonstrate that the prompts that result in the highest adherence to constraints do not necessarily indicate higher perceived scores for Fluency, Relevance, or Helpfulness. In a direct comparison, all of the non-baseline LLM responses were ranked higher than the actual teacher responses in the corpus over 50% of the time.

## 1 Introduction

Large language models (LLMs) have shown great promise across many applications including recommendation systems, social chatbots, writing code, and summarizing documents (Zhang et al., 2023). Many of these applications benefit from the generative capabilities of LLMs such as ChatGPT (Brown et al., 2020). However, when these models are deployed without further constraints in open-domain dialogue systems, they may generate outputs that do not adhere to the desired agent behavior (Kann et al., 2022). The risk of not adhering to desired

*Teacher* How do you know number two is not a straight angle?

*Student* Because a straight angle goes on, on both sides.

*Teacher* It goes on forever. There's a line, and it goes on forever on both sides, absolutely. Of course, it's a ray so it only goes in one direction.

Table 1: A sample interaction between a teacher and a student from the NCTE Corpus (Demszky and Hill, 2023)

agent behavior is even higher when we consider the application of these models to the education domain (Williams et al., 2023).

Well before the advent of LLMs, educational dialogue systems have been used to provide support to students in online classes, act as 1:1 subject-specific tutors, and provide professional training to teachers and tutors (Kuhail et al., 2023). Regardless of the exact learning application, they are more beneficial to students when the systems themselves are based on the same pedagogical frameworks that a teacher would use (Järvelä and Hadwin, 2013). As such, many of these systems are built using extensive dialogue frameworks that dictate when, and how the teacher should intervene in a particular scenario. Frequently, they are designed such that a teacher can modify the exact responses to apply to a given classroom age, subject, and lesson-specific goals. Although this makes it possible to create highly relevant responses from the conversational agent, it also means that the agent will not generalize well to new situations. Furthermore, these agents are typically designed for interaction with a single student working in an online classroom.

The educational dialogue system that we present here is designed to provide support to students in a real classroom setting. Consider the interaction in

Table 1. In this case, the teacher confirms that the student’s understanding is correct and elaborates on what has been said to avoid confusion in the definition. This ensures that they are still encouraging the student to elaborate on their reasoning while providing a fully correct assessment of the problem. This highlights the degree of nuance required when responding to a student in a classroom setting. We want to focus on how to constrain the output of an LLM in a manner consistent with how a teacher would respond. Constraints are surfaced to the LLM through strategic prompt engineering. The constraints are based on dialogue acts that capture speaker intentions from the preceding dialogue and a domain-specific dialogue policy. The dialogue policy consists of three conversational states identified in collaborative task conversations and corresponding actions that a supportive teacher would take in that particular scenario. We evaluate the effectiveness of imposing different forms of constraints on the NCTE Corpus (Demszky and Hill, 2023) of math classroom transcripts. The output from each of the prompt variations is assessed for Fluency, Relevance, Helpfulness, and adherence to the desired constraints. We also perform further annotations for overall best response and whether it was considered better than the teacher.

The dialogue system we present here provides a blueprint for how future researchers, and teachers themselves, can modify prompts to better engage with students in a classroom setting. This design is intended to echo the highly relevant and helpful nature of early rule-based education dialogue systems while allowing for more flexibility with LLMs to surface the response to the student. The major contributions of this work are as follows:

- As far as we are aware, this is the first application of dialogue-policy informed LLM response generation to the education domain.
- Adding more context-specific constraints (Dialogue States, Student and Teacher Dialogue Acts) can improve the perceived Helpfulness, Fluency and Relevance of responses produced by a conversational agent.
- LLMs can adhere to dialogue state constraints in up to 95% of samples
- Annotators rated the non-baseline LLM-based prompt variations to be better than the actual Teacher response in > 50% of samples

The results demonstrate the potential of LLM-based chatbots to interact with students in a helpful manner. There is much research to be done in exploring how to balance additional pedagogical constraints while maintaining a high degree of fluency in the responses. In future work, we intend to integrate these models with a speech-to-text interface to see how they perform in live classroom interactions.

## 2 Related Work

### 2.1 Dialogue Act Segmentation and Classification

To respond in a succinct manner, dialogue systems need to be able to differentiate different dialogue acts (DAs) such as statements, types of questions, and acknowledgements. There are different dialogue act schemas and datasets for dialogue act tagging such as: ViGGO a video game corpus tagged with dialogue acts meant for open domain systems (Juraska et al., 2019), MultiWoz a multi-domain and topic dataset meant for task-oriented dialogue systems (Budzianowski et al., 2018), Switchboard (SWDA) is a large multi-speaker dataset consisting of two-sided telephone calls (Stolcke et al., 2000a) and the AMI meeting corpus that is multi-modal corpus consisting of 100 hours of meeting recordings (Shang et al., 2018). ISO dialogue acts schema are mapped to other corpora such as SWDA, AMI, Maptask to then be used for training a SVM model on ISO tags (Mezza et al., 2018a; Thompson et al., 1993).

Previous work, has used this ISO dialogue act model off-the-shelf model to enrich TopicalChats with dialogue acts to then use the dialogue acts for response generation (Hedayatnia et al., 2020). LLMs such as gpt-3.5-turbo have been used for data augmentation and annotation of dialogue acts on outputs within the education domain (Shan et al., 2023). We expand on these papers by tagging dialogue acts using a combination of ISO tags and Switchboard feedback dialogue acts and by using using gpt-3.5-turbo on the NCTE dataset to produce a silver set of annotated dialogue acts as this corpus does not contain tags. Then we leverage this new tagged dataset as a constraint to provide the model for controllable response generation (Ramirez et al., 2023; Hedayatnia et al., 2020).

Dialogue act tagging and segmentation can be split into two separate tasks such as dialogue act classifiers (Stolcke et al., 2000a; Webb and Wilks,

2005), or a joint approach (Guz et al., 2010; Zhao and Kawahara, 2019, 2017). We take inspiration from joint models by combining the two tasks within the same prompt. The description used to prompt the model resembles annotation manuals that request annotators to consider both segmenting and tagging each segment with a dialogue act (Mezza et al., 2018a).

## 2.2 Pedagogical Conversational Agents

Pedagogical conversational agents are the subset of language models that can engage in dialogues to support learning. They vary greatly in terms of their role, their interaction style, and their functional purpose. Recent reviews of pedagogical conversational agents have found that they are frequently used as Teaching Agents in the context of online classroom settings, with the majority focused on Computer Science classroom courses (Kuhail et al., 2023). Although early research focused more on rule-based dialogue systems to power these agents, the surge in popularity of generative language models has led to more research examining the ability of language models to generate teacher-like responses. Tack and Piech (2022) proposed the AI Teacher Test to measure the effectiveness of a language model to engage in dialogues with a student based on the ability of the model to “speak like a teacher, understand the student, and help the student”. They evaluate several BlenderBot (Roller et al., 2020) and GPT-3 (Brown et al., 2020) models on the Teacher-Student Chatroom Corpus (TSCC) (Caines et al., 2022) as well as the Educational Uptake Dataset (Demszky et al., 2021) which is a subset of the NCTE corpus that is used in this paper. Their findings found that Blender outperformed GPT-3 across all metrics, and outperformed the actual Teacher response ratings in terms of levels of uptake from student responses. However, this did not translate to outperforming teachers in terms of levels of helpfulness or the ability to create responses similar to a teacher. The BEA 2023 shared task motivated researchers to expand on this work by focusing on generating responses to student utterances from the TSCC dataset. All of the submissions were evaluated using a set of automatic dialogue evaluation metrics from (Yeh et al., 2021) and the top three models were evaluated with pairwise comparisons from human raters based on the Tack’s original three categories. The best performing model, NAISTEACHER (Vasselli et al., 2023) was built on a pre-trained GPT 3.5 Turbo

(Brown et al., 2020). Their approach generates multiple teacher utterances in the form of either continuations of a previous utterance or replies to a student utterance. The responses are then re-ranked with DialogRPT (Gao et al., 2020). As part of this shared task, Hicke et al. Hicke et al. (2023) explored the use of GPT-4, as well as fine-tuning DialoGPT (Zhang et al., 2020), FlanT5 (Chung et al., 2022) and GPT-2 (Radford et al., 2019). GPT-4 outperformed their other variants. Other submissions focused on prompt engineering with open-source language models (Baladón et al., 2023) including Opt-2.7B (Gao et al., 2020) and Alpaca (Taori et al., 2023). The proposed system combines zero-shot prompts with a state-of-the-art LLM with previous research in controllable text generation and dialogue act classification to create teacher responses.

## 3 Methods

### 3.1 Overall System Design

The proposed pedagogical dialogue system consists of three major components: (1) a dialogue act segmentation and classification module (2) a dialogue policy that dictates when and how the conversational agent should intervene and (3) a response generation module that aggregates the output of components (1) and (2) into a prompt for an LLM. The dialogue system receives the most recent student utterance and the conversation history annotated with speaker labels by turn as seen in **Figure 1**. The conversation turns are fed to (1) where they are segmented into individual utterances and labeled with a dialogue act. In parallel, the most recent student utterance is fed to the dialogue policy to determine the dialogue state. The conversation state and dialogue act annotated conversation history are then aggregated into a coherent prompt for an LLM.

### 3.2 Dialogue Policy

A dialogue policy specifies, for each dialogue state, the actions the dialogue system can take in that state (Walker et al., 1998; Levin et al., 2000; Rieser and Lemon, 2011). Here, we leverage a dialogue policy based on an analysis of collaborative dialogue scenarios between college-aged students, where we frame the policy in terms of the dialogue acts available in each state. The students were working on a sensor-based task, in which each student becomes an expert in a particular kind of sensor - moisture, environmental, and sound, and then has to share

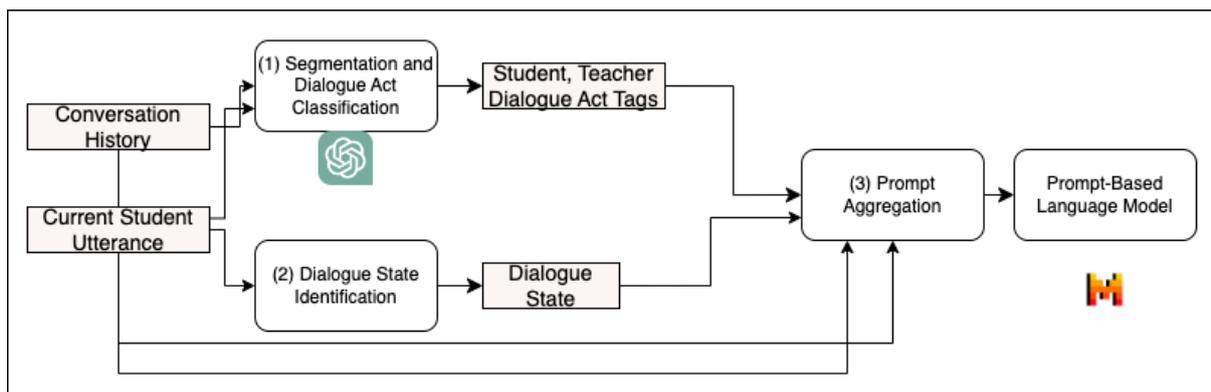


Figure 1: **Overall System Diagram:** The conversation history and current student utterance are given as input to both (1) segmentation and dialogue act classification and (2) a dialogue state identification module. They are then aggregated with the outputs of (1) and (2) into a prompt for a large language model.

the information they have learned with their group.

Education and human teamwork experts identified three broad categories of unproductive talk dialogue states during the collaboration sessions: *CONTENT*, *COLLABORATION*, and *PROCEDURAL*. The aim of our dialogue policy is to intervene in the student conversations when the conversation is in an unproductive state. Thus, each dialogue state has a set of potential dialogue actions, which are surfaced to the LLM as an additional component of the response generation prompt. The prompts are detailed in Appendix C and the complete steps for state identification are detailed in Appendix A. *PROCEDURAL* issues are identified when there is a literal String match indicating the students are confused about the next step to take to complete the task. Next, *CONTENT* issues are distinguished based on the topic feature of the NCTE dataset. The policy further divides the collaboration issues based on the number of speakers represented in the recent conversation history including *NO SPEAKER* and *SINGLE SPEAKER* categories. The NCTE dataset is segmented into only two speaker labels “student” and “teacher” without explicitly distinguishing which student is talking so we are not using these states in the experiments. For future work with multi-party conversations where speaker identification is available these are critical states to track. Additionally, since we are explicitly using non-blank utterances for the analysis, we do not have instances of the *NO SPEAKER* state. If none of the above conditions are met, the conversation is considered to be in a *FLOW* state and the student should be encouraged to continue with their reasoning.

### 3.3 Dialogue Act Segmentation and Classification

We use dialogue acts as a constraint for response generation as dialogue acts have been used for controllable response generation in different domains and dialogue systems (Hedayatnia et al., 2020; Ramirez et al., 2023). Following Shan et al., dialogue act labels for the NCTE corpus were obtained through few shot prompting using the gpt-3.5 turbo model. We combine the task of dialogue act classification and segmentation within the same prompt to handle longer utterances. For dialogue acts, we use the 10 ISO dialogue acts (Mezza et al., 2018b) (e.g., Inform, Set Question, Apology, etc.) and 7 Switchboard (Stolcke et al., 2000b) feedback-based dialogue acts (e.g. Backchannel, Sympathy, Appreciation) to classify and segment the NCTE dataset.

In the prompt, we provide a description of the task and instructions for the model, then for each dialogue act we give the definition and an optionally an example for certain difficult dialogue acts, and lastly we provide 19 examples of utterances that are segmented and annotated with dialogue act tags the final prompt can be found in the Appendix B. These examples came either directly from the Switchboard and ISO annotation manuals, or were handcrafted to contain the characteristics of the dialogue acts. To verify the performance of the prompt(s), we ran the prompt on a set of 100 examples, and would annotate for dialogue act and segmentation appropriateness on the teacher text and would adjust the number of examples or the definitions. Once we had the final prompt, we found that 95% of the time the model would choose an appropriate dialogue act(s). We note that the teacher dialogue acts were not always appropriate for the

conversation that could impact the performance in the downstream task.

### 3.4 Response Generation with Prompt Engineering

We take a modular approach to prompt engineering, wherein we dynamically construct a unique prompt for each response from component templates. These templates are injected with contextual information and desired constraints before being aggregated into a single prompt as input to the LLM. The complete set of templates can be found in **Appendix C**.

First, the baseline prompt has four components: *Preamble*, *Setting*, *Formatting*, and *Task*. We consider these to be the minimum configuration necessary to reliably produce teacher-like responses without applying any further constraints on the model’s behavior. Definitions for these basic components are as follows:

- **Preamble:** This primes the model by describing the task it will ultimately be assigned after first being given the requirements and restrictions contained in the other components.
- **Setting:** This describes a persona that the model should assume, the responsibilities of its role, and environmental details, like the grade level of the students and the subject of the current lesson.
- **Formatting:** This outlines the exact format the response should take, including a soft- and hard-cap on length, expectations of attribution and quotation marks, prohibitions of explained reasoning or word count, etc.
- **Task:** This simply instructs the model to produce a response and provides the utterance to which it will respond.

Then, we add one or more of the following experimental components: *Context*, *Student DA*, *Teacher DA*, and/or *Dialogue Policy*. Unlike the baseline components, these aim to constrain model behavior. When active, they are inserted into the prompt between *Formatting* and *Task*. Definitions for these additional components are as follows:

- **Context:** This gives the model a brief conversation history, consisting of three prior utterances and who said them.

- **Student DA:** This gives the model a version of the student utterance that is segmented by the classifier (see **Section 3.3**) and the definition for each of the resulting dialogue act labels. It then instructs the model to use the given segmentation and definitions to determine what the student meant.
- **Teacher DA:** This gives the model a list of dialogue act labels (and associated definitions) found within the segmented teacher utterance, **but it does not give the actual segmented utterance**. It then instructs the model to formulate its response to match the definitions of the given labels.
- **Dialogue Policy:** This gives the model a description of the current dialogue state (see **Appendix A**) and the consequences of allowing that state to continue. It then instructs the model to formulate a response that includes the expert-recommended intervention appropriate for the given state.

Finally, we have one *implicit* constraint: the language used in the templates mimics the language used by teachers when assigning schoolwork to students. The intent is to bias the model toward seeking similar language from its training data.

## 4 Experiments

### 4.1 Dataset

For this work, we want to demonstrate the appropriateness of different prompt variations for teacher-like responses. We use the NCTE Transcripts corpus (Demszky and Hill, 2023)- a dataset consisting of 1,660 anonymized elementary school mathematics lessons. This is one of the few publicly available datasets with annotated classroom conversations. We use a subset of the dataset that is broken down into student-teacher utterance pairs where each teacher response is associated with the immediately preceding student utterance. We also provide the three prior dialogue turns for conversation history. When utterances were within the first three dialogue turns of a particular lesson, only the available turns were provided as part of the context. For evaluation, we classified the utterances and ended up with three distinct state labels present in the dataset: *CONTENT*, *FLOW*, and *COLLABORATION*. We randomly selected 100 samples consisting of 33 *FLOW* turns, 33 *COLLABORATION*,

<b>Quality + Adherence Metrics</b>	0.29606
<b>Comparison to Teacher</b>	0.24689
<b>Best Model Response</b>	0.21557

Table 2: Interrater Reliability based on Krippendorff’s Alpha. The Quality and Adherence metrics are judged on 100 samples for all prompt variations individually. The Response Comparison metrics are based on a subset of 30 samples.

and 34 CONTENT, which is indicative of the highest subset of the datasets. All of the analysis sample utterances are then automatically segmented and classified for dialogue acts.

## 4.2 Prompt Variations

As detailed in **Section 3.4**, we aggregate different constraints into the prompts for the Mistral model. The goal is to discern which of these constraints leads to the most helpful, pedagogically informed outputs from the system. We start with the *Baseline* classroom response generation prompt, then provide additional constraints including: *Context*, *Dialogue State (DS)*, *DS + Student Dialogue Act (DA)*, *Student DA + Teacher DA*, and a prompt combining all the constraints (*DS, Student DA, Teacher DA*). All of the non-baseline prompts including the conversation history.

## 4.3 Evaluation

The generated responses are evaluated by raters on three categories: *Fluency*, *Relevance*, and *Helpfulness* using a scale from 0 to 2. Fluency describes how natural the generated response appears to be. This is meant to be comparable to prior work that evaluates text for *Naturalness* and/or *Humanness*. *Relevance* indicates how relevant the response is to the conversation history and current student utterance, with a 1 indicating that the response is vague. Raters also annotate for *Helpfulness* to indicate whether the response helps the student or helps move the activity forward. Finally, we have a binary value to indicate how well the generated output adheres to the dialogue state-specific prompt and whether it is an acceptable based on the ISO definition of the desired output DA. We calculate Krippendorff’s (Krippendorff, 2004) Alpha to gauge interrater reliability across the 600 samples evaluated by each of the three raters. The results are in **Table 2**.

## 5 Results Analysis

The annotation analysis is divided into three sections: output quality metrics, adherence to constraints, and teacher comparison. For the first two sections, 100 samples were annotated across all six prompt variations by three annotators. The low interrater reliability (IRR) scores in Table 1, based on Krippendorff’s Alpha (Krippendorff 2004), demonstrate that this type of evaluation was difficult to come to a consensus on. This could be due to the fact that the output quality metrics were all on a score of 0 to 2, as opposed to being binary values. These three values could also be ambiguous even with the specifications provided to the annotators. To avoid this in future studies, it would be beneficial to specify a larger number of metrics that capture more fine-grained linguistic details to improve agreement across annotators. We believe the results of the study to still be valuable when considering the goal is to present an approach to designing a conversational system specifically for the educational context.

Given the low IRR values, we report the mean and standard deviation for these values to get a better sense of how the agent responses were perceived by the annotators. We see the highest agreement levels across annotators for the adherence to discourse policy. However, the adherence to dialogue act constraint did not demonstrate as high agreement. We believe this can be attributed to the fact that raters considered this to be an “acceptability” annotation - i.e. is the output aligned with an “acceptable” dialogue act for a teacher response in the given context. Additionally, when annotating the outputs were compared to the actual teacher output dialogue act tags which were annotated with our classification prompt as opposed to gold standard tags. Furthermore, the teacher utterances were often extremely lengthy making it difficult for both a model or a human to identify a single correct dialogue act.

### 5.1 Quality Metrics: Fluency, Relevance, Helpfulness

One of the more interesting findings was that the impact of dialogue acts on the perceived fluency of response outputs. Annotators considered responses to be less fluent (a 1 over a 2) when the model contradicted itself, provided overly formal responses, or the phrasing was considered awkward. The inclusion of Student DAs with the state led to 117%

Prompt	Fluency	Helpfulness	Relevance	DA	DS
Baseline	1.410 $\pm$ 0.596	0.917 $\pm$ 0.755	1.143 $\pm$ 0.670	0.340	0.313
+ Context	1.370 $\pm$ 0.583	1.053 $\pm$ 0.777	1.260 $\pm$ 0.626	0.490	0.360
+ DS	1.390 $\pm$ 0.564	1.070 $\pm$ 0.765	1.297 $\pm$ 0.665	0.617	0.893
+ DS, Student DA	1.613 $\pm$ 0.500	1.237 $\pm$ 0.775	1.437 $\pm$ 0.648	0.603	<b>0.957</b>
+ Student DA, Teacher DA	1.580 $\pm$ 0.563	1.287 $\pm$ 0.803	<b>1.530</b> $\pm$ 0.585	0.623	0.567
+ DS, Student DA, Teacher DA	<b>1.653</b> $\pm$ 0.503	<b>1.320</b> $\pm$ 0.798	1.450 $\pm$ 0.659	<b>0.670</b>	0.903

Table 3: Experiment results for quality metrics and adherence metrics. Fluency, Helpfulness, and Relevance are scored based on the average mean across all 100 samples and three annotators plus or minus the standard deviation. Adherence to Dialogue Act (DA) and Dialogue State (DS) constraint is based on the percentage of the time that the raters marked samples as adhering to the constraint.

Prompt	BTT	OR
Baseline	0.500	0.122
Context	<b>0.733</b>	<b>0.222</b>
DS	0.633	0.156
DS, Student DA	0.689	0.200
Student DA, Teacher DA	0.70	0.144
DS, Student DA, Teacher DA	0.678	0.156

Table 4: **BTT** represents the percentage of the time the prompt output was rated as "Better than the Teacher" response to the student utterance. **OR** represents the percentage of the samples that the prompt was rated as the best overall response of the prompt outputs.

increase in the mean Fluency score over the inclusion of the dialogue state alone. The combination of Student DAs, Teacher DAs, and Dialogue State had the highest Fluency value, suggesting that when provided more constraints, the model produced more natural responses. The prompt variation including Teacher DAs and student DAs has a much higher Fluency rating than the baseline, but lower than combining student DAs with a dialogue state suggesting the state itself contributes to higher fluency. Annotators did frequently see outputs that included dialogue act tags when the Teacher DA was included, occasionally ones that were not even mentioned in the prompt, which would need to be removed before surfacing a response in an end-to-end system. The prompt including student and teacher DAs plus state information was evaluated as having the highest mean Helpfulness score. Additionally, when the results were broken down into 2-value pairs, with [0,1] rated as "not Helpful" or 0 and 2 being scored as Helpful, this prompt was marked as helpful over 50% of the time by raters. As with the Helpfulness and Fluency ratings, we notice that the prompts that included any type of DA

information had much higher mean scores for Relevance than those without suggesting that grounding the conversation in speaker intentions leads to more relevant responses. Unsurprisingly, the baseline has the worst performance in this category given that it does not include the conversation history in the prompt. The prompt that did not include the state information was scored higher in terms of Relevance over the version with this constraint. This could be due to the model trying to attend to too many constraints at once, or this could be related to the quality of the states themselves. There are an infinite number of hypothetical states that a classroom conversation could be in that a teacher may respond to differently and in future versions of the system we intend to explore more fine-grained state-action pairs. In general, all of the response quality metrics indicate that the inclusion of DA information does lead to better output responses from an LLM.

## 5.2 Adherence to Dialogue Act + Policy Constraints

For annotation, raters considered an output to 'adhere' to the dialogue act constraint if it was considered an acceptable dialogue act in the context of the conversation history. The adherence rate goes up even when just including the conversation history in the prompt over the baseline suggesting that there is some implicit dialogue flow information that the model is able to learn from the history itself. However, there is a meaningful jump in performance when additional constraints are applied, including the dialogue state without any DAs. Including the DA and discourse policy information resulted in the highest rates of adherence to this constraint. However, the adherence rate is still notably less than the agent's ability to

adhere to the dialogue state constraint. In the case of our dialogue policy, the adherence rate is higher than 89% of the time when the state intervention is included in the prompt. This is also the set of annotations with the highest rates of interrater reliability. Overall, this suggests that the Mistral model was extremely good at adhering to our discourse policy when provided with the appropriate information. This is especially true in comparison to the responses when no constraints are provided, and even in the drop in performance when using the student DA and teacher DA without the discourse policy constraint. These results suggest that as we refine the discourse policy to cover a wider range of classroom situations that we will be able to output responses that will adhere to it appropriately. This finding is incredibly important when considering the need to constrain agents in the classroom to be consistent with teacher behavior and reduce the risk of providing unsafe outputs to students.

### 5.3 Teacher and Prompt Comparisons

In addition to evaluating each of the prompts individually, we wanted to compare them to each other, and the actual teacher responses. We selected 30 samples from the 100 annotated samples above stratified to 10 per dialogue states. The 6 model outputs were shuffled to prevent the annotators from being biased towards a particular prompt style. For these, the raters considered two questions (1) Is this response better than the teacher response? (2) What is the best model response of the 6 provided? The results in Table 5 show that all of the non-baseline prompts were considered better than the teacher response in over half of the cases selected. Unexpectedly, the condition with only the conversation history, was rated as better than the teacher most frequently and considered the best response to the student utterance most frequently. The next highest rated prompt is the combination of student DAs and the dialogue state. In further discussion of the annotations, the raters mentioned that brief responses were considered better, and that the models typically provided explicit supportive feedback such as “Keep up the good teamwork” to students more often than the teacher did. Additionally, the actual teacher responses may have been addressing earlier conversation topics or other students as opposed to the most recent student utterance. The fact that one type of constraint did not inherently improve the overall perception of the responses suggests that the system may benefit

from the use of an over-generate and rank approach in which we provide an output from all of these variations and select one to provide to the student based on a set of criteria informed by a teacher.

## 6 Conclusion and Future Work

The goal of these experiments was to compare the effectiveness of different constraints in the context of zero-shot prompting a language model to provide teacher-like responses to real student utterances. We evaluated the inclusion of student and teacher dialogue acts annotations as well as dialogue states in the prompts provided to Mistral. A sample of 100 utterances was selected and evaluated across 6 different prompt settings by three annotators for Fluency, Relevance, Helpfulness, and adherence to the provided constraints. The inclusion of any type of constraints showed a positive impact on all of the utterance quality metrics, but could benefit from post processing to ensure that erroneous tags are not included in the agent response to a student. We saw the highest ratings for Fluency and Helpfulness when student dialogue acts, teacher dialogue acts, and dialogue states were provided in the prompt. The prompt version with all the constraints also had the second highest value for Relevance, suggesting that more contextual dialogue information in the prompt leads to higher quality responses from the agent. When provided a dialogue state in the input, the best performing prompt adhered to the constraint in over 95% of cases. This suggests that as future research is done identifying key pedagogical dialogue states and the ability to distinguish them from one another, LLMs can be very successful in adhering to the recommended states. Furthermore, a subset of 30 utterances were compared to the teacher response from the NCTE corpus. All of the LLM prompt variations that included the conversation history were considered to be better than the actual teacher response over 50% of the time. This was largely because the annotators found the model was more likely to be directly addressing the most recent student utterance as opposed to another student, and frequently included additional supportive phrases in the response. Ultimately, these findings suggest that LLM-based conversational agents have a lot of potential for providing learners with additional support in the classroom, when provided the appropriate constraints. In future work, we aim to refine the set of states from the three present in the NCTE

dataset along with educators to cover a broader set of scenarios. Additionally, we would like to explore model-based approaches for identifying the dialogue states themselves. The experiments here were limited to transcript-based annotations, in future work we intend to evaluate the responses in real scenarios with students.

## 7 Limitations

The scope of this paper is limited by a number of factors, including the types of models used and the types of constraints evaluated. We focused on a limited set of possible dialogue states based on initial expert analysis of classroom conversations. There are more possible dialogue states that would require a different type of support in the classroom. Additionally, the annotations were based on a set of classroom transcripts. Future work should evaluate the performance of such an agent in a live learning setting.

## 8 Ethics Statement

The experiments that we have conducted here are intended to improve the responses generated by LLMs for the classroom setting. However, the models that we use in our experiments are trained with large datasets that may be subject to unknown biases due to the exact content of the original training materials. Our research is intended to be used as a classroom support but this assumes that teachers will not use the information collected from dialogues to assess students' grades.

## Acknowledgments

This research was supported by NSF National AI Institute for Student-AI Teaming (iSAT) under grant DRL 2019805. The opinions expressed are those of the authors and do not represent views of NSF.

## References

Alexis Baladón, Ignacio Sastre, Luis Chiruzzo, and Aiala Rosá. 2023. RETUYT-InCo at BEA 2023 shared task: Tuning open-source LLMs for generating teacher responses. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 756–765, Stroudsburg, PA, USA. 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 [cs.CL]*.

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.

Andrew Caines, Helen Yannakoudakis, Helen Allen, Pascual Pérez-Paredes, Bill Byrne, and Paula Buttery. 2022. The teacher-student chatroom corpus version 2: more lessons, new annotation, automatic detection of sequence shifts. In *Proceedings of the 11th Workshop on NLP for Computer Assisted Language Learning*, pages 23–35, Louvain-la-Neuve, Belgium. LiU Electronic Press.

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. *arXiv [cs.LG]*.

Dorottya Demszky and Heather Hill. 2023. The NCTE transcripts: A dataset of elementary math classroom transcripts. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 528–538, Toronto, Canada. Association for Computational Linguistics.

Dorottya Demszky, Jing Liu, Zid Mancenido, Julie Cohen, Heather Hill, Dan Jurafsky, and Tatsunori Hashimoto. 2021. Measuring Conversational Uptake: A Case Study on Student-Teacher Interactions. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, and Bill Dolan. 2020. Dialogue response ranking training with large-scale human feedback data. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 386–395, Online. Association for Computational Linguistics.

- Umit Guz, Gokhan Tur, Dilek Hakkani-Tür, and Sébastien Cuendet. 2010. [Cascaded model adaptation for dialog act segmentation and tagging](#). *Computer Speech & Language*, 24(2):289–306.
- Behnam Hedayatnia, Karthik Gopalakrishnan, Seokhwan Kim, Yang Liu, Mihail Eric, and Dilek Hakkani-Tur. 2020. [Policy-driven neural response generation for knowledge-grounded dialog systems](#). In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 412–421, Dublin, Ireland. Association for Computational Linguistics.
- Yann Hicke, Abhishek Masand, Wentao Guo, and Tushaar Gangavarapu. 2023. Assessing the efficacy of large language models in generating accurate teacher responses. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 745–755, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Juraj Juraska, Kevin Bowden, and Marilyn Walker. 2019. [ViGGO: A video game corpus for data-to-text generation in open-domain conversation](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 164–172, Tokyo, Japan. Association for Computational Linguistics.
- Sanna Järvelä and Allyson F Hadwin. 2013. New frontiers: Regulating learning in CSCL. *Educ. Psychol.*, 48(1):25–39.
- Katharina Kann, Abteen Ebrahimi, Joewie Koh, Shiran Dudy, and Alessandro Roncone. 2022. Open-domain dialogue generation: What we can do, cannot do, and should do next. In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 148–165, Dublin, Ireland. Association for Computational Linguistics.
- K Krippendorff. 2004. Reliability in content analysis: Some common misconceptions and recommendations. *Hum. Commun. Res.*, 30(3):411–433.
- Mohammad Amin Kuhail, Nazik Alturki, Salwa Alramlawi, and Kholood Alhejori. 2023. Interacting with educational chatbots: A systematic review. *Education and Information Technologies*, 28(1):973–1018.
- Esther Levin, Roberto Pieraccini, and Wieland Eckert. 2000. A stochastic model of human-machine interaction for learning dialog strategies. *IEEE Transactions on speech and audio processing*, 8(1):11–23.
- Stefano Mezza, Alessandra Cervone, Evgeny Stepanov, Giuliano Tortoreto, and Giuseppe Riccardi. 2018a. [ISO-standard domain-independent dialogue act tagging for conversational agents](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3539–3551, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Stefano Mezza, Alessandra Cervone, Evgeny Stepanov, Giuliano Tortoreto, and Giuseppe Riccardi. 2018b. [ISO-standard domain-independent dialogue act tagging for conversational agents](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3539–3551, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Samuel L Pugh, Shree Krishna Subburaj, Arjun Ramesh Rao, Angela E B Stewart, Jessica Andrews-Todd, and Sidney K D’Mello. 2021. Say what? automatic modeling of collaborative problem solving skills from student speech in the wild. *International Educational Data Mining Society*.
- Alec Radford, Jeff Wu, R Child, D Luan, Dario Amodei, and I Sutskever. 2019. Language models are unsupervised multitask learners.
- Angela Ramirez, Kartik Agarwal, Juraj Juraska, Utkarsh Garg, and Marilyn Walker. 2023. Controllable generation of dialogue acts for dialogue systems via few-shot response generation and ranking. In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 355–369.
- Verena Rieser and Oliver Lemon. 2011. *Reinforcement learning for adaptive dialogue systems: a data-driven methodology for dialogue management and natural language generation*. Springer Science & Business Media.
- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Kurt Shuster, Eric M Smith, Y-Lan Boureau, and Jason Weston. 2020. Recipes for building an open-domain chatbot. *arXiv [cs.CL]*.
- Dapeng Shan, Deliang Wang, Chenwei Zhang, Ben Kao, and Carol Chan. 2023. [Annotating Educational Dialog Act with Data Augmentation in Online One-on-One Tutoring](#), pages 472–477.
- Guokan Shang, Wensi Ding, Zekun Zhang, Antoine Tixier, Polykarpos Meladianos, Michalis Vazirgiannis, and Jean-Pierre Lorré. 2018. [Unsupervised abstractive meeting summarization with multi-sentence compression and budgeted submodular maximization](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 664–674, Melbourne, Australia. Association for Computational Linguistics.
- Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000a. [Dialogue act modeling for automatic tagging and recognition of conversational speech](#). *Computational Linguistics*, 26(3):339–374.
- Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000b. Dialogue act modeling for automatic tagging and recognition of conversational speech. *Comput. Linguist. Assoc. Comput. Linguist.*, 26(3):339–373.

- Anaïs Tack and Chris Piech. 2022. The AI teacher test: Measuring the pedagogical ability of blender and GPT-3 in educational dialogues. In *Proceedings of the 15th International Conference on Educational Data Mining*, pages 522–529, Durham, United Kingdom. International Educational Data Mining Society.
- 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](https://github.com/tatsu-lab/stanford_alpaca).
- Henry S. Thompson, Anne Anderson, Ellen Gurman Bard, Gwyneth Doherty-Sneddon, Alison Newlands, and Cathy Sotillo. 1993. [The HCRC map task corpus: Natural dialogue for speech recognition](#). In *Human Language Technology: Proceedings of a Workshop Held at Plainsboro, New Jersey, March 21-24, 1993*.
- Justin Vasselli, Christopher Vasselli, Adam Nohejl, and Taro Watanabe. 2023. NAISTeacher: A prompt and rerank approach to generating teacher utterances in educational dialogues. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 772–784, Toronto, Canada. Association for Computational Linguistics.
- Marilyn Walker, Jeanne C Fromer, and Shrikanth Narayanan. 1998. Learning optimal dialogue strategies: A case study of a spoken dialogue agent for email. In *COLING 1998 Volume 2: The 17th International Conference on Computational Linguistics*.
- Nick Webb and Yorick Wilks. 2005. [Dialogue act classification based on intra-utterance features](#).
- Tom Williams, Cynthia Matuszek, Kristiina Jokinen, Raj Korpan, James Pustejovsky, and Brian Scasselati. 2023. Voice in the machine: Ethical considerations for language-capable robots. *Commun. ACM*, 66(8):20–23.
- 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.
- Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. 2023. A survey of controllable text generation using transformer-based pre-trained language models. *ACM Comput. Surv.*, 56(3):1–37.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. DIALOGPT : Large-scale generative pre-training for conversational response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278, Online. Association for Computational Linguistics.
- Tianyu Zhao and Tatsuya Kawahara. 2017. [Joint learning of dialog act segmentation and recognition in spoken dialog using neural networks](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 704–712, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Tianyu Zhao and Tatsuya Kawahara. 2019. [Joint dialog act segmentation and recognition in human conversations using attention to dialog context](#). *Computer Speech & Language*, 57:108–127.

## A Dialogue State Identification

1. We classify the utterances for Collaborative Problem Solving (CPS) codes that score them on three facets: Constructing Shared Knowledge, Negotiation/Coordination, and Maintaining Team Function (Pugh et al., 2021)
2. We check if the current utterance is on-topic for the lesson based on the “ontopic” feature of the NCTE dataset
3. Check for a literal match to one of our procedural issue sentences such as “What do we do next?” then return *PROCEDURAL* issue
4. Check the utterance is ontopic and the CPS codes are above a given threshold then return *CONTENT* issue
5. Check if the CPS codes is lower than a given threshold then return *COLLABORATION* issue
6. If there are no speakers, then we consider this to be a *NO SPEAKER* collaboration issue
7. If there is a single speaker, then we consider this to be a *SINGLE SPEAKER* issue
8. If there are multiple speakers then this is a general *COLLABORATION* issue
9. If none of the above conditions are met, the conversation is considered to be in a *FLOW* state

## B Dialogue Act Segmentation Prompts and Dialogue Acts

See [Table 5](#) below.

Within the actual prompt we used 19 examples. All examples and full prompt can be found within the github repository<sup>12</sup>.

<sup>12</sup><https://github.com/aramir62/constraints-generative-supportive-teacher-responses>

The dialogue acts used from ISO are: Set Question, Propositional Question, Choice Question, Inform, Commissive, Directive, Thanking, Apology, Salutation, and Feedback (default act).

From Switchboard, the feedback dialogue acts utilized are: Signal-Not Understanding, Appreciation, Sympathy, Summarize/Reformulate, Repeat-Phrase, Acknowledge (backchannel), and Acknowledge Answer.

### **C Response Generation Prompts**

See **Table 6** below.

### **D Generated Outputs**

See Table 6 below.

<b>Dialogue Act Segmentation Prompt</b>
<p>You are given an utterance to label a dialogue act with, and certain utterances need to be segmented if needed. You'll segment the utterance into as many chunks needed to capture all the dialogue acts, but if there are two chunks in a row with the same dialogue act combine them. Use punctuation and clause separators as a way to consider if this is a new topic or idea that would be labeled with a different dialogue act. Choose only one dialogue act per segment and you have to choose one. Definitions and examples of each given below. Only choose dialogue acts from this set, dialogue act names are contained in &lt;&gt;.</p> <p>&lt;Set Question&gt;: A question that focuses on the speaker wanting to know certain information often contains "wh" at the beginning such as "Where are you going?", "What did Kevin do yesterday?"</p> <p>&lt;Propositional Question&gt;: A question where the speaker wants to know if something is true or false, such as "Do you know what time it is?"</p> <p>&lt;Choice Question&gt;: A question that provides a list of options for someone to choose from, usually contains or and requires the other speaker to choose something, such as "Do you like peanut butter or chocolate more?"</p> <p>&lt;Inform&gt;: A statement that is providing information. Described as informative information to provide context, knowledge, and information about a topic.</p> <p>&lt;Commissive&gt;: The speaker will express a commitment to performing an action, such as "I will go home at 1pm"</p> <p>&lt;Directive&gt;: The speaker is directing or suggesting what will need to be committed to, such as "You need to go left then right down the hall"</p> <p>&lt;Thanking&gt;: A sentence that is expressing gratitude, such as "Thanks"</p> <p>&lt;Apology&gt;: An expression acknowledging regret or remorse towards an individual's actions, such as "I'm sorry"</p> <p>&lt;Salutation&gt;: An utterance made as a greeting or acknowledgment of another's arrival or departure, such as "Hello"</p> <p>&lt;Signal-Not Understanding&gt;: The speaker is expressing that they did not understand what was said prior. Or is using some form of non verbal language.</p> <p>&lt;Appreciation&gt;: The speaker is expressing feedback by providing appreciation towards the other speaker.</p> <p>&lt;Sympathy&gt;: The speaker is sympathetic towards the other speaker and is expressing this within the utterance.</p> <p>&lt;Summarize/Reformulate&gt;: The speaker summarizes or reformulates what was said before to demonstrate understanding of what was heard.</p> <p>&lt;Repeat-phrase&gt;: The speaker repeats back what was said beforehand.</p> <p>&lt;Acknowledge (backchannel)&gt;: The speaker expressing acknowledgement by using backchannels such as "Mmm hmm", "Mm yeah", "Uh huh"</p> <p>&lt;Acknowledge Answer&gt;: The speaker provides acknowledgement of what was said or asked prior by providing an answer.</p> <p>&lt;Feedback&gt;: An utterance that provides or elicit information about the type of understanding and processing of what was said prior, use this as a default if the other sets do not match and put as the lowest priority.</p> <p>Utterance: You need to give me your ideas and then I need to see whether that would sell in the marketplace. Output: You need to give me your ideas&lt;Directive&gt; and then I need to see whether that would sell in the marketplace&lt;Commissive&gt;</p> <p>Utterance: {utterance} Output:</p>

Table 5: **Dialogue Act Segmentation Prompt** with 1 example (prompt used for labeling used 19 examples.)

<b>Preamble:</b>
You will be given the task of generating a realistic response to a given statement or question. In order to complete this task successfully, you must pay careful attention to the following requirements and limitations. Failure to do so will result in failure of the task.
<b>Setting:</b>
For context, the given statement or question will have been spoken by a {LEVEL} student in the setting of a {SUBJECT} classroom. The student is working on an in-class assignment with a group of their peers. You will generate a response to what the student said as if you are the teacher of this {SUBJECT} class. As a teacher, your response must be: (1) kind, (2) helpful, (3) encouraging of collaboration, (4) relevant to the subject of {SUBJECT}, (5) factual, and (6) befitting the role of a professional classroom teacher. Please limit your response to the content of the student statement or question. These guidelines may help you:
<b>Formatting:</b>
Please format your response to the given statement or question as a SINGLE direct quote from the {SUBJECT} teacher whose role you are playing, including quotation marks. Please try to limit your response to {SOFT_CAP} words or less words if you can. However, if you are unable to accurately and coherently respond to the student's statement or question in {SOFT_CAP} words or less, you may use up to {HARD_CAP} words. Any response longer than {HARD_CAP} words will be considered a failure for the purpose of this task. Please do NOT include a word count in your response. Please do NOT provide any additional reasoning, explanations, or context along with your response. Please preface your response with the text "Response:" OUTSIDE of the quotation marks, but do NOT include any additional or alternative speaker attributions. Please do NOT preface or follow your response with ANY text other than the quotation marks and attribution required above.
<b>Student DA &amp; Teacher DA:</b>
The following are definitions for dialogue act labels, which are useful for understanding text: {LABEL_DEFINITIONS} To help you understand the student, their statement or question has been annotated with dialogue act labels. These labels have been inserted AFTER the section of text they describe, and each label will correspond to one of the definitions shown above. What follows is the annotated student quote: {STUDENT_SEGMENTATION} Furthermore, you should formulate your response such that it matches the definition(s) for the following label(s): {TEACHER_LABELS} However, the dialogue act labels are exclusively for your own use in understanding what the student has said and in formulating your response. Please do NOT include ANY actual dialogue act labels in your response.
<b>Dialogue Policy:</b>
{STATE_DESCRIPTION} <sup>3</sup> {CONSEQUENCES} {INTERVENTIONS}
<b>Context:</b>
To give you some insight into the on-going discussion, the following conversation history is provided to you: {HISTORY}
<b>Task:</b>
The following student quote is what you will be responding to: {STUDENT_UTTERANCE} Please generate the teacher's response according to the requirements and limitations provided above.

Table 6: **Response Generation Prompt** with all 8 component templates.

<b>State</b>	<b>Action/Prompt</b>
<b>Content Issue</b>	This group of students appears to be struggling with some aspect of the lesson material. The purpose of this activity is for them to learn the material; however, if they are struggling to understand the new information or have forgotten prior learning, they may be unable to finish the assignment. In order to help them make progress toward their current objective, your response should politely encourage them to compare notes and make sure that they agree on the lesson material.
<b>Collaboration Issue</b>	This group of students appears to be actively speaking, but none of them are making meaningful progress toward their current objective. This could be a problem, because class time is limited. If they spend too much time stuck or off-task, they may not complete their group assignment by the end of the class period. Your response should gently remind them of their current objective, politely encourage them to work collaboratively to achieve it, and perhaps suggest a way for them to get back on task
<b>Procedural Issue</b>	The students are struggling to understand the steps required to complete the task. Please encourage them to refer to the Sensor Immersion tutorials. This group of students appears to have found a good working rhythm and are making meaningful progress toward their current objective. They are doing a good job of collaborating with one another, and you do not need to intervene. Therefore, your response should praise their enthusiasm and encourage them to remain on-task and continue to make steady progress.

Table 7: A breakdown of the action prompts used for the content, collaboration, and procedural issue dialogue states. The procedural issue is not present in the NCTE dataset and has reference to the Sensor Immersion task used in the dialogue policy creation sessions.

<p><b>Student Utterance</b> [Flow] So, not in the envelope?&lt;Feedback&gt;</p>	<p><b>Conversation History</b></p> <p>teacher: "Anybody? Nobody? What did you use them for, Student C."  student: "To help us with our fractions."  teacher: "To help you with your fractions. Do you guys understand that sometimes it's easier to just see it and touch it and hold it? Today, the job that we have is that we are going to learn to simplify fractions and we are going to use these pieces to help us to simplify them. So, by the time I am done, I have 1, 2, 3, 4, 5, 6, 7, 8, 9, piles of fraction pieces. And one of them is one whole strip. And all of your scraps need to go into the blue recycling bin, please, which I just put right here. I'm not so worried about your names going on the backs of them, boys and girls, because you have your own little toolbox, the envelope you are going to keep them in, so your name does need to go on the envelope."</p>
<p><b>Prompt Variation</b> Actual Teacher Response  Baseline  + Context  + DS  + DS, Student DA  + Student DA, Teacher DA + DS, Student DA, Teacher DA</p>	<p><b>Generated Response</b></p> <p>I'm sorry, what?&lt;Signal-Not Understanding&gt;/ It does go on the envelope.&lt;Inform&gt;/ We need to get moving here.&lt;Directive&gt;/ It's the noise.&lt;Inform&gt;/ We don't need to do that.&lt;Directive&gt;/</p> <p>"That's correct, the problem does not involve finding the area or volume of an envelope."  "That's correct. The fraction pieces themselves don't need names on them, but be sure to label your envelopes."  "That's correct, Student. The names don't need to be on the fraction pieces themselves, only on the envelopes. Keep up the good work."  "Great question! Yes, you're correct, any scraps or unused pieces don't go in the envelope. Keep up the good work, team!"  You're correct, those pieces that don't fit in the envelope should be placed in the recycling bin. Keep up the good work, team!  "Yes, that's correct. The scraps don't need to go in the envelope, they can go in the blue bin." &lt;Inform, Directive&gt;</p>

<p><b>Student Utterance [Collaboration Issue]</b> That's not how I write my fractions though so&lt;Feedback&gt;/</p>	<p><b>Conversation History</b> teacher: "Okay, all right. And then if you finish, just ask yourselves. Come up with a fraction and ask yourself, you know, how you can compare." teacher: "Okay, hard to write on those little lines, isn't it?" teacher: "Bless you."</p>
<p><b>Prompt Variation</b> Actual Teacher Response Baseline  + Context  + DS  + DS, Student DA  + Student DA, Teacher DA  + DS, Student DA, Teacher DA</p>	<p><b>Generated Response</b> All right,&lt;Acknowledgement&gt;/ are you girls finished?&lt;Set Question&gt;/  "I understand that everyone has their own way of representing fractions. Let's discuss how we can all use the same convention for this assignment to ensure clarity." "I understand, everyone has their own way of writing fractions. Let's compare our methods and learn from each other." "Understood. Let's try using the standard fraction format for consistency within your group to ensure effective comparisons." "I understand that everyone has different ways of representing fractions. Let's collaborate and share our methods to help us all learn. Let's focus on understanding the concept of comparing fractions." "It's perfectly fine if we each have slightly different ways of writing fractions as long as the meaning is clear. Let's discuss how we can compare them effectively." "It's great that we all have different ways of expressing ideas. Let's see if we can use each other's methods to help us understand how to compare fractions effectively. Let's continue working together."</p>

<b>Student Utterance</b> <b>[Content Issue]</b>	<b>Conversation History</b>
Because a straight angle goes on, on both sides.<Inform>/	teacher: "Let's see. All right. Good. Straight angles. Let's get our pencils ready because this is one you need to do. Some of you may have heard of this. How many of you all have heard of a straight angle? Very small. All right. A straight angle is a line that goes infinitely in both directions. It measures 180 degrees. Okay. Yes, Student L? All right. So again, which is not a straight angle? Any others? What do you think, based on definition?" student: "Number two." teacher: "How do you know number two is not a straight angle?"
<b>Prompt Variation</b>	<b>Generated Response</b>
Actual Teacher Response	It goes on forever. There's a line, and it goes on forever on both sides, absolutely.<Inform>/ Of course, it's a ray so it only goes in one direction.<Inform>/ Now we can talk about our main guy today.<Inform>/ Our main guy today is a triangle.<Inform>/ We all know triangles, right?<Acknowledge Answer>/
Baseline	"That's an interesting way to think about it. Can you clarify what you mean by 'a straight angle goes on, on both sides' in the context of geometry?"
+ Context	"That's correct, Student. A straight angle does indeed go on indefinitely in both directions. Well done for applying the definition."
+ DS	"That's correct, Student, a straight angle does indeed extend indefinitely in both directions. Let's discuss with your group to ensure everyone agrees on this definition."
+ DS, Student DA	"That's correct, a straight angle extends infinitely in both directions. Let's compare our definitions and work together to find other angles in the diagram that are not straight angles."
+ Student DA, Teacher DA	"That's correct, Student. A straight angle extends indefinitely, meaning it continues on both sides, which is why it measures 180 degrees."
+ DS, Student DA, Teacher DA	"That's correct. A straight angle continues indefinitely on both sides, making it an important concept in geometry. Let's discuss this with your group and ensure everyone agrees on this definition."

Table 8: A comparison of the generated utterances from the different prompt variations for each of the dialogue states.

# HelloThere: A Corpus of Annotated Dialogues and Knowledge Bases of Time-Offset Avatars

Alberto Chierici, Nizar Habash

Computational Approaches for Modeling Languages (CAMEL) Lab

New York University Abu Dhabi

{alberto.chierici, nizar.habash}@nyu.edu

## Abstract

A Time-Offset Interaction Application (TOIA) is a software system that allows people to engage in face-to-face dialogue with previously recorded videos of other people. There are two TOIA usage modes: (a) creation mode, where users pre-record video snippets of themselves representing their answers to possible questions someone may ask them, and (b) interaction mode, where other users of the system can choose to interact with created avatars. This paper presents the HelloThere corpus that has been collected from two user studies involving several people who recorded avatars and many more who engaged in dialogues with them. The interactions with avatars are annotated by people asking them questions through three modes (card selection, text search, and voice input) and rating the appropriateness of their answers on a 1 to 5 scale. The corpus, made available to the research community, comprises 26 avatars' knowledge bases and 317 dialogues between 64 interrogators and the avatars in text format.

## 1 Introduction

Time-Offset Interaction Applications (TOIAs) have evolved as an innovative dialogue system, bridging the interaction between individuals and pre-recorded video representations of others, hence enabling users to hold conversations outside real-time constraints (Artstein et al., 2015; Traum et al., 2015; Abu Ali et al., 2018). We built on an open-source project's application, offering a dual interface targeting two distinct user groups: (a) avatar creators, individuals interested in generating their time-offset personas, and (b) interactors, those who engage with these avatars.

However, designing a robust TOIA is a challenging endeavor. The goal is to mirror human-to-human interactions as authentically as possible. This demands seamless integration from an engineering standpoint, such as flawless video clip

transitions and numerous linguistic and dialogue-turns complexities that intrigue dialogue system researchers. Central to a TOIA's functionality are the avatar's Knowledge Bases (KBs), repositories of questions paired with corresponding video responses and their transcriptions. One of the inherent challenges is devising an optimal strategy for populating this KB. Should it be intuition-driven, or should it stem from authentic dialogue transcripts? Furthermore, what data sets can be useful for training models to retrieve the right answer for an interrogator interacting with the avatars? While we explored such questions in other research (Chierici et al., 2020; Chierici and Habash, 2021, 2023), here we focus on building on such body of work and present the language resources generated in the process. We explored KBs created in three distinct ways: intuition-guided (brainstormed), led by automatic suggestions (generated by GPT-3), and led by human suggestions. We used GPT-3 because our software and study were designed and set up between 2022 and 2023 before newer versions were available. The **HelloThere Corpus** offers a unique resource for dialogue researchers, enabling studies on multi-modal interactions, user engagement patterns, and the effectiveness of time-offset avatar responses. By providing annotated dialogues across different interaction modes, this corpus supports research into natural language understanding, response retrieval and generation, and user experience in asynchronous communication systems.

## 2 Related Work

We categorize pertinent literature on Time-Offset Interaction Applications (TOIA) into three primary areas: System Approaches, Data Sources, and Evaluation Methodologies.

## 2.1 System Approaches

Our work builds upon the foundations laid in [Chierici et al. \(2020\)](#); [Chierici and Habash \(2021\)](#); [Chierici et al. \(2021\)](#), whose initial inspiration stemmed from the work of [Traum et al. \(2015\)](#) in their New Dimensions in Testimony project. While Traum et al. created a time-offset interaction with Holocaust survivor Pinchas Gutter, we extend their approach to different contexts and focus on system scalability. The TOIA open-sourced in ([Chierici et al., 2021](#)) aims to operate with fewer recorded statements, adapt to multiple users, and facilitate getting to know a stranger in a 10- to 15-minute interaction.

Following the taxonomy we proposed in [Chierici et al. \(2020\)](#), we work on a novel subcategory of ‘self-narrative bots,’ which can be seen as an intermediate between social and task-driven bots, leveraging both structured and unstructured training data ([Gao et al., 2019](#)). Retrieving the appropriate video from a TOIA Knowledge Base (KB) shares similarities with FAQ retrieval, a dichotomous problem. While its single-turn question-answer (q-a) mechanism may seem rudimentary, tasks like search and Retrieve-And-Generate – where a model retrieves relevant information and generates a response based on it – introduce complexities due to the dynamic nature of dialogue ([Mass et al., 2020](#); [Yehudai et al., 2023](#)).

As the dataset scales, classification approaches may falter, highlighting the presence of long-tail problems and the challenges of chit-chat scenarios, where queries can have subtle differences (e.g., “What is your name?” vs. “What is your parent’s name?”). Technologies involved range from traditional RNN models and word embeddings to newer language models like OpenAI’s GPT families, Mistral, Llama and Nomic ([Radford et al., 2018](#); [Zhang et al., 2022](#); [Touvron et al., 2023](#); [Jiang et al., 2023](#); [Nussbaum et al., 2024](#)).

Recent advancements in neural architectures have led to cutting-edge performance in answer retrieval tasks, but the limited scale of our dialogue datasets—and those of similar scope—does not readily support deep learning approaches. This limitation does not preclude using pre-trained large language models for sentence similarity tasks, leveraging or not transfer and few-shot learning techniques.

While TOIAs share some similarities with recent advancements in speech and video synthesis tech-

nologies, they differ in their focus on preserving authentic human responses. Unlike synthetic systems that generate responses in real-time, TOIAs rely on pre-recorded human responses, maintaining the nuances of human communication. However, the retrieval mechanisms in TOIAs can benefit from advancements in natural language processing used in synthetic systems, particularly for improving response selection accuracy.

## 2.2 Data Sources

Various datasets have been employed to tackle problems related to chit-chat and question answering in dialogue systems, such as SQuAD ([Rajpurkar et al., 2016](#)), the Ubuntu dialogue corpus ([Lowe et al., 2015](#)), and bAbI ([Weston et al., 2015](#)). However, these works address tasks like time-based reasoning and logical induction, which differ from the context of TOIAs. The landscape of dialogue-focused datasets is evolving to capture complexities absent in earlier reading comprehension collections. Datasets like CoQA and HUMOD are designed with human dialogues and annotations in mind, enhancing natural conversational elements ([Reddy et al., 2019](#); [Merdivan et al., 2020](#)). Similarly, the Douban Conversation Corpus offers insights into real-world social discussions on various topics ([Wu et al., 2016](#)).

While large-scale datasets serve various purposes, dialogue systems often operate with far smaller datasets. For instance, the Margarita Dialogue Corpus (MDC) features a Knowledge Base (KB) with only 431 answers and complete annotated dialogues ([Chierici et al., 2020](#)). The nuanced context of dialogue in TOIAs demands different, more tailored datasets. The MDC offers a unique blend of structured and unstructured dialogues for time-offset interactions. While influential for our work, it is limited to a single avatar and real person-to-person transcripts, not mediated through a TOIA interface. This work extends the MDC by incorporating more avatars and collecting extensive real-world interactions with them, addressing identified limitations and enriching the corpus.

In previous work ([Chierici and Habash, 2023](#); [Chierici, 2023](#)), we addressed a key challenge in TOIA development – the daunting task of creating extensive video-anchored question-answer (q-a) pair databases without overwhelming the avatar maker, improving upon [Chierici et al. \(2020\)](#). We introduced Question Suggester (QS), a GPT-3-based intelligent service designed to alleviate

this problem by dynamically suggesting relevant follow-up questions based on the existing conversation history, significantly reducing the effort required to populate the video database and enhancing user experience.

### 2.3 Evaluation

We acknowledge that evaluating dialogue systems is a complex task, as traditional metrics often fail to correlate with human judgment, which itself is challenging to quantify (Li et al., 2019). The corpus we present addresses some gaps highlighted in Chierici and Habash (2021), where we performed a human evaluation study with a fictional TOIA interface and Amazon Mechanical Turk raters. We deployed the open-source software described in Chierici et al. (2021), with updates to the dialogue systems module and user interface, and built datasets using real TOIA-interactions. Participants were tasked with getting to know the avatar creator within a 10-minute interaction, evaluating each response as they interacted with the tool.

## 3 Data Acquisition and Annotation

Our work resulted in collecting and annotating dialogue data comprising 2.2 million words. This effort was part of a large user study involving 90 individuals, some who built the avatars, and others evaluated their interaction quality, along with testing and evaluating a few software features and related research questions discussed in (Chierici and Habash, 2023). Ethical considerations were upheld as our institution’s Institutional Review Board approved the experiments, and participants consented to release data transcriptions, annotations, and video recordings for research purposes only. In both parts of the study, participants were university students recruited via an online form that included informed consent and details about the study. In the first part, 26 individuals aged 18-24 participated, with 14 females and various international provenance. They are fluent in English and major in various fields, mostly science. In the second part, 64 people participated. They were mostly between the ages of 18-23, and 35 were female. All are also fluent in English, though 80% consider it their second language. The majority were science majors, and a subset of 16 had participated in the previous part of the study. To clarify how data is collected, we describe the user interface used in the extensive user study that generated the corpora.

## 3.1 User Interface

The user interface (UI) components are: (see Fig. 1)

**1. User Account** (Fig. 1 (a)): This is the initial page that users see after creating an account. It displays a button to create new videos, suggested questions for creating new videos, and videos previously recorded by the user.

**2. Recorder** (Fig. 1 (b)): This page is accessed by clicking on the buttons to add a new video or edit a previously recorded video or a suggested question in the User Account page. This is where users can create new videos by typing a question and hitting the record button. The system automatically transcribes what the user says, and the user can edit the transcriptions before saving the video. Once a video is saved, the user interface shows a pop-up menu (Fig. 1 (c)) with the command for creating a new video and follow-up question suggestions.

**3. Player** (Fig. 1 (d)): Here, users can interact with previously recorded videos of public *TOIA avatars*. The player interface comprises a video looping different ‘filler’ videos—clips without audio, where the *TOIA avatar* does not speak. Users can click on suggested questions displayed on the right side of the video, triggering an immediate response from the *TOIA avatar*. We call this interaction type ‘CARD’ in our later data description. Users can also ask questions verbally using a voice input button, and they are then transcribed and matched to appropriate responses. There’s a button to interact with the *TOIA avatar* by voice (marked as ‘VOICE’ in the data), and below that button, a text input field allows users to type in their questions, which are then matched to the most relevant pre-recorded response (interaction labeled as ‘TYPE’ in the data). These interaction modes offer flexibility in how users engage with the avatars, catering to different preferences and contexts. <

## 3.2 Creating Avatars

The first step of our user study focused on evaluating the methodology for creating avatars, using both qualitative and quantitative approaches. A key aspect of this evaluation was examining the impact of different question generation methods (for a more detailed discussion of this, we refer readers to the publication presenting the user experience study, Chierici and Habash (2023)). Metrics include the efficiency of avatar creation, the quality

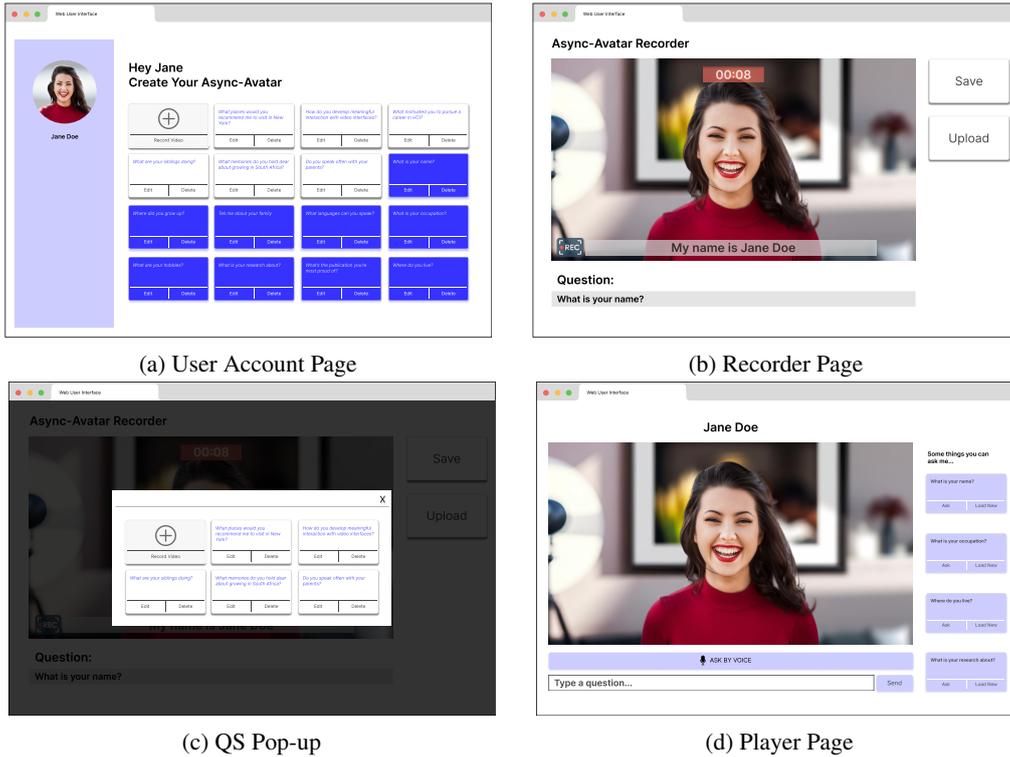


Figure 1: User Interface (UI) designs. These are similar to what we used when collecting data, though the actual UI has since evolved. (a) is the user account page showing the QS in white backgrounds and previously recorded questions (and videos) shaded in blue; (b) is the recorder page; (c) shows suggestions appearing in a pop-up window once the user completes a recording on the Recorder page; and (d) is the player page.

of suggested questions, and the influence of the creator’s personality traits on user acceptability and interface interaction. Three experimental conditions were examined when creating *TOIA avatars*’ KBs: 1) GPT-3-based question suggestions (GPT-3 QS), 2) human-curated questions (Human-QS), and 3) a no-suggestion, brainstorming condition (QS-off). As a result, 26 avatars were crafted: 10 through GPT-3 QS, 8 via Human-QS, and 8 using the QS-off approach.

### 3.3 Avatar Interaction

In the second step of the user study, to investigate key interaction metrics, including the minimum number of videos needed for a satisfying experience, variants of the original 26 avatars were created. These variants were based on three conditions concerning video count (first 30, first 60, or all recorded videos) and two filler videos (attentive or inattentive) types. Thus, each original avatar spawned 6 distinct interaction variants, leading to 156 unique avatars. We aimed to collect at least two evaluations for robust statistical analysis for each, totaling 312 unique dialogue interactions (to satisfy some experimental constraints and replace

participants who withdrew, we ended up with 317 dialogues in total).

### 3.4 Single-turn Answer Retrieval

We employ the GPT-3 model family from OpenAI for the retrieval task, specifically geared for semantic similarity-based text search (Neelakantan et al., 2022).<sup>1</sup> This choice was informed by the model’s superior performance tested on the Margarita Dialogue Corpus (Chierici et al., 2020). In our setup, q-a pairs are documents and converted into 1024-dimensional vector embeddings using the ‘text-search-ada-doc-001’ model. Incoming user queries are similarly transformed into 1024-dimensional vector embeddings through the ‘text-search-ada-query-001’ model. The Dialogue Manager (DM) suggests an answer when the cosine similarity between the query and document vectors exceeds a threshold of 0.29. If the similarity falls below this cutoff, the DM defaults to a predetermined set of videos intended for situations where no appropriate answer exists, such as “I haven’t recorded an answer for that question.” Our dia-

<sup>1</sup>For implementation guidelines, see <https://beta.openai.com/docs/guides/embeddings>.

logue data set also reports the similarity measure for each answer played out as a response to the interactors’ questions.

### 3.5 Annotations

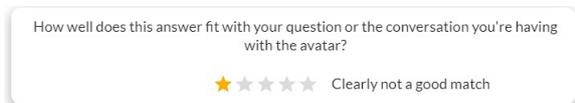


Figure 2: On the Player interface, a pop-up appears after every answer is played. The interactor has to rate the answer before going ahead with asking the next question.

We have three kinds of annotations. First, the Knowledge Base (KB) of each avatar who linked a question with an answer. Second, we collect the questions the automated and human QS suggested and mark them as selected or rejected by the avatar maker when creating their video recordings. Third, we have 64 human subjects who conversed with an avatar variant for a minimum of 10 minutes. We employed a 5-point rating scale, triggered by a pop-up after each video-based answer, to collect user assessments (Figure 2). Participants interacted with at least four different avatars (barring a few exceptions, who interacted with eight and one person just with one avatar).

Key conditions for the experimental design include:

- Each avatar variant received evaluations from at least two different participants.
- Participants never interacted with the same avatar more than once.
- Variants with different numbers of videos require separate evaluations.
- Filler video types were not considered separate conditions, allowing for collective evaluations.
- Interaction methods were flexible: participants could ask questions through text, voice, or preset options shown on the right of the player page by clicking on them (Fig 1d (d)).

## 4 Data Description and Exploration

Data for this study is accessible on NYUAD CAMEL Lab’s Resource page.<sup>2</sup> We present the

<sup>2</sup><http://resources.camel-lab.com/>

summary statistics of the two main language resources, ‘Knowledge Base’ and ‘Dialogues’, in Tables 1 and 2. We then discuss the agreement between annotations, a baseline retrieval evaluation, and a qualitative assessment of the topics covered in the corpora.

### 4.1 Avatar Knowledge Bases

In the first part of the human subject study (Table 1), the data generated encompasses 26 distinct avatars, each with a unique set of q-a pairs and dialogues. The data is structured into three cohorts: GPT-3-QS, Human-QS, and QS-Off, providing us with a rich platform to compare avatar behavior and performance across different conditions. The choice to create 26 distinct avatars was made to balance depth and breadth in our corpus. This number allows for a diverse range of personalities and interaction styles while remaining manageable for detailed analysis and within budget and time constraints. The distribution across different question suggestion methods (10 GPT-3 QS, 8 Human-QS, and 8 QS-off) enables comparative studies on the effectiveness of these approaches in creating engaging and comprehensive avatar knowledge bases. Here, we describe general insights and patterns observed across the three cohorts.

The corpus comprises 3,548 q-a pairs across all 26 subjects, with an average of 136.5 per subject. The data set encompasses 606,458 words, with an average of 43.1 words per question and longer answers (127.9 words on average).

The ‘answer’ category is overwhelmingly prevalent, constituting 2,407 of the q-a pairs—averaging about 92.6 per subject. This dominance underscores the avatars’ primary role: to deliver informative and substantive responses. The Human-QS cohort exhibits the highest word count per answer, indicative of more elaborate and nuanced responses.

The Human-QS cohort answers are the longest, followed closely by those of the QS-Off cohort. Categories like ‘exit,’ ‘greeting,’ ‘no-answer,’ and ‘y/n-answer’ are relatively (and obviously) rare across all cohorts. However, they exhibit diversity in terms of average word count. These categories might be infrequent but serve specific roles within the dialogic interaction and should not be overlooked.

### 4.2 Dialogues

Dialogues offer a more dynamic measure of conversational capabilities and limitations, allowing

	Total	By video-type					
		answer	exit	filler	greeting	no-answer	y/n-answer
<i>All (N=26 Subjects)</i>							
# q-a pairs	3,548	2,407	47	696	49	157	192
(Avg./subject)	136.5	92.6	1.8	26.8	1.9	6.0	7.4
# words	606,458	536,318	2,600	43,784	1,645	12,560	9,551
Avg. # words/question	43.1	40.1	31.4	59.3	22.8	32.7	38.2
Avg. # words/answer	127.9	182.8	23.9	3.6	10.7	47.3	11.5
<i>GPT-3-QS Cohort (N=10 Subjects)</i>							
# q-a pairs	1,538	1,067	20	284	18	70	79
(Avg./subject)	153.8	106.7	2.0	28.4	1.8	7.0	7.9
# words	251,522	223,504	1,127	17,518	561	4,815	3,997
Avg. # words/question	43.0	40.9	31.2	58.8	22.9	28.7	35.9
Avg. # words/answer	120.5	168.6	25.2	2.9	8.3	40.1	14.7
<i>Human-QS Cohort (N=8 Subjects)</i>							
# q-a pairs	1,094	791	12	198	16	41	36
(Avg./subject)	136.8	98.9	1.5	24.8	2.0	5.1	4.5
# words	218,935	197,552	739	13,555	641	4,269	2,179
Avg. # words/question	45.2	41.0	36.5	64.2	25.3	41.0	50.9
Avg. # words/answer	154.9	208.8	25.1	4.3	14.8	63.1	9.6
<i>QS-Off Cohort (N=8 Subjects)</i>							
# q-a pairs	916	549	15	214	15	46	77
(Avg./subject)	114.5	68.6	1.9	26.8	1.9	5.8	9.6
# words	136,001	115,262	734	12,711	443	3,476	3,375
Avg. # words/question	40.5	37.1	27.7	55.4	20.1	31.3	34.7
Avg. # words/answer	108.0	172.8	21.2	4.0	9.4	44.2	9.2

Table 1: Summary statistics on the data sets collected in the user study on the avatar creation. Statistics for the various *TOIA avatars*’ knowledge bases are also shown for each video type and by the experimental condition cohort (Question Suggester powered by GPT-3, by a human, and switched off).

	Tot	By Interaction Type		
		CARD	SEARCH	VOICE
# dialogues	317			
# q-a pairs	9,684	2,955	2,579	4,150
# no-answers	792	17	182	593
(in %)	8.2%	0.6%	7.1%	14.3%
# words	1,602,582	581,826	426,964	593,792
Avg. # turns/dialogue	30.5	9.3	8.1	13.1
Avg. # words/question	32.5	38.8	31.9	28.3
Avg. # words/answer	133.0	158.1	133.7	114.8

Table 2: Summary statistics on the dialogues collected from the interaction user study’s chat logs. Statistics are also shown for each type of interaction with the player interface (CARD, SEARCH, VOICE).

Mode	#	%	Mean	StDev	Min	25%	50%	75%	Max
CARD	2,851	31.3	4.6	0.9	1.0	5.0	5.0	5.0	5.0
SEARCH	2,459	27.0	3.9	1.6	1.0	3.0	5.0	5.0	5.0
VOICE	3,790	41.7	3.5	1.6	1.0	2.0	4.0	5.0	5.0
Total	9,100	100.0	4.0	1.5	1.0	3.0	5.0	5.0	5.0

Table 3: Distribution of interactors’ ratings by mode of interaction from the conversation log data of our *TOIA*.

for deeper understanding beyond individual, single-turn questions and answers. The data on dialogues is grouped into two key tables: Table 2 captures metrics by interaction type, while Table 3 focuses on annotations results (retrieval ratings) by mode.

The data set encompasses 317 dialogues, unfolding over 9,684 q-a pairs. These pairs are distributed across CARD (2,955), SEARCH (2,579), and VOICE (4,150) interactions. The 'No-Answers' account for 792 pairs or 8.2% of the total interactions. The dialogues encompass just over 1.6 million words, with an average of 30.5 turns per dialogue, 32.5 words per question, and 133 words per answer. The average of 30.5 turns per dialogue implies that the conversations are not just transactional but likely complex and multilayered.

VOICE-based interactions comprise the bulk of the dataset with the highest number of q-a pairs and a 14.3% 'No-Answers' rate. This suggests that voice interactions are frequent and more susceptible to information gaps or misunderstandings. The exceptionally low 'No-Answers' rate in CARD interactions (0.6%) is a consequence of the more scripted or straightforward engagement due to a deterministic retrieval (it is not 100% deterministic because the suggested cards are retrieved using prompting GPT-3 text completion and not always the underlying questions are reproduced verbatim).

CARD interactions have the highest average words per answer at 158.1, indicating a propensity for asking questions with more detailed responses in this particular mode of interaction –perhaps these are less trivial or less mundane questions that users wouldn't ask if they didn't see the suggestion on the card.

Looking at Table 3, the mean rating stands at 4.0 across all interactions with a standard deviation of 1.5. The scores range from a minimum of 1.0 to a maximum of 5.0. While VOICE accounts for 41.7% of all interactions, it has the lowest mean score of 3.5 and the same standard deviation as SEARCH. This follows from VOICE being the interaction that mostly depends on answer retrieval algorithms to provide answers. In contrast, CARD interactions have the highest mean score of 4.6 and a low standard deviation of 0.9. SEARCH interactions yield a mean of 3.9 and a slightly higher standard deviation, indicating a middle ground between VOICE and CARD. A mean score of 4.0 suggests that while the system performs reasonably well, raters may be particularly generous, and there remains scope for targeted improvements. Given the

Coefficient	Value (C.I.)	p-value
Gwet's AC1	0.82 (0.64, 1.00)	$1.66 \times 10^{-13}$
Fleiss Kappa	0.79 (0.61, 0.97)	$1.85 \times 10^{-13}$
Brennan-Prediger	0.81 (0.63, 1.00)	$8.35 \times 10^{-14}$
Conger's kappa	0.76 (0.57, 0.94)	$6.26 \times 10^{-12}$

Table 4: Inter-annotator agreement computed using coefficients of agreement that are all relevant in our scenarios where we have multiple raters using ordinal ratings.

high volume but variable quality, the VOICE category could benefit from refined natural language understanding algorithms to reduce 'No-Answers' and improve consistency.

### 4.3 Retrieval Evaluation Results

The interaction experiment yielded a total of 9,100 q-a pairs, with the summary statistics and answer ratings across different interaction modalities presented in Table 3. The data show that the voice modality was the most frequently utilized method of interaction, accounting for 41.6% of the cases. This was followed by clicking on suggested questions (31.3%) and typing (27.0%). However, frequency of use does not necessarily indicate user preference. Collectively, quicker interaction modalities like clicking and typing were used more often, comprising 58.4% of the interactions.

Anomalies in the CARD mode were observed despite its deterministic nature. Although it garnered the highest average rating, some users still rated answers poorly. Closer observation revealed that misclicks and inattentiveness during ratings were the primary causes of these anomalies. The SEARCH mode revealed similar variability in user ratings, echoing the patterns observed in the VOICE mode. Due to limitations in our log data, we restricted our analysis to the SR@1 performance in VOICE interactions. Qualitative insights suggest that participants often switched between the three modalities during a conversation, primarily initiating voice interactions.

We measured retrieval success with Success Rate@1 (SR@1) based on two scenarios: including neutral ratings (3, 4, and 5), which resulted in an SR@1 of 68.2%, and only considering high ratings (4 and 5), which yielded an SR@1 of 54.5%.

### 4.4 TOIA Interaction Rater Agreement

Inter-rater agreement was assessed on a small sample and is reported in Table 4. To identify equal instances rated by multiple interactors, paraphrased

Theme (Short Name) & Sample Question
<b>Opinion and personal beliefs (Opinion)</b> <i>Do you believe in second chances?</i>
<b>Reflection, Self-Awareness, Goals (Reflection)</b> <i>If you were to die this evening with no opportunity to communicate with...</i>
<b>Student Life, Major, Education (Education)</b> <i>What are you studying right now?</i>
<b>Food (Food)</b> <i>What is your favorite dish at Circle Cafe?</i>
<b>Preferences, Interests, and Lifestyle (Lifestyle)</b> <i>What is the most crucial element in a balanced life?</i>
<b>Cities, Countries and Travel (Travel)</b> <i>Are you interested in traveling to Australia?</i>
<b>Music (Music)</b> <i>Can you recommend some songs you like?</i>
<b>Books, Movies, TV (Media)</b> <i>What was the last tv series you binge watched?</i>
<b>Personal experiences, opinions, and advice (Advice)</b> <i>When was your first kiss?</i>
<b>Name, Age, Birthplace, Location (Identity)</b> <i>How old are you?</i>
<b>Family (Family)</b> <i>What is your family like?</i>
<b>Hobbies, Pastimes (Hobbies)</b> <i>What's your favorite way to spend a day off?</i>
<b>Animals (Animals)</b> <i>If you could have an animal sidekick, what would it be and why?</i>
<b>Abu Dhabi (AbuDhabi)</b> <i>How is living in Abu Dhabi?</i>
<b>Sports (Sports)</b> <i>Are you involved in sports?</i>
<b>Job, Career Aspirations, Plans After Graduation (Career)</b> <i>What do you want to do after graduation?</i>
<b>People Qualities and Characteristics (Traits)</b> <i>What do you value in people?</i>
<b>Greetings (Greetings)</b> <i>Hello!</i>
<b>Missing Home (Home)</b> <i>Do you miss home?</i>
<b>Time (Time)</b> <i>What time do you...</i>
<b>Miscellaneous, Trivia (Trivia)</b> <i>Morgan supporting in the World Cup...</i>
<b>Language (Language)</b> <i>How many languages do you speak?</i>

Table 5: Summary of the topic clustering for questions asked by voice.

questions were grouped using cosine similarity of their sentence embeddings and checked manually to identify groups of the same question asked. A heuristically inspected threshold of 0.87 +/- 0.003 was used to cluster similar questions, leaving us with 86 comparable instances.

We computed four coefficients, namely Gwet’s AC1, Fleiss Kappa, Brennan-Prediger, and Conger’s kappa, to measure the agreement level. All coefficients indicated significant levels of agreement (see Table 4 for numerical results and p-values).

Lastly, we observed a correlation coefficient 0.44 (p-value:  $1.03 \times 10^{-153}$ ) between the retrieval results and the interactors’ ratings. This stronger correlation compared with the work of (Chierici and Habash, 2021) underscores a higher agreement between the retrieved responses and human opinions in our setup.

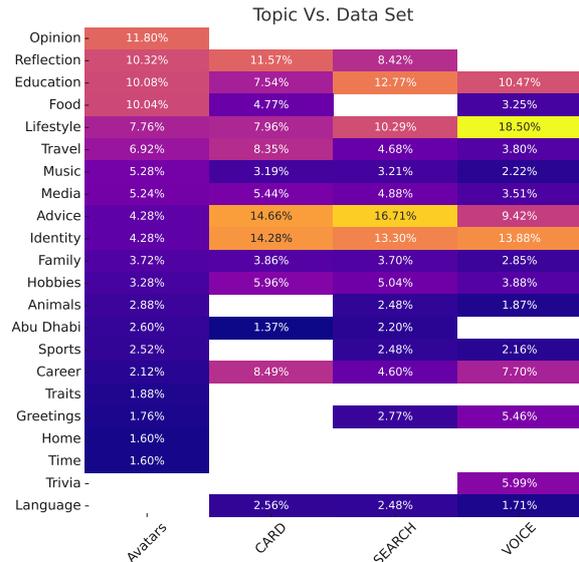


Figure 3: Heatmap of Topic Groups vs. Corpus Subset: The heatmap visualizes the distribution of questions across various topic groups ("Topic") and a subset of the HelloThere Corpus ("Data Set")—Avatars (the KBs of the recorded avatars), and (dialogue interactions by) CARD, SEARCH, and VOICE. The color intensity represents the proportion of questions, with brighter shades indicating higher proportions. Topics are ordered by higher coverage in the avatars’ KBs.

#### 4.5 What do People Ask?

We carried out topic clustering by leveraging the embeddings generated from GPT-3.5 Turbo. Specifically, we utilized the k-means clustering algorithm to group similar questions and tune the number of clusters until we identified recurring themes and could group them together sensibly. While we acknowledge this is a subjective labeling process, the clustering helped identify common themes across the avatars’ KBs and the dialogues, providing insights into the types of questions present in the corpus. We describe the topics in Table 5 and map their occurrence in the corpus in Figure 3. The heatmap visualization allows us to identify and quantify the prevalence of different topic clusters across the corpus subsets. The color intensity represents the proportion of questions in each topic-subset combination, offering an intuitive view of user interests and avatar knowledge distribution. This visualization helps identify potential gaps in avatar knowledge bases (*Avatars* on the X-Axis) and areas of high user engagement, informing future improvements in TOIA system design.

The heatmap presents several key observations

about how different topics fare across the HelloThere Corpus subsets. For instance, ‘Identity’ and ‘Advice’ are standout topics in the dialogues. The ‘Lifestyle’ topic is the most common in the VOICE channel, suggesting a focus on personal and day-to-day queries in voice-based (free-form) interactions. Interestingly, ‘Education’ and ‘Reflection’ topics are pretty evenly distributed across all modalities but VOICE and the avatars’ KBs, signifying their universal appeal to users. Contrarily, the localized topic of ‘Abu Dhabi’ seems less prevalent than in previous sub-sets. Some topics, such as ‘Home’ and ‘Time,’ lag in user engagement across all sets. Furthermore, a newly added ‘Trivia’ category shows particular traction in the VOICE channel, hinting at various questions that don’t necessarily slot into the existing categories. Lastly, it’s worth noting that there are visible data gaps in topics like ‘Opinion’ and ‘Traits,’ which appear exclusively in the Avatars channel. This could signify a lack of user engagement for these topics in the dialogues.

## 5 Conclusion and Further Work

In this paper, we presented the HelloThere corpus, which includes two main categories of datasets: 26 single-turn knowledge bases and multi-turn dialogue corpora featuring annotated chat logs. To ensure consistency, we have standardized our terminology throughout, using “q-a pairs” to refer to question-answer pairs in the knowledge bases and dialogues. All q-a pairs are rated by Human interactors and benchmarked for answer retrieval.

The HelloThere Corpus offers a multifaceted resource for the SIGDial community. It is beneficial for benchmarking conversational agents, studying user behavior, and conducting multimodal analysis. It allows for focused studies on dialogue complexity, retrieval failures, and localized or general user interests, providing a comprehensive foundation for future research in natural language interactions.

The key future directions we plan to work on include: (a) expanding the corpus with more data to support diverse research applications; (b) refining models to enhance answer retrieval efficiency and engagement in multi-turn dialogues; and (c) providing and evaluating model performance under multilingual conditions.

## References

- Dana Abu Ali, Muaz Ahmad, Hayat Al Hassan, Paula Dozsa, Ming Hu, Jose Varias, and Nizar Habash. 2018. A bilingual interactive human avatar dialogue system. In *Proceedings of the 19th Annual SIGDial Meeting on Discourse and Dialogue*, pages 241–244.
- Ron Artstein, Anton Leuski, Heather Maio, Tomer Mor-Barak, Carla Gordon, and David Traum. 2015. How many utterances are needed to support time-offset interaction? In *The Twenty-Eighth International Flairs Conference*.
- Alberto Chierici and Nizar Habash. 2021. A view from the crowd: Evaluation challenges for time-offset interaction applications. In *Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 75–85.
- Alberto Chierici and Nizar Habash. 2023. Tell me more, tell me more: Ai-generated question suggestions for the creation of interactive video recordings. In *2023 32nd IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*, pages 1725–1730. IEEE.
- Alberto Chierici, Nizar Habash, and Margarita Bicec. 2020. The margarita dialogue corpus: A data set for time-offset interactions and unstructured dialogue systems. In *Proc. of Language Resources and Evaluation Conference*.
- Alberto Chierici, Tyece Kiana Fredorcia Hensley, Wahib Kamran, Kertu Koss, Armaan Agrawal, Erin Meekhof, Goffredo Puccetti, and Nizar Habash. 2021. A cloud-based user-centered time-offset interaction application. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 265–268.
- Alberto Maria Chierici. 2023. *Scalable, Human-Like Asynchronous Communication*. Ph.D. thesis, New York University Tandon School of Engineering.
- Jianfeng Gao, Michel Galley, Lihong Li, et al. 2019. Neural approaches to conversational ai. *Foundations and Trends® in Information Retrieval*, 13(2-3):127–298.
- 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*.
- Margaret Li, Jason Weston, and Stephen Roller. 2019. Acute-eval: Improved dialogue evaluation with optimized questions and multi-turn comparisons. *arXiv preprint arXiv:1909.03087*.
- Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. *arXiv preprint arXiv:1506.08909*.

- Yosi Mass, Boaz Carmeli, Haggai Roitman, and David Konopnicki. 2020. Unsupervised faq retrieval with question generation and bert. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 807–812.
- Erinc Merdivan, Deepika Singh, Sten Hanke, Johannes Kropf, Andreas Holzinger, and Matthieu Geist. 2020. Human annotated dialogues dataset for natural conversational agents. *Applied Sciences*, 10(3):762.
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. 2022. Text and code embeddings by contrastive pre-training. *arXiv preprint arXiv:2201.10005*.
- Zach Nussbaum, John X Morris, Brandon Duderstadt, and Andriy Mulyar. 2024. Nomic embed: Training a reproducible long context text embedder. *arXiv preprint arXiv:2402.01613*.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. OpenAI.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Siva Reddy, Danqi Chen, and Christopher D Manning. 2019. Coqa: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- 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*.
- David Traum, Andrew Jones, Kia Hays, Heather Maio, Oleg Alexander, Ron Artstein, Paul Debevec, Alesia Gainer, Kallirroi Georgila, Kathleen Haase, et al. 2015. New Dimensions in Testimony: Digitally preserving a Holocaust survivor’s interactive storytelling. In *Proceedings of the International Conference on Interactive Digital Storytelling*, pages 269–281.
- Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart van Merriënboer, Armand Joulin, and Tomas Mikolov. 2015. Towards ai-complete question answering: A set of prerequisite toy tasks. *arXiv preprint arXiv:1502.05698*.
- Yu Wu, Wei Wu, Chen Xing, Ming Zhou, and Zhoujun Li. 2016. Sequential matching network: A new architecture for multi-turn response selection in retrieval-based chatbots. *arXiv preprint arXiv:1612.01627*.
- Asaf Yehudai, Matan Vetzler, Yosi Mass, Koren Lazar, Doron Cohen, and Boaz Carmeli. 2023. Qaid: Question answering inspired few-shot intent detection. *arXiv preprint arXiv:2303.01593*.
- 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*.

# It Couldn't Help But Overhear: On the Limits of Modelling Meta-Communicative Grounding Acts with Supervised Learning

Brielen Madureira<sup>1</sup>

David Schlangen<sup>1,2</sup>

<sup>1</sup>Computational Linguistics, Department of Linguistics  
University of Potsdam, Germany

<sup>2</sup>German Research Center for Artificial Intelligence (DFKI), Berlin, Germany  
{madureiralasota,david.schlangen}@uni-potsdam.de

## Abstract

Active participation in a conversation is key to building common ground, since understanding is jointly tailored by producers and recipients. Overhearers are deprived of the privilege of performing grounding acts and can only conjecture about intended meanings. Still, data generation and annotation, modelling, training and evaluation of NLP dialogue models place reliance on the *overhearing paradigm*. How much of the underlying grounding processes are thereby forfeited? As we show, there is evidence pointing to the impossibility of properly modelling human meta-communicative acts with data-driven learning models. In this paper, we discuss this issue and provide a preliminary analysis on the variability of human decisions for requesting clarification. Most importantly, we wish to bring this topic back to the community's table, encouraging discussion on the consequences of having models designed to only “listen in”.

## 1 Is Grounding “Supervisable”?

“What are you looking at?” asked Bob. “Magpies are building a nest outside!” Alice replied. If you were Bob, how would you continue that conversation? He could for instance say “Awesome!” or “I saw that”. Whatever you say, it will probably differ from how he continued: “Building what?”. The decision to request clarification depends on mutual understanding, which is contingent on *e.g.* the current situation, the familiarity between interlocutors and the previous utterances. Or, more formally, it depends on the clarification potential of these utterances (Ginzburg, 2012) and how they are assimilated into their *common ground* (Clark, 1996).

The one-to-many property of dialogue continuations is well-known in NLP (Zhao et al., 2017; Yeh et al., 2021; Towle and Zhou, 2022; Liu et al., 2023). There is a combinatorial explosion of possibilities for any interaction (Bates and Ayuso, 1991; Dingemans and Enfield, 2023), and individual

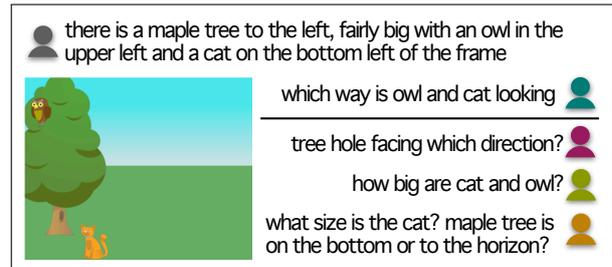


Figure 1: Variability of clarification requests produced by three overhearers in comparison to the original one, in an instance of the instruction-following CoDraw dialogue game (CC BY-NC 4.0), with cliparts from Zitnick and Parikh (2013).

human behaviour may vary at each point. This variability is hard to measure, since arguably no two people will ever be in the exact same situation with the same conversation history to react to (Yeomans et al., 2023).

Still, the prevailing end-to-end deep learning methods commonly rely on supervised learning (SL) from a sample of human behaviour instantiating the reaction of *a single human* at each observed context. Besides the issue of multiplicity of valid continuations, this paradigm faces another conceptual contention: dialogue models are trained to react upon a conversational history produced by someone else. In other words, they act as *overhearers*<sup>1</sup> of a dialogue in which they did not participate.

The suitability of data-driven methods and fixed corpora for modelling strategies and *conversational* grounding phenomena like Clarification Requests (CR) has been questioned (Schatzmann et al., 2005; Benotti and Blackburn, 2021b). Static datasets of human observations have empirically failed to provide enough information to define a human-like CR policy (Testoni and Fernández, 2024; Madureira and Schlangen, 2024). Moreover, chat-optimised

<sup>1</sup>We will use this term to also mean reading or seeing signs. Also called *observers* by Georgila et al. (2020).

LLMs mostly do not engage in grounding acts and, when they do, it does not fully align with human behaviour (Kuhn et al., 2022; Deng et al., 2023; Shaikh et al., 2023). The latter is not necessarily a problem: one can use other effective methods when it comes to building applications. But the first is: grounding is essential for human communication, and lack of it can lead to undesired breakdowns (Benotti and Blackburn, 2021a).

Since modelling human dialogue strategies and the use of meta-communicative acts remains an unsolved problem, we hereby wish to re-open the discussion on the consequences of overhearing, focusing on two grounding devices: backchannels and interactive repair (Fusaroli et al., 2017).

## 2 Overhearers in a Conversation

As Clark (1996) defined it, in addition to speakers and addressees,<sup>2</sup> a conversation can have *side-participants*, who are part of it but at a given moment are neither of the those two, and *overhearers*, who are spectators without any rights or responsibilities, e.g. a silent audience or a minute-taker who lacks the opportunity to interfere (Peters, 2010). They are further divided into *bystanders*, if one is aware of their presence, or *eavesdroppers*, who listen secretly (or at a later time). There is evidence that the very process of understanding differs between addressees and overhearers: while interlocutors actively construct mutual understanding with each other, overhearers only passively consume the product of that process (Schober and Clark, 1989).

Speakers can design their utterances while taking different attitudes towards overhearers when they are aware of their presence (Clark, 1992; Liu et al., 2016), but covert overhearers are not acknowledged at all in the conversation, and can only conjecture about the intended meanings (Clark, 1992). Although the grounding acts they witness, like backchannels, and the availability of multiple perspectives may indeed aid their comprehension (Tolins and Fox Tree, 2016; Tree and Mayer, 2008), the original interaction was opportunistically produced to be understood against the original participants' common ground (Schober and Clark, 1989).

In their corpus analysis of common ground in multi-party interactions, Eshghi and Healey (2007) showed evidence that overhearers reach lower levels of understanding than ratified side participants, who in their turn are not very different from di-

---

<sup>2</sup>Or producers and recipients.

rect addressees, in what they call *collective states of understanding*. Related to that, Georgila et al. (2020) showed that observers and participants perceive interactions differently and the experiments by Fox Tree (1999) provided evidence that overhearers can more easily comprehend instructions while listening to dialogues than to monologues. Clark (1992) even argued that most psycholinguistic subjects are actually overhearers, so theories of language processing may actually be theories of overhearing, due to their lack of interactivity.

Separating addressees from side participants and accommodating overhearers are salient problems in research on multi-party dialogue (Jovanovic and op den Akker, 2004; Ginzburg and Fernández, 2005; Traum et al., 2018; Parisse et al., 2022; Ganesh et al., 2023).

## 3 Are NLP Models Only Listening In?

More than a decade ago, Rieser and Lemon (2011) already discussed the limitations of using supervised approaches for learning dialogue strategies. They flagged up three concerns: textual data does not contain the underlying uncertainty measures, instances are treated as local point-wise estimates (instead of the sequences they really are) and exploration of novel strategies is not possible, since the model has access only to the outcomes of the chosen dialogue trajectory originally perpetrated by the humans. This reflects the (offline) *overhearing paradigm*: a person or agent interpreting a pre-existing conversation and deciding what to do if they were in the original participants' shoes.

In NLP, this paradigm is widely used in various modelling steps. Let us look closer at four main practices, which may have cascaded effects.

**Data Collection** Given the extra cost of coordinating the presence of more than one subject for generating dialogical data, especially in crowdsourcing campaigns, many strategies have been proposed to bypass that with overhearing. For instance, this happens when the data collection procedure is framed as a dialogue continuation task (Frommherz and Zarcone, 2021). To name a few related to grounding, we have Zhou et al. (2022) who extracted dialogue contexts from existing datasets and presented them in a two-stage approach for some workers to generate common ground inferences and, separately, others generated a continuation as a response. Variations of overhearing manifest in techniques to generate CRs or

their responses (Aliannejadi et al., 2021; Gao et al., 2022; Addlesee and Eshghi, 2024) and are even embedded in data collection tools that allow dialogues to be constructed without persistent workers (Cascante-Bonilla et al., 2019).

**Annotation and Analysis** Corpus studies of interactive linguistic use can only be performed from an overhearer perspective, without full evidence of what participants intended and understood or the reasons for their decisions (Brennan, 2000; Brennan et al., 2005). This is particularly challenging for research on common ground. For instance, Rodríguez and Schlangen (2004) and Schlöder and Fernández (2015) were confronted with the limitations of overhearers having only indirect access to the intentions of interlocutors when annotating CRs, partly remediating that by making a long dialogue context available. Niekrasz and Moore (2010) annotated references to conversation participants, joint actions that also serve to build common ground, emphasising that annotators were overhearers instructed to judge the speaker’s intended purpose. Other annotations of grounding acts and common ground states had to rely on overhearers (Markowska et al., 2023; Zhang et al., 2023; Mohapatra et al., 2024).

**Modelling** Prototypical data-driven models trained with supervised learning to *process* dialogue, and possibly continue it, can, by design, be regarded as overhearers. This fact was made clear, for instance, in the CR model by Schlangen (2004). Traum (2017) differentiated between the perspective of an observer in *dialogue modelling* and of a participant in *dialogue management*, stating that the main difference lies in the decision-making process of the latter, although some specific applications also exist for the first.

**Evaluation** In human evaluation, overhearer experiments (Whittaker and Walker, 2005) are very common, even though it limits the judgements and measurements to user’s *perceptions* of the dialogue (rather than the actual behaviour) (Whittaker and Walker, 2005; Foster and White, 2005; Moore, 2011) and restricts assessment of metrics like effectiveness and efficiency (Paksima et al., 2009). It has historically been a ubiquitous approach due to advantages like having control on one aspect of the evaluation while avoiding navigational and timing aspects of real interactions (Villalba et al., 2017), avoiding interference from ASR and other

technical problems (Buß et al., 2010), allowing the collection of feedback about alternative system responses (Walker et al., 2004) and avoiding natural language interpretation problems (Wärnestal et al., 2007). Demberg et al. (2011) contrasted text overhearers with speech overhearers, pointing out that reading dialogues is artificially simplified, since participants can go back to difficult portions and choose the pace, and the two setups may also impact how evaluators rate the system. The available context may also have to be adjusted (Spanger et al., 2010). Cercas Curry and Rieser (2019) explicitly addressed the limitations of evaluation by overhearing and advocated for interaction with users. For a recent overview of works that use similar forms of *static* evaluation, see (Finch and Choi, 2020).

As we have seen, the overhearing paradigm (fairly silently) permeates fundamental phases of dialogue modelling. The choice of this paradigm used to be a salient concept, with authors showing awareness of its limitations when it was employed. Kousidis and Schlangen (2015) even modelled a ratified side participant and had evaluators “overhear the overhearer”. In recent publications, however, it is often taken for granted, as if it was the only natural way to go. What can be the consequences when it comes to cognitive models of conversational grounding?

#### 4 Variability in Human Grounding Acts

As humans speak, they can provide positive and negative evidence of mutual understanding (Clark and Brennan, 1991; Roque and Traum, 2008), but modelling their timing and decision-making is challenging. Traum (2017) claimed that “it can be very difficult to efficiently capture regularities in behavioral patterns that lead to similar, but not identical structures”. In connection to that, people may take various paths in similar conversational situations (Bates and Ayuso, 1991). It is thus an open question how far data-driven supervised learning can get given the inherent variability of explicit (not to mention the latent) collateral signs of grounding.

Backchannels, a positive evidence of grounding, were demonstrated to involve individual variability, and even idiosyncrasy, possibly due to personality, gender or randomness (Huang and Gratch, 2012; Blomsma et al., 2024). Although those works showed some regularity in their timing, the SotA for the backchannel prediction task is not very high (.66 weighted F1) (Liermann et al., 2023).

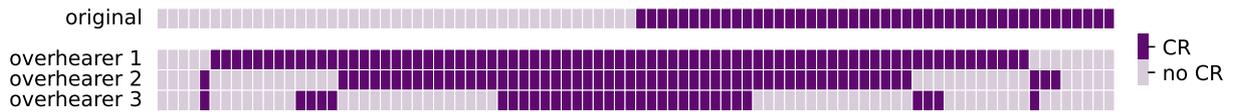


Figure 2: Variability in the decision of when to request clarification, comparing the decision of the original player with those of three overhearer annotators over 90 instances (horizontal axis) of the CoDraw game. Each cell is a data point and columns correspond to decisions on the same instance.

Findings on the variability of human decisions to initiate a CR, a negative sign of grounding, are still sparse. [Stoyanchev et al. \(2013\)](#) measured an absolute agreement of 39% among three annotators for *scripted* dialogues with missing ASR information. As another reference, [Shaikh et al. \(2023\)](#) reported a Cohen’s  $\kappa$  of 48.45 for clarification in emotional support conversations, which, they claimed, may even be inflated. The task of deciding when to request clarification in collaborative instruction following is under active investigation, but models’ performance is still suboptimal ([Shi et al., 2022](#); [Li et al., 2023](#); [Madureira and Schlangen, 2023](#); [Mohanty et al., 2023](#)). Recent works on the multi-modal CoDraw dialogue game ([Kim et al., 2019](#)) argued that this may be due to the variability in human decisions and the limitations of using supervised learning ([Testoni and Fernández, 2024](#); [Madureira and Schlangen, 2024](#)).

## 5 A Brief Analysis of Regularity in CRs

In CoDraw, an instruction follower receives instructions to reconstruct a scene using cliparts (as in Figure 1). Only the instruction giver sees the target scene. [Madureira and Schlangen \(2023\)](#) identified all CRs (around 11% of the instruction follower’s utterances) and defined the task of deciding when to request clarification, where models reached only up to .41 binary F1. What is missing as evidence for the claim that data-driven models cannot fully succeed in learning a “when policy” from human data is the actual human performance on this NLP task, i.e. what *overhearers* predict.

For an initial analysis, we collected a convenience sample with three annotators performing a similar task as the trained models: given a dialogue history and the current state of the reconstructed scene, decide which actions to take and, if needed, request clarification (details in Appendix). We randomly selected a sample with 90 instances; in half of them, the original player had produced a CR.

The average binary F1 of overhearers with respect to the original decision was .51, not much

above what SotA models achieve. But the proportion of CRs widely ranged from 36 to 85%. Among the three annotators, the Krippendorff’s  $\alpha$  was 0.10 and the mean pairwise Cohen’s  $\kappa$  was 0.18. That is already low, but if we consider the original decision as a fourth annotator, measures are even lower:  $\alpha$  was 0.02 and  $\kappa$  was 0.06. This indicates that there was slightly more agreement among overhearers than among addresses and overhearers, but in general there was little agreement on deciding when a CR should be realised. Figure 2 presents the main binary decision (whether to request clarification or not) for each of the 90 annotation instances, serving to provide a visual overview of such variability.

In terms of surface forms, the average BLEU score was 0.11 (std= 0.10) using the original CR as a source and the produced utterances as a reference. The mean cosine similarity between the embedding of the produced and the original CRs was 0.38, 0.29 and 0.36 for the three overhearers. Figure 1 shows an example of how diverse the produced clarifications can be, both in form and in content, even when all subjects made the same decision to clarify at a given point.

These are preliminary insights from a pilot study. Further standardised experiments with a larger sample must be conducted. Still, the results are already useful to strengthen the argument that, like backchannelling, human CR decisions lack regularity and overhearers have a much harder task trying to interpret and act upon someone else’s grounding acts. Decisions depend on how interlocutors distribute grounding costs, as per the principle of least collaborative effort ([Clark and Brennan, 1991](#)). Besides, there might be adaptive behaviours that models are not capturing ([Dideriksen et al., 2023](#)).

To continue this investigation, we propose distinguishing between the clarification potential ([Ginzburg, 2012](#); [Benotti, 2009](#)) and the clarification need. The first is a larger set of possibilities for clarification of a given utterance, while the latter refers to the decision of whether and what to clarify taken by a given individual operating with that ut-

terance and identifying something worth clarifying. Or, in other words, the clarification need, which is *in the agent*, refers to what was asked among all that could be asked. It is challenging to design experiments that can capture the clarification need among individuals, in particular due to the difficulty in replicating a given dialogue context for different subjects if they are not acting as overhearers. A possible next step is to turn the CR decision into an acceptability task, regarding it as a contrast. For each instance, the annotator would see a set of CRs. The actual CR observed in the data should ideally be accepted, but possibly others too. If the original CR falls into the empirical potential, there should be a plausible need for it at that point. Such experiment could also aim to measure uncertainty at each turn.

## 6 Discussion

Mutual understanding is crafted by “interacting minds” (Dingemanse et al., 2023). In dialogue, “interlocutors share or synchronise aspects of their private mental states and act together in the world” (Brennan et al., 2010). On the other hand, we have shown that the current NLP methodology mostly limits us to learning how *overhearers* predict discourse representations without the actual joint decision making facet, due to the way that data is produced and annotated, the assumptions behind training mechanisms and the evaluation protocols, each adding a layer of overhearing.

What can be a better setup to learn human dialogue behaviour, realising it as a truly interactive process? One needs to move on from one-off supervised learning to sequential models that not only *understand* dialogues but also *participate* in them.<sup>3</sup> Reinforcement learning provides that framing with a fully accessible and explorable environment (Rieser and Lemon, 2011), but somewhat circularly requires a good simulation of an user or interlocutor (Schatzmann et al., 2005; Georgila et al., 2006; Li et al., 2020). Although LLMs can serve as speaker simulators, so far they cannot fully model all dialogue phenomena. Another possibility are hybrid combinations of supervised and reinforcement learning (Henderson et al., 2008), as well as further improvements in techniques like RLHF,

<sup>3</sup>See (Min et al., 2022) for a related discussion on the limitations of imitation learning and behaviour cloning for embodied agents. See also (Ortega et al., 2021) for a discussion on supervised learning and the sequential aspect of an interaction.

PPO and DPO. But independently of the learning regime, data-driven approaches, which rely on extracting latent patterns and regularities in a corpus, stumble upon the individual variability of some dialogue phenomena, so that tasks may be ill-defined in datasets. Besides, although transcribed dialogue contain clues about the decision making during a conversation, they provide only limited evidence of what participants understood or intended, or their internal states (Brennan et al., 2005), which are pertinent for modelling some dialogue decisions and meta-communicative acts.

Indeed, interfaces do not necessarily have to conform to human behaviour, as long as they can sustain *graceful interaction* (Hayes, 1980). But from a cognitive perspective, the current NLP resources do not seem to satisfactorily meet our needs for modelling grounding mechanisms. To study the human mind, do we want cognitive models of how meaning and common ground are constructed or only of how they can be reverse engineered from someone else’s interactions?

**To conclude** With this argumentative paper, we wish to encourage more studies on the variability of human grounding acts and its impact in modelling human dialogue strategies. Besides, we advocate making the overhearing paradigm explicit whenever it is used in future publications and discussing how it can have influenced reported findings.

## Acknowledgements

We thank the anonymous reviewers for their valuable feedback. We are also thankful to the three annotators who took part on the pilot study.

## References

- Angus Addlesee and Arash Eshghi. 2024. [You have interrupted me again!: making voice assistants more dementia-friendly with incremental clarification](#). *Frontiers in Dementia*, 3:1343052.
- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2021. [Building and evaluating open-domain dialogue corpora with clarifying questions](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4473–4484, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Madeleine Bates and Damaris Ayuso. 1991. [A proposal for incremental dialogue evaluation](#). In *Speech and Natural Language: Proceedings of a Workshop Held at Pacific Grove, California, February 19-22, 1991*.

- Luciana Benotti. 2009. [Clarification potential of instructions](#). In *Proceedings of the SIGDIAL 2009 Conference*, pages 196–205, London, UK. Association for Computational Linguistics.
- Luciana Benotti and Patrick Blackburn. 2021a. [Grounding as a collaborative process](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 515–531, Online. Association for Computational Linguistics.
- Luciana Benotti and Patrick Blackburn. 2021b. [A recipe for annotating grounded clarifications](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4065–4077, Online. Association for Computational Linguistics.
- Peter Blomsma, Julija Vaitonyté, Gabriel Skantze, and Marc Swerts. 2024. [Backchannel behavior is idiosyncratic](#). *Language and Cognition*, page 1–24.
- Susan E. Brennan. 2000. [Invited talk: Processes that shape conversation and their implications for computational linguistics](#). In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, pages 1–11, Hong Kong. Association for Computational Linguistics.
- Susan E Brennan, Alexia Galati, and Anna K Kuhlen. 2010. [Two minds, one dialog: Coordinating speaking and understanding](#). In *Psychology of learning and motivation*, volume 53, pages 301–344. Elsevier.
- Susan E Brennan et al. 2005. How conversation is shaped by visual and spoken evidence. *Approaches to studying world-situated language use: Bridging the language-as-product and language-as-action traditions*, pages 95–129.
- Okko Buß, Timo Baumann, and David Schlangen. 2010. [Collaborating on utterances with a spoken dialogue system using an ISU-based approach to incremental dialogue management](#). In *Proceedings of the SIGDIAL 2010 Conference*, pages 233–236, Tokyo, Japan. Association for Computational Linguistics.
- Paola Cascante-Bonilla, Xuwang Yin, Vicente Ordonez, and Song Feng. 2019. [Chat-crowd: A dialog-based platform for visual layout composition](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 138–142, Minneapolis, Minnesota. Association for Computational Linguistics.
- Amanda Cercas Curry and Verena Rieser. 2019. [A crowd-based evaluation of abuse response strategies in conversational agents](#). In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 361–366, Stockholm, Sweden. Association for Computational Linguistics.
- Herbert H Clark. 1992. *Arenas of language use*. University of Chicago Press.
- Herbert H. Clark. 1996. *Common ground*, page 92–122. “Using” Linguistic Books. Cambridge University Press.
- Herbert H Clark and Susan E Brennan. 1991. [Grounding in communication](#). In *Perspectives on socially shared cognition.*, pages 127–149. American Psychological Association.
- Vera Demberg, Andi Winterboer, and Johanna D. Moore. 2011. [A strategy for information presentation in spoken dialog systems](#). *Computational Linguistics*, 37(3):489–539.
- Yang Deng, Lizi Liao, Liang Chen, Hongru Wang, Wenqiang Lei, and Tat-Seng Chua. 2023. [Prompting and evaluating large language models for proactive dialogues: Clarification, target-guided, and non-collaboration](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10602–10621, Singapore. Association for Computational Linguistics.
- Christina Dideriksen, Morten H Christiansen, Kristian Tylén, Mark Dingemanse, and Riccardo Fusaroli. 2023. [Quantifying the interplay of conversational devices in building mutual understanding](#). *Journal of Experimental Psychology: General*, 152(3):864.
- Mark Dingemanse and NJ Enfield. 2023. [Interactive repair and the foundations of language](#). *Trends in Cognitive Sciences*.
- Mark Dingemanse, Andreas Liesenfeld, Marlou Rasenberg, Saul Albert, Felix K Ameka, Abeba Birhane, Dimitris Bolis, Justine Cassell, Rebecca Clift, Elena Cuffari, et al. 2023. [Beyond single-mindedness: A figure-ground reversal for the cognitive sciences](#). *Cognitive science*, 47(1):e13230.
- Arash Eshghi and Patrick G.T. Healey. 2007. [Collective states of understanding](#). In *Proceedings of the 8th SIGdial Workshop on Discourse and Dialogue*, pages 2–9, Antwerp, Belgium. Association for Computational Linguistics.
- Sarah E. Finch and Jinho D. Choi. 2020. [Towards unified dialogue system evaluation: A comprehensive analysis of current evaluation protocols](#). In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 236–245, 1st virtual meeting. Association for Computational Linguistics.
- Mary Ellen Foster and Michael White. 2005. [Assessing the impact of adaptive generation in the comic multimodal dialogue system](#). In *Proceedings of the IJCAI 2005 Workshop on Knowledge and Reasoning in Practical Dialogue Systems*, pages 24–31.
- Jean E Fox Tree. 1999. [Listening in on monologues and dialogues](#). *Discourse processes*, 27(1):35–53.

- Yannick Frommherz and Alessandra Zarcone. 2021. Crowdsourcing ecologically-valid dialogue data for german. *Frontiers in computer science*, 3:686050.
- Riccardo Fusaroli, Kristian Tylén, Katrine Garly, Jakob Steensig, Morten H Christiansen, and Mark Dingemans. 2017. Measures and mechanisms of common ground: Backchannels, conversational repair, and interactive alignment in free and task-oriented social interactions. In *the 39th Annual Conference of the Cognitive Science Society (CogSci 2017)*, pages 2055–2060. Cognitive Science Society.
- Ananya Ganesh, Martha Palmer, and Katharina Kann. 2023. A survey of challenges and methods in the computational modeling of multi-party dialog. In *Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)*, pages 140–154, Toronto, Canada. Association for Computational Linguistics.
- Xiaofeng Gao, Qiaozi Gao, Ran Gong, Kaixiang Lin, Govind Thattai, and Gaurav S Sukhatme. 2022. Dalfred: Dialogue-enabled agents for embodied instruction following. *IEEE Robotics and Automation Letters*, 7(4):10049–10056.
- Kallirroi Georgila, Carla Gordon, Volodymyr Yanov, and David Traum. 2020. Predicting ratings of real dialogue participants from artificial data and ratings of human dialogue observers. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 726–734, Marseille, France. European Language Resources Association.
- Kallirroi Georgila, James Henderson, and Oliver Lemon. 2006. User simulation for spoken dialogue systems: learning and evaluation. In *Ninth International Conference on Spoken Language Processing*.
- Jonathan Ginzburg. 2012. *The interactive stance: Meaning for conversation*. Oxford University Press.
- Jonathan Ginzburg and Raquel Fernández. 2005. Scaling up from dialogue to multilogue: Some principles and benchmarks. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 231–238, Ann Arbor, Michigan. Association for Computational Linguistics.
- Phil Hayes. 1980. Expanding the horizons of natural language interfaces. In *18th Annual Meeting of the Association for Computational Linguistics*, pages 71–74, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- James Henderson, Oliver Lemon, and Kallirroi Georgila. 2008. Hybrid reinforcement/supervised learning of dialogue policies from fixed data sets. *Computational Linguistics*, 34(4):487–511.
- Lixing Huang and Jonathan Gratch. 2012. Crowdsourcing backchannel feedback: understanding the individual variability from the crowds. In *Feedback behaviors in dialog*.
- Natasa Jovanovic and Rieks op den Akker. 2004. Towards automatic addressee identification in multi-party dialogues. In *Proceedings of the 5th SIG-Dial Workshop on Discourse and Dialogue at HLT-NAACL 2004*, pages 89–92, Cambridge, Massachusetts, USA. Association for Computational Linguistics.
- Jin-Hwa Kim, Nikita Kitaev, Xinlei Chen, Marcus Rohrbach, Byoung-Tak Zhang, Yuandong Tian, Dhruv Batra, and Devi Parikh. 2019. CoDraw: Collaborative drawing as a testbed for grounded goal-driven communication. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6495–6513, Florence, Italy. Association for Computational Linguistics.
- Spyridon Kousidis and David Schlangen. 2015. The power of a glance: Evaluating embodiment and turn-tracking strategies of an active robotic overhearer. In *Proceedings of AAAI Spring Symposium on Turn-taking and Coordination in Human-Machine Interaction*.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2022. Clam: Selective clarification for ambiguous questions with generative language models. *arXiv preprint arXiv:2212.07769*.
- Haau-Sing (Xiaocheng) Li, Mohsen Mesgar, André Martins, and Iryna Gurevych. 2023. Python code generation by asking clarification questions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14287–14306, Toronto, Canada. Association for Computational Linguistics.
- Ziming Li, Julia Kiseleva, and Maarten de Rijke. 2020. Rethinking supervised learning and reinforcement learning in task-oriented dialogue systems. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3537–3546, Online. Association for Computational Linguistics.
- Wencke Liermann, Yo-Han Park, Yong-Seok Choi, and Kong Lee. 2023. Dialogue act-aided backchannel prediction using multi-task learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15073–15079, Singapore. Association for Computational Linguistics.
- Kris Liu, Jean Fox Tree, and Marilyn Walker. 2016. Coordinating communication in the wild: The art-walk dialogue corpus of pedestrian navigation and mobile referential communication. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 3159–3166, Portorož, Slovenia. European Language Resources Association (ELRA).
- Yongkang Liu, Shi Feng, Daling Wang, Yifei Zhang, and Hinrich Schütze. 2023. PVGRU: Generating diverse and relevant dialogue responses via pseudo-variational mechanism. In *Proceedings of the 61st Annual Meeting of the Association for Computational*

- Linguistics (Volume 1: Long Papers)*, pages 3295–3310, Toronto, Canada. Association for Computational Linguistics.
- Brielen Madureira and David Schlangen. 2023. [Instruction clarification requests in multimodal collaborative dialogue games: Tasks, and an analysis of the Co-Draw dataset](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2303–2319, Dubrovnik, Croatia. Association for Computational Linguistics.
- Brielen Madureira and David Schlangen. 2024. [Taking action towards graceful interaction: The effects of performing actions on modelling policies for instruction clarification requests](#). In *Proceedings of the Third Workshop on Understanding Implicit and Underspecified Language*, pages 1–21, Malta. Association for Computational Linguistics.
- Magdalena Markowska, Mohammad Taghizadeh, Adil Soubki, Seyed Mirroshandel, and Owen Rambow. 2023. [Finding common ground: Annotating and predicting common ground in spoken conversations](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8221–8233, Singapore. Association for Computational Linguistics.
- So Yeon Min, Hao Zhu, Ruslan Salakhutdinov, and Yonatan Bisk. 2022. [Don’t copy the teacher: Data and model challenges in embodied dialogue](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9361–9368, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Shrestha Mohanty, Negar Arabzadeh, Julia Kiseleva, Artem Zholus, Milagro Teruel, Ahmed Awadallah, Yuxuan Sun, Kavya Srinet, and Arthur Szlam. 2023. [Transforming human-centered ai collaboration: Redefining embodied agents capabilities through interactive grounded language instructions](#). *arXiv preprint arXiv:2305.10783*.
- Biswesh Mohapatra, Seemab Hassan, Laurent Romary, and Justine Cassell. 2024. [Conversational grounding: Annotation and analysis of grounding acts and grounding units](#). In *Proceedings of LREC-COLING 2024*.
- Johanna D. Moore. 2011. [Language generation for spoken dialogue systems \[invited talk\]](#). In *Proceedings of the 13th European Workshop on Natural Language Generation*, page 132, Nancy, France. Association for Computational Linguistics.
- John Niekrasz and Johanna D. Moore. 2010. [Annotating participant reference in English spoken conversation](#). In *Proceedings of the Fourth Linguistic Annotation Workshop*, pages 256–264, Uppsala, Sweden. Association for Computational Linguistics.
- Pedro A Ortega, Markus Kunesch, Grégoire Delétang, Tim Genewein, Jordi Grau-Moya, Joel Veness, Jonas Buchli, Jonas Degraeve, Bilal Piot, Julien Perolat, et al. 2021. [Shaking the foundations: delusions in sequence models for interaction and control](#). *arXiv preprint arXiv:2110.10819*.
- Taghi Paksima, Kallirroi Georgila, and Johanna Moore. 2009. [Evaluating the effectiveness of information presentation in a full end-to-end dialogue system](#). In *Proceedings of the SIGDIAL 2009 Conference*, pages 1–10, London, UK. Association for Computational Linguistics.
- Christophe Parisse, Marion Blondel, Stéphanie Caët, Claire Danet, Coralie Vincent, and Aliyah Morgenstern. 2022. [Multidimensional coding of multimodal languaging in multi-party settings](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2781–2787, Marseille, France. European Language Resources Association.
- Stanley Peters. 2010. [Listening in](#). In *Proceedings of the 24th Pacific Asia Conference on Language, Information and Computation*, pages 31–31, Tohoku University, Sendai, Japan. Institute of Digital Enhancement of Cognitive Processing, Waseda University.
- 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 and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Verena Rieser and Oliver Lemon. 2011. *Reinforcement learning for adaptive dialogue systems: a data-driven methodology for dialogue management and natural language generation*. Springer Science & Business Media.
- Kepa Joseba Rodríguez and David Schlangen. 2004. [Form, intonation and function of clarification requests in german task-oriented spoken dialogues](#). In *Proceedings of Catalog (the 8th workshop on the semantics and pragmatics of dialogue; SemDial04)*.
- Antonio Roque and David Traum. 2008. [Degrees of grounding based on evidence of understanding](#). In *Proceedings of the 9th SIGdial Workshop on Discourse and Dialogue*, pages 54–63, Columbus, Ohio. Association for Computational Linguistics.
- Jost Schatzmann, Kallirroi Georgila, and Steve Young. 2005. [Quantitative evaluation of user simulation techniques for spoken dialogue systems](#). In *Proceedings of the 6th SIGdial Workshop on Discourse and Dialogue*, pages 45–54, Lisbon, Portugal. Special Interest Group on Discourse and Dialogue (SIGdial).
- David Schlangen. 2004. [Causes and strategies for requesting clarification in dialogue](#). In *Proceedings of the 5th SIGdial Workshop on Discourse and Dialogue at HLT-NAACL 2004*, pages 136–143, Cambridge, Massachusetts, USA. Association for Computational Linguistics.

- Julian J. Schlöder and Raquel Fernández. 2015. [Clarifying intentions in dialogue: A corpus study](#). In *Proceedings of the 11th International Conference on Computational Semantics*, pages 46–51, London, UK. Association for Computational Linguistics.
- Michael F Schober and Herbert H Clark. 1989. [Understanding by addressees and overhearers](#). *Cognitive psychology*, 21(2):211–232.
- Omar Shaikh, Kristina Gligorić, Ashna Khetan, Matthias Gerstgrasser, Diyi Yang, and Dan Jurafsky. 2023. [Grounding or guesswork? large language models are presumptive grounders](#). *arXiv preprint arXiv:2311.09144*.
- Zhengxiang Shi, Yue Feng, and Aldo Lipani. 2022. [Learning to execute actions or ask clarification questions](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2060–2070, Seattle, United States. Association for Computational Linguistics.
- Philipp Spanger, Ryu Iida, Takenobu Tokunaga, Asuka Terai, and Naoko Kuriyama. 2010. [Towards an extrinsic evaluation of referring expressions in situated dialogs](#). In *Proceedings of the 6th International Natural Language Generation Conference*. Association for Computational Linguistics.
- Svetlana Stoyanchev, Alex Liu, and Julia Hirschberg. 2013. [Modelling human clarification strategies](#). In *Proceedings of the SIGDIAL 2013 Conference*, pages 137–141, Metz, France. Association for Computational Linguistics.
- Alberto Testoni and Raquel Fernández. 2024. [Asking the right question at the right time: Human and model uncertainty guidance to ask clarification questions](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 258–275, St. Julian’s, Malta. Association for Computational Linguistics.
- Jackson Tolins and Jean E Fox Tree. 2016. [Overhearers use addressee backchannels in dialog comprehension](#). *Cognitive science*, 40(6):1412–1434.
- Benjamin Towle and Ke Zhou. 2022. [Learn what is possible, then choose what is best: Disentangling one-to-many relations in language through text-based games](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4955–4965, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- David Traum. 2017. [Computational approaches to dialogue](#). *The Routledge Handbook of Language and Dialogue*, 1:143–161.
- David Traum, Cassidy Henry, Stephanie Lukin, Ron Artstein, Felix Gervits, Kimberly Pollard, Claire Bonial, Su Lei, Clare Voss, Matthew Marge, Cory Hayes, and Susan Hill. 2018. [Dialogue structure annotation for multi-floor interaction](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Jean E Fox Tree and Sarah A Mayer. 2008. [Overhearing single and multiple perspectives](#). *Discourse Processes*, 45(2):160–179.
- Martín Villalba, Christoph Teichmann, and Alexander Koller. 2017. [Generating contrastive referring expressions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 678–687, Vancouver, Canada. Association for Computational Linguistics.
- Marilyn A Walker, Stephen J Whittaker, Amanda Stent, Preetam Maloor, Johanna Moore, Michael Johnston, and Gunaranjan Vasireddy. 2004. [Generation and evaluation of user tailored responses in multimodal dialogue](#). *Cognitive Science*, 28(5):811–840.
- Pontus Wärnestal, Lars Degerstedt, and Arne Jönsson. 2007. [Emergent conversational recommendations: A dialogue behavior approach](#). In *Proceedings of the 8th SIGdial Workshop on Discourse and Dialogue*, pages 63–66, Antwerp, Belgium. Association for Computational Linguistics.
- Steve Whittaker and Marilyn Walker. 2005. [Evaluating dialogue strategies in multimodal dialogue systems](#). *Spoken Multimodal Human-Computer Dialogue in Mobile Environments*, pages 247–268.
- 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.
- Michael Yeomans, F Katelynn Boland, Hanne K Collins, Nicole Abi-Esber, and Alison Wood Brooks. 2023. [A practical guide to conversation research: How to study what people say to each other](#). *Advances in Methods and Practices in Psychological Science*, 6(4):25152459231183919.
- Xuanming Zhang, Rahul Divekar, Rutuja Ubale, and Zhou Yu. 2023. [GrounDialog: A dataset for repair and grounding in task-oriented spoken dialogues for language learning](#). In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 300–314, Toronto, Canada. Association for Computational Linguistics.
- Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017. [Learning discourse-level diversity for neural dialog models using conditional variational autoencoders](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 654–664, Vancouver, Canada. Association for Computational Linguistics.

Pei Zhou, Hyundong Cho, Pegah Jandaghi, Dong-Ho Lee, Bill Yuchen Lin, Jay Pujara, and Xiang Ren. 2022. [Reflect, not reflex: Inference-based common ground improves dialogue response quality](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10450–10468, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

C Lawrence Zitnick and Devi Parikh. 2013. Bringing semantics into focus using visual abstraction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3009–3016.

## A Appendix

**Annotation Task** The decisions from the over-hearer perspective were performed by 3 annotators. Two of them are student assistants employed in our lab and one is a volunteer acquainted with the first author. A simple GUI interface showed the dialogue history (from 1 to 3 turns) up to the last instruction giver instruction, the current state of the reconstructed scene and the gallery of available cliparts. They could select up to 4 high level, discrete actions (add, move, resize, flip, delete) and the corresponding cliparts from dropdown lists. Besides, they could type a clarification request to continue the dialogue if they wished (otherwise, the next utterance field should be left blank). In future studies, a full interface similar to the original game should be used, i.e. giving the opportunity for cliparts to be moved around and edited in the scene. Here, the selection of actions was just used to enforce that the overhearers reflected on the pertinent actions while deciding whether to request clarification. Note that the step of action taking makes annotators more privileged than plain overhearers that just process the dialogue, but it better approximates the decision of the iCR-Action-Taker models in [Madureira and Schlangen \(2024\)](#). In this case, they are overhearers of the dialogue context, but try to minimally act as a player doing the next step. The results work as an upper bound for plain overhearers.

**Additional Details** The inter-annotator agreement metrics were computed with `nltk` using `chencherry.method3` for smoothing. The sentence embeddings for the CR utterances were computed with model `sentence-transformers/all-MiniLM-L6-v2` from `SentenceTransformers` ([Reimers and Gurevych, 2019](#)).

# Data Augmentation Integrating Dialogue Flow and Style to Adapt Spoken Dialogue Systems to Low-Resource User Groups

Zhiyang Qi

The University of  
Electro-Communications  
1-5-1, Chofugaoka, Chofu,  
Tokyo, Japan  
qizhiyang@uec.ac.jp

Michimasa Inaba

The University of  
Electro-Communications  
1-5-1, Chofugaoka, Chofu,  
Tokyo, Japan  
m-inaba@uec.ac.jp

## Abstract

This study addresses the interaction challenges encountered by spoken dialogue systems (SDSs) when engaging with users who exhibit distinct conversational behaviors, particularly minors, in scenarios where data are scarce. We propose a novel data augmentation framework to enhance SDS performance for user groups with limited resources. Our approach leverages a large language model (LLM) to extract speaker styles and a pre-trained language model (PLM) to simulate dialogue act history. This method generates enriched and personalized dialogue data, facilitating improved interactions with unique user demographics. Extensive experiments validate the efficacy of our methodology, highlighting its potential to foster the development of more adaptive and inclusive dialogue systems.

## 1 Introduction

As an innovative technology at the forefront of artificial intelligence and speech processing, spoken dialogue systems (SDSs) have attracted significant interest from both academia and industry (Kawahara, 2018; Si et al., 2023; Abdul-Kader and Woods, 2015; Kim et al., 2021). Despite the powerful capabilities of large language models (LLMs), traditional SDS remain a focal point of research due to their superior control and interpretability (Singh et al., 2024). These systems are predominantly trained using data from human-to-human interactions, which highlight varying speaking styles, such as clarity of intentions, as depicted in Figure 1. This variability necessitates that human speakers adjust their dialogue strategies when engaging with different users. For instance, compared to adults, minors often exhibit less clarity in their intentions and give ambiguous responses, requiring more confirmatory language or additional inquiries to better adapt to the unique speaking styles of younger users. This adaptive approach is crucial for enhanc-



**Speaking Styles of User A**

- There may be a clear direction intention regarding destinations and activities.
- There can be a demand for more detailed information or specific proposals.

**Speaker :** [*DirectionQuestion*] So, do you have any plans for a trip? I mean, is there a particular place you want to visit?  
**User A :** Yes, I'd like to go to Hokkaido.  
**Speaker :** [*SeasonQuestion*] Oh, I see. Do you have a preferred season for your trip?  
**User A :** I'm thinking about going in the autumn.  
**Speaker :** [*PeopleQuestion*] How many people will be traveling with you?  
**User A :** Oh, it's just me.  
**Speaker :** [*RequestQuestion*] Okay, is there anything specific you want to do or see?  
**User A :** Yes. Well, I'd like to visit a place where I can see beautiful autumn leaves.  
**Speaker :** [*RequestConfirm, SearchInform*] Alright. Let me look up some places with beautiful autumn foliage for you.

**Speaking Styles of User B**

- There may be a lack of clear intention regarding destinations and activities.
- Responses to questions and suggestions can sometimes be limited.

**Speaker :** [*DirectionQuestion*] Now, have you decided on a specific place you'd like to visit?  
**User B :** Yes.  
**Speaker :** [*DirectionQuestion*] Where would you like to go?  
**User B :** Hiroshima.  
**Speaker :** [*RequestConfirm, DirectionQuestion*] Hiroshima, got it. Do you have a specific area in Hiroshima in mind?  
**User B :** No, I haven't decided yet.  
**Speaker :** [*RequestQuestion*] Okay. Is there anything in particular you'd like to do there? Any activities or food you want to try?  
**User B :** I want to visit my grandparents.  
**Speaker :** [*RequestConfirm*] Oh, your grandparents. So, they live around that area.

Figure 1: Real human-to-human conversations. Speakers adopt various dialogue acts when interacting with users employing diverse speaking styles.

ing the effectiveness and user-friendliness of SDS in real-world scenarios.

However, adapting SDSs to these distinctive speaking styles typically requires a wealth of annotated dialogue data, which can be challenging to

obtain in abundance owing to the minority status of users employing unique conversational strategies or behaviors. To enhance the performance of dialogue systems when engaging with data-scarce user groups exhibiting distinct speaking styles, targeted data augmentation is imperative, enabling the system to better cater to their needs.

This study introduces a tailored data augmentation framework designed specifically for low-resource user groups exhibiting distinctive conversational behaviors. Recognizing the unique conversational behaviors and challenges associated with minors and the inherent difficulty in obtaining their data (Aydin et al., 2021), our study conducts experiments utilizing dialogue data from minors to facilitate targeted data augmentation for this demographic.

As depicted in Figure 1, the unique speaking style of users directly influences the speaker’s dialogue acts (DAs) and indirectly shape response content. Therefore, our data augmentation framework focuses on the speaking styles of users and the trajectory of DAs.

Specifically, we utilized a LLM to extract the speaking styles of such users and speakers interacting with them. We then fine-tuned a pre-trained language model (PLM) using all available data in a low-resource setting to create varied histories of DAs for speakers interacting with these user groups. The resulting speaker styles and DA histories were input into the LLM to produce customized training dialogue data for these users. The primary goal is to enhance the model’s ability to predict DAs when interacting with low-resource groups with unique speaking styles, as controlling the content of generated responses through DAs is deemed effective (Kawano et al., 2021).

This study’s contributions are outlined below.

- We introduced a data augmentation method to enhance the performance of the DA prediction model when dealing with users who have limited data and unique conversational behaviors and styles.
- Through multiple experiments conducted in a low-resource setting, we have discovered that the difficulty of DA prediction varies across different users and demonstrated the adaptability and effectiveness of our proposed method.

## 2 Related Work

The scarcity of annotated data and the challenge of data imbalance are persistent issues in various artificial intelligence domains (Shorten and Khoshgoftaar, 2019; Shi et al., 2020; Ahmad et al., 2021; Hedderich et al., 2021; Kim et al., 2023). To address these effectively, data augmentation techniques have been employed, as demonstrated in prior research across different tasks (Feng et al., 2021; Bayer et al., 2022). For instance, Schick and Schütze (2021) generated text similarity datasets from scratch by instructing a large PLM. Similarly, Liu et al. (2022) and Chen and Yang (2021) enhanced data by manipulating individual utterances within dialogues—such as adding, deleting, changing their order, or regenerating them—while preserving the original meaning, which improved model performance in dialogue summarization tasks. While the abovementioned methods focus on generating individual sentences, our study aims to create coherent dialogues comprising multiple sentences tailored for specific target groups.

Mohapatra et al. (2021) utilized GPT-2 (Radford et al., 2019) to develop user and agent bots, generating comprehensive task-oriented dialogues through bot interactions, demonstrating notable enhancements in low-resource scenarios with datasets MultiWOZ (Budzianowski et al., 2018) and PersonaChat (Zhang et al., 2018). Recently, with the advanced text generation capabilities of LLMs, researchers have started using LLMs for data augmentation (Pan et al., 2023; Kim et al., 2023; Wang et al., 2023). For instance, Kim et al. (2023) guided LLMs to generate a broad spectrum of social dialogues using social commonsense knowledge from a knowledge graph. Pan et al. (2023) generated domain-specific, task-oriented dialogues by extracting dialogue paths from out-of-domain conversations. The concept of dialogue paths in their work aligns with the concept of DA history in our research. However, the key distinction is that while they extract DA paths from existing data, we generate tailored DA histories based on existing data, specifically optimized for target user groups.

## 3 The Proposed Framework

In this study, we aim to enhance the DA prediction performance of the system when dealing with low-resource user groups that exhibit unique dialogue strategies, by generating training data through the proposed data augmentation framework. In the

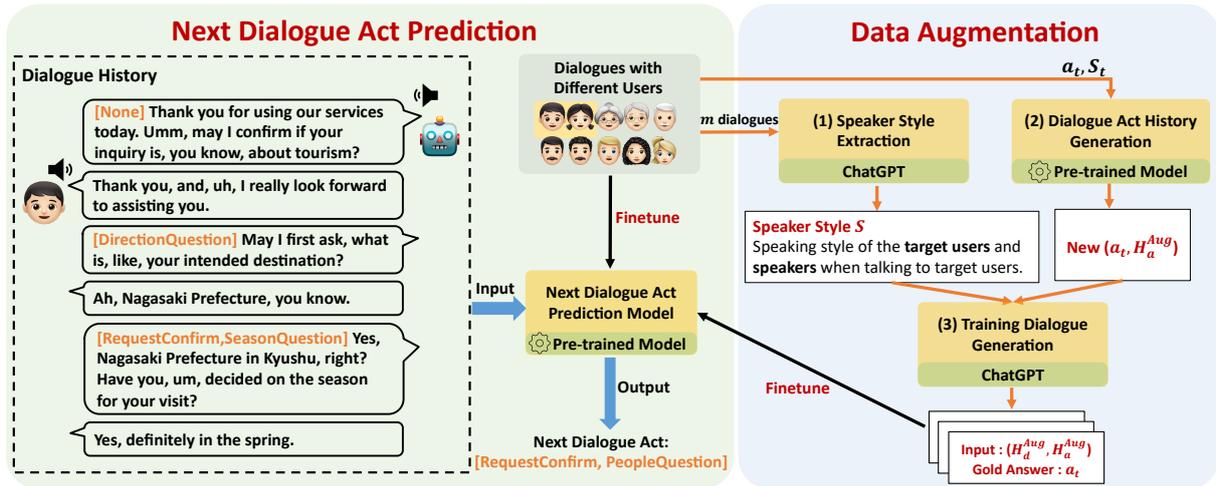


Figure 2: Our data augmentation framework is designed to improve the performance of the PLM in predicting DA when interacting with low-resource users who exhibit unique speaking styles. Beginning with dialogues that involve specific target users, we: (1) extract speaker styles, (2) generate DA histories of system interactions with these users, and (3) input this information into ChatGPT for tailored data augmentation.

construction of SDSs, accurate DA prediction is crucial as it facilitates dialogue state tracking and guides response generation, thereby reducing erroneous responses (Chen et al., 2017). The task depicted in the left portion of Figure 2 is defined as follows. Assuming the current turn of the dialogue is turn  $t$ , we utilize the dialogue history  $H_d = (S_{t-n}, U_{t-n}, \dots, S_{t-1}, U_{t-1})$  from the previous  $n$  turns, along with the system’s DA history  $H_a = (a_{t-n}, \dots, a_{t-1})$  from these turns, as the input. The output is the system’s DA  $a_t$  for the current turn.

Since we predict the current turn’s DA based on the dialogue history and the system’s DA history, it becomes crucial to generate dialogue and system DA histories that closely align with the target user group. To achieve this, we control the generation of dialogue data by capturing the speaking style of dialogue participants and generating dialogue flows that mimic real human interactions with the target user group. The importance of this approach lies in the fact that the model can effectively understand and adapt to unique dialogue strategies only when the training data realistically simulates complex dialogue scenarios. In real human interactions, users with unique dialogue strategies are in the minority and exhibit considerable diversity. Due to the limitations in data scale, traditional training datasets often fail to cover this diversity, which limits the model’s adaptability and accuracy when dealing with such users. By simulating the dialogue styles and processes of specific user groups, we can gener-

ate more diverse and precise training data, thereby enhancing the model’s generalizability and adaptability to diverse users.

As illustrated in Figure 2, our data augmentation framework comprises three components: (1) employing ChatGPT<sup>1</sup> to extract the speaker’s styles  $S$ , (2) finetuning a pre-trained model to generate the system’s DA history  $H_a^{Aug} = (a_{t-n}^{Aug}, \dots, a_{t-1}^{Aug})$ , and (3) inputting the extracted speaking styles  $S$  and the generated system’s DA history  $H_a^{Aug}$  into ChatGPT to generate the training dialogue data  $H_d^{Aug} = (S_{t-n}^{Aug}, U_{t-n}^{Aug}, \dots, S_{t-1}^{Aug}, U_{t-1}^{Aug})$ .

### 3.1 Speaker Styles Extraction

Since the unique speaking styles employed by the target user group significantly influence the content of conversations, it’s crucial to capture the speaking styles of this group by comparing dialogues from the target user group with those from non-target groups. This helps guide the subsequent generation of dialogues specifically tailored to the target user group. To facilitate this, we employ ChatGPT to extract speaker styles from conversations involving target users.

Specifically, we input a set of  $m$  dialogues, half of which involve users from the target group and the other half from non-target user groups. This balanced approach allows for an effective comparison, helping to identify and differentiate prominent speaking characteristics unique to the target group. Subsequently, ChatGPT is utilized to generate out-

<sup>1</sup><https://openai.com/blog/chatgpt>

puts representing the speaking style of the target users, as well as the speaking style of speakers when engaging with the target user group. Notably, our primary focus is on extracting abstract styles, such as "target users often exhibit ambiguous intentions towards destinations and activities." These styles are crucial because they significantly influence the direction of the dialogue, thereby enhancing the realism and relevance of the generated dialogues to actual human conversations. The prompt and extracted speaker styles are presented in Appendix C.

### 3.2 DA History Generation

As depicted in Figure 1, the unique conversational strategies employed by the target group also significantly influence the DAs of those engaging with them. Our objective at this stage is to generate a diverse and realistic DA history  $H_a^{Aug}$  that is specifically optimized for groups with distinctive speaking strategies. As shown in Figure 3, we achieve this by finetuning a PLM using existing data to generate the system’s DA history  $H_a^{Aug}$  for the previous  $n$  turns.

In particular, we utilize the DA  $a_t$  and utterance  $S_t$  from the current turn  $t$  as inputs, with the DA history  $H_a$  from the previous  $n$  turns as the desired output to establish training data. These data are then divided into two sets: one for training and the other for generation. Initially, we finetune the PLM using all available training data to capture DA histories that closely resemble real human conversations. Subsequently, we conduct a secondary finetuning utilizing training data exclusively from the target user group. This dual finetuning approach ensures that the model can generate DA histories that closely mimic real human dialogues and align with the unique speaking strategies of the target users. The first finetuning, which employs a relatively large dataset, enables the model to produce DA histories that mirror authentic human interactions. The second finetuning, focused on a smaller dataset specific to the target user group, allows the model to better tailor the DA histories to their unique characteristics.

During the generation phase, we input the the DA  $a_t$  and utterance  $S_t$  from the current turn  $t$  and generate the DA history  $H_a^{Aug}$  from the previous  $n$  turns. To ensure diversity, we simultaneously generate multiple outputs, selecting only those  $(a_t, H_a^{Aug})$  combinations that have not been previously observed.

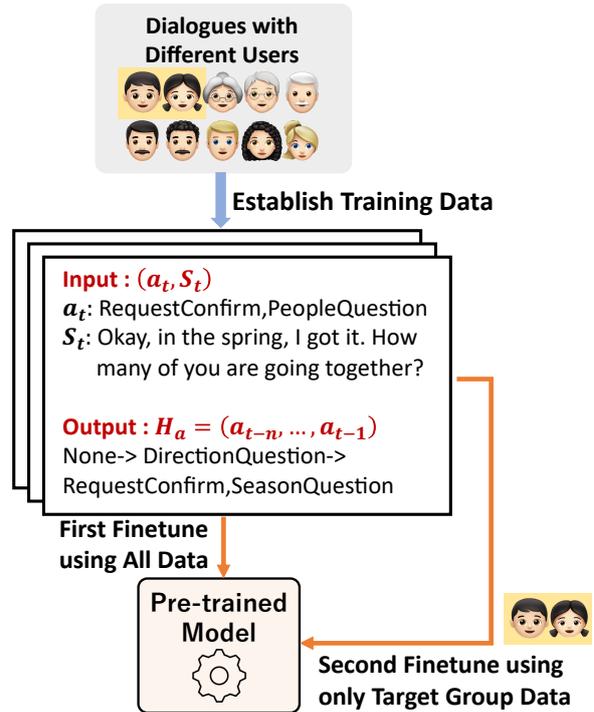


Figure 3: DA History Generation. We conduct two rounds of finetuning: the first round using all available data, and the second round using only data from the target user group, to ensure the generated DA history more closely aligns with the target demographic.

### 3.3 Dialogue Generation

Having obtained speaker styles and DA history tailored to users employing unique dialogue strategies, our ultimate goal is to generate dialogues corresponding to these styles and histories to enrich the training data for DA prediction. At this stage, we leverage ChatGPT’s powerful generation capabilities to create dialogue data for training purposes. Utilizing a few-shot prompt, we input the extracted speaking styles  $S$  and the DA histories  $H_a^{Aug}$  into ChatGPT to generate dialogues  $H_d^{Aug}$  that reflect the conversational style of the target users. Subsequently, we use the generated dialogues  $H_d^{Aug}$  and DA histories  $H_a^{Aug}$  as inputs, with  $a_t$  as the gold-standard answer, to construct the training data. The prompts used for generating these dialogues are detailed in Appendix D.

This approach aims to enhance the model’s ability to predict DAs when interacting with target users who exhibit unique conversational strategies. It effectively addresses the challenge of data scarcity by employing data augmentation.

## 4 Experiment

To evaluate the effectiveness of the proposed data augmentation framework, we conducted experiments using data from minors who employed unique conversational styles and strategies in actual dialogues within the dataset. These experiments were carried out in a low-resource setting across multiple splits, each utilizing different subsets of data from minors. We trained multiple DA prediction models on datasets of varying sizes, including models trained with augmented data added to the existing datasets.

### 4.1 Dataset

This study utilized a multimodal dialogue Japanese dataset known as the “Travel Agency Task Dialogue Corpus” (Inaba et al., 2022, 2024), which features conversations from users of various age groups, with detailed annotations of DAs. This dataset contains 115 hours of dialogue, spanning 330 conversations, with each averaging about 20 minutes. The dialogues were facilitated via Zoom video calls, involving six operators and 55 customers, including 20 minors (ages 7-17), 25 adults (ages 20-60), and 10 seniors (ages 65-72). Each customer participated in six dialogues.

The dialogues revolve around recommending travel destinations to users across various age groups. The dataset authors employed a Hidden Markov Model (HMM) (Rabiner, 1989) to analyze the transitions in dialogue among different age groups using sequences of DAs. A notable observation was that minors often used unique dialogue strategies compared to other age groups, typically expressing fewer independent opinions. The annotation of DAs was performed by functional segment, a unit smaller than an utterance. Each operator’s segment is annotated as one of the 28 predefined DAs related to travel destination recommendations, or as “None”. Examples of these DAs include asking about the travel season (Season-Question) and summarizing the travel plan (Travel-Summary), all of which are detailed in Appendix A. Since segments labeled “None” primarily consist of non-informative responses such as “Yeah” or “Uh-huh,” and our objective is to guide the system to generate accurate and meaningful responses using DA tags, we selectively included only those training instances where the gold-standard responses were not labeled “None” in this study. Additionally, we employed text-based human transcriptions

rather than audio recordings for our research.

### 4.2 Low-Resource Setting

We trained five DA prediction models using datasets of varying scales: Minors-Only, Zero-Shot, Low-Resource, Full-Resource, and Low-Resource+Augmentation(Ours). To simulate low-resource conditions for specific user demographics, we used dialogue data from only 3 minors out of a group of 20, totaling 18 dialogues for training. For evaluation, we used 60 dialogues from 10 minors.

- **Minors-Only:** Employed only 18 dialogues from 3 minors.
- **Zero-Shot:** Utilized all data from adults and seniors, amounting to 210 dialogues.
- **Low-Resource:** Combined the 18 dialogues from the Minors-Only with all 210 dialogues from adults and seniors, totaling 228 dialogues.
- **Full-Resource:** Included dialogues from 10 minors (60 dialogues), encompassing those from the 3 minors in the low-resource setting, plus all 210 dialogues from adults and seniors, totaling 270 dialogues.
- **Low-Resource + Aug(mentation) (Ours):** Used the 228 dialogues from the Low-Resource and supplemented them using our proposed augmentation framework. Additional data was generated until the dataset size matched that of the Full-Resource for a direct comparison.

### 4.3 Setup and Details

In the process of extracting speaker styles, we fed  $m = 6$  dialogues into GPT-4-0125-preview, where three were from minors in a low-resource setting, and the other three involved different adults or seniors. For generating training dialogues, GPT-3.5-turbo-0125 was employed.

During the DA history generation phase, we utilized Japanese T5-Large<sup>2</sup> as the PLM. We conducted two rounds of finetuning to ensure the model is capable of generating DA histories that not only closely mimic real human conversations but also align with the unique conversational strategies of minors during interactions. During the first training phase, the learning rate was set at 1e-4, and

<sup>2</sup><https://huggingface.co/retrieva-jp/t5-large-long>

Table 1: Training data quantity for DA prediction across four splits: MO (Minors-Only), ZS (Zero-Shot), LR (Low-Resource), FR (Full-Resource)

Split	Valid	MO-Valid	MO	ZS	LR	FR	Ours	Test
1	2,027	307	1,662	21,011	22,980	26,375	26,375	6,004
2	2,027	199	1,117	21,011	22,327	26,434	26,434	5,945
3	2,027	262	1,578	21,011	22,851	26,712	26,712	5,667
4	2,027	271	1,574	21,011	22,856	26,961	26,961	5,418

for the subsequent phase exclusively involving data from minors, it was set at  $5e-5$ . We utilized 210 adult and elderly conversations for generating DA histories, dividing them into 120 for training and 90 for generation purposes. To ensure data diversity and novelty, we retained only those  $(a_t, H_a^{Aug})$  combinations that had not previously existed; all 18 dialogues from 3 minors were included in both training and generation phases. To ensure diversity, we set the `num_return_sequences=3` when generating DA histories, meaning that for each data point, three DA histories are generated simultaneously.

In the DA prediction phase, Japanese T5-base<sup>3</sup> and Japanese GPT-NeoX<sup>4</sup> were used as the PLMs to validate the effectiveness of the generated data. We reconstructed the training and evaluation sets for the same DA prediction task to optimize hyperparameters, with specific details provided in Appendix B. Regarding the distribution of training and validation sets, the validation sets for all settings, except Minors-Only, are identical, comprising 21 dialogues from adults and seniors. The Minors-Only validation set consists of 3 dialogues from minors in the low-resource scenario. To validate the generalizability of our method, we conducted experiments across four splits, each using data from three different minors for training under a low-resource setting, while also varying the test data. Details on the data points for each split, after removing entries with a gold-standard answer of "None," are outlined in Table 1.

Considering that a single utterance may consist of multiple segments (see Figure 1 and Figure 2), each potentially be labeled with a different DA, there may be more than one gold-standard DA label for the current turn. Therefore, we employed both **exact match and partial match rates** as evaluation metrics. The exact match rate is a strict metric requiring the predicted set of labels to completely align with the true set of gold labels, measuring the model’s ability to fully grasp the dialogue con-

text and predict all relevant DA labels accurately. The partial match rate assesses the model’s performance in predicting some correct labels. This metric is more lenient, recognizing that in real conversations, capturing the main intent or action of the dialogue, even if not every label is precisely predicted, is still valuable. Therefore, the partial match rate helps understand the model’s robustness in practical use. Combined, these two metrics offer a balanced approach to evaluating the model’s DA prediction capabilities, providing a more accurate reflection of the model’s performance.

## 5 Results and Analysis

Table 2 shows the mean and standard deviation after five runs using seeds ranging from 1 to 5 across four different splits. While the **Minors-Only** solely comprised data from minors, its performance was inferior to the **Zero-Shot** model trained only with adult and elderly dialogue data due to the limited amount of training data. Therefore, we also used all available adult and elderly dialogue data in other setups to enhance the model’s generalization capabilities.

Additionally, since **Zero-Shot** does not use minor’s dialogues, the training data remains consistent across the four different splits. The variation in **Zero-Shot**’s performance across the splits further underscores the differences in the model’s adaptability to different minors, with the third split proving most challenging.

Across the four splits, the performance of our proposed data augmentation framework, **Low-Resource + Aug (Ours)**, almost all surpassed that of **Low-Resource** on both T5 and GPT-NeoX in terms of mean exact and partial match rates. This demonstrates that even in a low-resource setting, our method successfully captures the characteristics of minor speakers and generates dialogue flows that align with minor speaking behaviors, thereby guiding the generation of training dialogues.

However, even though we augmented the data to match the quantity of the **Full-Resource** in each split, **Full-Resource** typically showed superior per-

<sup>3</sup><https://huggingface.co/retrieval-ja/t5-base-long>

<sup>4</sup><https://huggingface.co/stockmark/gpt-neox-japanese-1.4b>

Table 2: Results across four different splits.

Split	Setting	Japanese GPT-NeoX		Japanese T5-base	
		Exact Match	Partial Match	Exact Match	Partial Match
1	Minors-Only	0.2451 ± 0.0117	0.3447 ± 0.0131	0.2533 ± 0.0083	0.3519 ± 0.0090
	Zero-Shot	0.2966 ± 0.0071	0.4049 ± 0.0092	0.3000 ± 0.0059	0.4066 ± 0.0053
	Low-Resource	0.3041 ± 0.0070	0.4228 ± 0.0073	0.3085 ± 0.0065	0.4232 ± 0.0064
	Low-Resource + Aug (Ours)	<b>0.3137 ± 0.0064</b>	<b>0.4320 ± 0.0094</b>	<b>0.3148 ± 0.0050</b>	<b>0.4244 ± 0.0056</b>
	Full-Resource	0.3190 ± 0.0074	0.4489 ± 0.0049	0.3125 ± 0.0029	0.4418 ± 0.0023
2	Minors-Only	0.2302 ± 0.0103	0.3677 ± 0.0105	0.2419 ± 0.0050	0.3311 ± 0.0079
	Zero-Shot	0.3162 ± 0.0069	0.4247 ± 0.0099	0.3200 ± 0.0039	0.4263 ± 0.0046
	Low-Resource	0.3220 ± 0.0071	0.4401 ± 0.0051	0.3257 ± 0.0019	0.4430 ± 0.0066
	Low-Resource + Aug (Ours)	<b>0.3290 ± 0.0083</b>	<b>0.4460 ± 0.0111</b>	<b>0.3270 ± 0.0029</b>	<b>0.4473 ± 0.0095</b>
	Full-Resource	0.3294 ± 0.0068	0.4526 ± 0.0074	0.3339 ± 0.0052	0.4486 ± 0.0075
3	Minors-Only	0.2329 ± 0.0033	0.3291 ± 0.0069	0.2528 ± 0.0038	0.3499 ± 0.0010
	Zero-Shot	0.2771 ± 0.0053	0.3878 ± 0.0075	0.2787 ± 0.0054	0.3889 ± 0.0054
	Low-Resource	0.2863 ± 0.0055	0.4070 ± 0.0019	0.2825 ± 0.0036	0.4010 ± 0.0156
	Low-Resource + Aug (Ours)	<b>0.2906 ± 0.0055</b>	<b>0.4077 ± 0.0067</b>	<b>0.2865 ± 0.0042</b>	<b>0.4097 ± 0.0090</b>
	Full-Resource	0.2889 ± 0.0069	0.4282 ± 0.0085	0.2986 ± 0.0058	0.4270 ± 0.0057
4	Minors-Only	0.2325 ± 0.0083	0.3336 ± 0.0093	0.2429 ± 0.0036	0.3480 ± 0.0091
	Zero-Shot	0.2900 ± 0.0066	0.4041 ± 0.0066	0.2947 ± 0.0047	0.4056 ± 0.0059
	Low-Resource	0.2925 ± 0.0067	0.4098 ± 0.0088	0.2983 ± 0.0031	<b>0.4156 ± 0.0120</b>
	Low-Resource + Aug (Ours)	<b>0.3005 ± 0.0069</b>	<b>0.4254 ± 0.0087</b>	<b>0.3000 ± 0.0056</b>	0.4144 ± 0.0096
	Full-Resource	0.3096 ± 0.0049	0.4425 ± 0.0098	0.3094 ± 0.0073	0.4336 ± 0.0019

formance. A possible explanation is the lack of quality control, which meant that subpar data was not filtered out, leading to poorer adaptation compared to **Full-Resource**, which used data exclusively from real human conversations. Additionally, the "Travel Agency Task Dialogue Corpus," derived from video calls and manually transcribed, may contain colloquial filler words and other informal elements in its complete utterances. In contrast, ChatGPT-generated dialogues tend to be more structured and fluid. This stylistic difference could also contribute to the observed performance disparity between **Low-Resource + Aug (Ours)** and **Full-Resource**.

### 5.1 Ablation

To evaluate the individual effectiveness of components in our proposed framework, we conducted ablation experiments using Japanese GPT-NeoX across four splits:

- **w/o DA History Gen:** In this model, we omitted the generation of new DA histories and instead randomly selected DA histories from the Low-Resource for data generation.
- **DA History Gen w/o Second Finetune:** This variant involved finetuning the DA history generation model only once, without a second round of finetuning tailored specifically for minors.
- **w/o Speaker Style:** This model utilized the

same DA histories as our complete method but did not use extracted speaker styles during dialogue data generation.

Table 3 shows the average results across the four splits, conducting five trainings for each model in every split with seed values set from 1 to 5. The findings indicate that both **w/o DA History Gen** and **w/o Speaker Style** variants achieved higher mean exact and partial match rates than the **Low-Resource**. This demonstrates that the training data generated through the independent use of style extraction and DA history generation components can also significantly improve performance.

Furthermore, although **DA History Gen w/o Second Finetune** did not use data from the target user group for a second fine-tuning during the training of the DA history generation model, its performance still surpassed that of **w/o DA History Gen**. This indicates that in generating DA history, even without a second finetuning to optimize the PLM for minors, the new  $(a_t, H_a^{Aug})$  combinations generated by a PLM trained with all available data can still enhance performance. Ultimately, **Ours** achieved the highest rates for both exact and partial matches, indicating that the combination of speaker styles extraction and DA history generation is most effective and underscores the necessity of targeted age-specific second finetuning when training the DA history generation model.

Table 3: Average results of the ablation experiments across four splits.

Setting	Exact Match	Partial Match
Low-Resource	0.3012	0.4199
w/o DA History Gen	0.3052	0.4263
DA History Gen w/o Second Finetune	0.3072	0.4269
w/o Speaker Style	0.3027	0.4274
Ours	<b>0.3085</b>	<b>0.4278</b>

### Speaker Style

#### Speaking style of the target users:

- Tends to lack a concrete image of the destination or activities, using vague expressions.
- More likely to have sensory-based desires rather than specific plans. For example, "want to enjoy nature" or "prefer a quiet place with few people". (...)

#### Speaking style of the interlocutors when conversing with target users:

- Make suggestions while checking the speaker's response.
- Ask questions to draw out the speaker's interests and desires. (...)

### Dialogue Act History

RequestQuestion -> RequestConfirm -> SearchConditionInform

#### w/o Speaker Style

**System**: 旅行先では何を楽しみにされていますか？(What are you looking forward to when traveling?)

**User**: あ、お寺や神社が見たいです。(Oh, I want to see temples and shrines.)

**System**: 了解しました。お寺や神社の観光地をお探しいたしますね。(Understood. I'll look for tourist spots with temples and shrines for you.)

**User**: よろしくお願ひします。(Thank you.)

**System**: では、京都の清水寺や金閣寺など、たくさんのお寺や神社があるエリアがおすすめです。どういたしまして。(In that case, I recommend areas like Kyoto, where there are many temples and shrines, such as Kiyomizu Temple and Kinkaku-ji. You're welcome.)

**User**: はい。(Yes.)

#### Ours

**System**: 旅行の際に特にこだわりや希望はありますか？(Do you have any particular preferences or desires for your trip?)

**User**: うーん、特にないです。(Hmm, not really.)

**System**: そうですね、何も特にこだわりがないということですね。(I see, no specific preferences then.)

**User**: はい。(Yes.)

**System**: その場合、近場で穏やかな雰囲気が楽しめる場所をおすすめします。どうでしょうか？(In that case, I recommend somewhere nearby with a calm atmosphere. How does that sound?)

**User**: いいですね。(That sounds nice.)

Figure 4: Dialogues generated by the variant without speaker styles and our approach.

## 5.2 Why did the Speaker Style work?

Figure 4 displays dialogues generated by **w/o Speaker Style** and **Ours**, using the same DA history. The DA history consists of first asking the user a travel-related request (RequestQuestion), then confirming the request (RequestConfirm), and finally indicating the content to be searched (SearchConditionInform). We observed that without the speaker style, the user in the **w/o Speaker Style** provided specific travel requirements, and the dialogue progressed smoothly. In contrast, the user in the **Ours** did not exhibit a clear intent. This indicates that the speaker style is effective, resulting in dialogues that more closely match the speaking styles of minors and aligning more closely with real human conversations.

## 5.3 Why did the DA History Generation work?

We compared the performance in generating DA histories between **DA History Gen w/o Second**

## Finetune and Ours on split 1.

For a direct comparison, we used 9,999 data points  $(a_t, S_t)$  from dialogues involving 90 adults and seniors to generate DA histories  $H_a^{Aug}$ , resulting in three DA histories per data point. This generation was conducted under the settings of top\_k=50, top\_p=0.9, and temperature=0.9. After removing duplicate  $(a_t, H_a^{Aug})$ , **DA History Gen w/o Second Finetune** produced 7,677 new  $(a_t, H_a^{Aug})$ , whereas **Ours** generated 10,412. We assessed how many of these combinations appeared in dialogues involving 17 minors (excluding those from the **Low-Resource**), finding 908 for **DA History Gen w/o Second Finetune** and 956 for **Ours**. Referencing Table 3, we can infer that compared to **w/o DA History Gen** which relied solely on existing DA histories, both **DA History Gen w/o Second Finetune** and **Ours** generated DAs that were present in the target user group, leading to improved performance. Notably, **Ours**, which underwent secondary finetuning for the target users,

produced more DA histories closely aligned with the target group, enhancing performance.

## 6 Conclusion

We introduced a data augmentation method designed to enhance the performance of the DA prediction model for users with limited data and unique conversational styles. Our experiments confirmed the reliability of the proposed method and the effectiveness of its components. While this study did not exhaustively explore the full potential for improvement of the proposed method, we plan to further evaluate this aspect in our future work.

## Acknowledgments

This work was supported by JSPS KAKENHI Grant Number 19H05692.

## References

- Sameera A. Abdul-Kader and Dr. John Woods. 2015. [Survey on chatbot design techniques in speech conversation systems](#). *International Journal of Advanced Computer Science and Applications*, 6(7).
- Tanveer Ahmad, Dongdong Zhang, Chao Huang, Hongcai Zhang, Ningyi Dai, Yonghua Song, and Huanxin Chen. 2021. [Artificial intelligence in sustainable energy industry: Status quo, challenges and opportunities](#). *Journal of Cleaner Production*, 289:125834.
- Selami Aydin, Leyla Harputlu, Özgehan Uştuk, Şeyda Savran Çelik, and Serhat Güzel. 2021. Difficulties in collecting data from children aged 7–12. *International Journal of Teacher Education and Professional Development (IJTEPD)*, 4(1):89–101.
- Markus Bayer, Marc-André Kaufhold, and Christian Reuter. 2022. [A survey on data augmentation for text classification](#). *ACM Comput. Surv.*, 55(7).
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017. [A survey on dialogue systems: Recent advances and new frontiers](#). *SIGKDD Explor. Newsl.*, 19(2):25–35.
- Jiaao Chen and Diyi Yang. 2021. [Simple conversational data augmentation for semi-supervised abstractive dialogue summarization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6605–6616, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. [A survey of data augmentation approaches for NLP](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, Online. Association for Computational Linguistics.
- Michael A. Hedderich, Lukas Lange, Heike Adel, Janik Strötgen, and Dietrich Klakow. 2021. [A survey on recent approaches for natural language processing in low-resource scenarios](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2545–2568, Online. Association for Computational Linguistics.
- Michimasa Inaba, Yuya Chiba, Ryuichiro Higashinaka, Kazunori Komatani, Yusuke Miyao, and Takayuki Nagai. 2022. [Collection and analysis of travel agency task dialogues with age-diverse speakers](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 5759–5767, Marseille, France. European Language Resources Association.
- Michimasa Inaba, Yuya Chiba, Zhiyang Qi, Ryuichiro Higashinaka, Kazunori Komatani, Yusuke Miyao, and Takayuki Nagai. 2024. [Travel agency task dialogue corpus: A multimodal dataset with age-diverse speakers](#). *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*
- Tatsuya Kawahara. 2018. Spoken dialogue system for a human-like conversational robot erica. In *International Workshop on Spoken Dialogue Systems Technology*.
- Seiya Kawano, Koichiro Yoshino, and Satoshi Nakamura. 2021. [Controlled neural response generation by given dialogue acts based on label-aware adversarial learning](#). *Transactions of the Japanese Society for Artificial Intelligence*, 36(4):E-KC9<sub>1</sub> – –14.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2023. [SODA: Million-scale dialogue distillation with social commonsense contextualization](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12930–12949, Singapore. Association for Computational Linguistics.
- Seokhwan Kim, Yang Liu, Di Jin, Alexandros Papangelis, Karthik Gopalakrishnan, Behnam Hedayatnia, and Dilek Z. Hakkani-Tür. 2021. How robust r u?: Evaluating task-oriented dialogue systems on spoken conversations. *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1147–1154.
- Yongtai Liu, Joshua Maynez, Gonçalo Simões, and Shashi Narayan. 2022. [Data augmentation for low-resource dialogue summarization](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages

- 703–710, Seattle, United States. Association for Computational Linguistics.
- Biswesh Mohapatra, Gaurav Pandey, Danish Contractor, and Sachindra Joshi. 2021. [Simulated chats for building dialog systems: Learning to generate conversations from instructions](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1190–1203, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yan Pan, Davide Cadamuro, and Georg Groh. 2023. Data-augmented task-oriented dialogue response generation with domain adaptation. In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, pages 96–106, Hong Kong, China. Association for Computational Linguistics.
- L.R. Rabiner. 1989. [A tutorial on hidden markov models and selected applications in speech recognition](#). *Proceedings of the IEEE*, 77(2):257–286.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Timo Schick and Hinrich Schütze. 2021. [Generating datasets with pretrained language models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6943–6951, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Wenzhong Shi, Min Zhang, Rui Zhang, Shanxiong Chen, and Zhao Zhan. 2020. [Change detection based on artificial intelligence: State-of-the-art and challenges](#). *Remote Sensing*, 12(10).
- Connor Shorten and Taghi M. Khoshgoftaar. 2019. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6:1–48.
- Shuzheng Si, Wentao Ma, Haoyu Gao, Yuchuan Wu, Ting-En Lin, Yinpei Dai, Hangyu Li, Rui Yan, Fei Huang, and Yongbin Li. 2023. [SpokenWOZ: A large-scale speech-text benchmark for spoken task-oriented dialogue agents](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Chandan Singh, Jeevana Priya Inala, Michel Galley, Rich Caruana, and Jianfeng Gao. 2024. [Rethinking interpretability in the era of large language models](#).
- Xi Wang, Hossein Rahmani, Jiqun Liu, and Emine Yilmaz. 2023. [Improving conversational recommendation systems via bias analysis and language-model-enhanced data augmentation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3609–3622, Singapore. Association for Computational Linguistics.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting*
- of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.

## A DA tags in Travel Agency Task Dialogue Corpus

In this study, we utilized the "Travel Agency Task Dialogue Corpus" collected by Inaba et al. (2024), which includes task specific DA annotations. The dataset defines DA tags for operators and customers in travel agency conversations, with 28 tags for operators and 8 tags for customers. In this study, only the operator’s tags were used, as shown in Table 4.

Table 4: Task Specific Dialogue Act Tags for Operator Segments.

Dialogue Act	Description	Example
DirectionQuestion	Question on areas for the desired travel	To which destination are you planning to travel?
SeasonQuestion	Question on the desired season	When will you go?
PeopleQuestion	Question about the number of people traveling and their relationships with the customer	How many people are traveling with you?
AgeQuestion	Question on the age of customers or their companions	How old are your children?
ExperienceQuestion	Question about the customer’s experience	Have you ever been to Osaka?
RequestQuestion	Question about the tourist spot request	What would you like to do there?
SearchAdvice	Questions or suggestions related to the tourist spot information retrieval system	Should I look for a restaurant there?
RequestConfirm	Confirmation of requests for tourist spots	You want to go to a Spa, don’t you?
DestinationConfirm	Confirmation of destination	Am I correct in Assuming that you are going to Yashi Park?
AddDestinationList	Addition to destination list by operator	I’ll add this location to the list.
TravelSummary	Summary of trip planning	Looking back, you plan to visit the Toshogu Shrine first.
SearchInform	Operator’s declaration of intent to search tourist spots in the system	I will now search.
PhotoInform	Provide information on photos displayed on the system	Here is a picture of a meal containing a lot of salmon roe.
SearchConditionInform	Provide information on search conditions	I can also filter by the time required.
NameInform	Provide information on the names of tourist spots	There is a commercial complex called the Sapporo Factory.
IntroductionInform	Provide information on tourist spots based on the system search results	It was established In 1876.
OfficeHoursInform	Provide information on hours of operation and closing dates	Our business hours span 10:00 a.m. to 10:00 p.m.
PriceInform	Provide information on fees and price range	The admission fee is 360 yen.
FeatureInform	Providing information about the characteristics of tourist spots	It is recommended for women even when it rains.
AccessInform	Provide information on access	This location is a five-minute walk from the railway station.
PhoneNumberInform	Provide information on telephone numbers	The phone number is 095 824.
ParkInform	Provide information on parking	There are three parking lots.
EmptyInform	Statement that there are no search results or specific description	I do not see anything in the search results.
MistakeInform	Correcting errors in tourist spot information	Sorry, this store is open on all days of the week.
OperatorSpotImpression	Subjective evaluations and assumptions about a tourist spot by operators	This restaurant looks nice and inexpensive.
SearchResultInform	Report overall search results	It appears there are numerous stores in this location.
OnScreenSuggest	Suggestions for tourist spots on the shared screen	How about this site?
OnScreenQuestion	Questions about tourist spots on the shared screen	Which one looks the best, number 1, 2, or 3?

## B Hyperparameter Optimization

During our experiments, we performed hyperparameter optimization.

For T5-base, we conducted a grid search with batch sizes of {8, 16, 32, 64}, warmup ratios of {0, 0.1, 0.2}, and learning rates of {3e-3, 2e-3, 1e-3, 9e-4, 8e-4}. The optimal configuration was identified as a batch size of 64, a warmup ratio of 0.1, and a learning rate of 1e-3.

Similarly, for GPT-NeoX, we conducted a grid search with batch sizes of {4, 8, 16}, warmup ratios of {0.1, 0.2, 0.3}, and a range of learning rates of {3e-4, 2e-4, 1e-4, 9e-5, 8e-5, 7e-5, 6e-5, 5e-5, 4e-5}. The best settings were determined to be a batch size of 8, a warmup ratio of 0.1, and a learning rate of 9e-5.

## C Details for Speaker Styles Extraction.

We utilized the prompt shown in Figure 5 to extract speaker styles using the GPT-4-0125-preview model, with six dialogues from different users, three from the target user group and three from a non-target user group. As the extraction was conducted with the default temperature setting (i.e., temperature=1), the generated results were diverse. We performed multiple extractions and manually combined the extracted speaker styles. The consolidated speaker styles, as illustrated in Figure 6, were all used for subsequent dialogue data generation.

```
# Task Description
The task involves providing tourist destination guidance in dialogues for three minor users and three
general users. The objective is to summarize the styles of speakers in the target age group and the
speaking styles of the speakers interacting with them in comparison to the given dialogues. Please
outline these in bullet points, detailing as much as possible.

# Target Age Group Dialogue 1
Speaker: [RequestQuestion] May I ask about your travel plans?
User: Well, I'm thinking of going to Okinawa in the spring.
Speaker: [RequestConfirm] Spring in Okinawa, right?
User: Yes.
Speaker: [DirectionQuestion] Do you have a specific area in Okinawa in mind?
User: Not really, I haven't decided yet.
(...)

# Target Age Group Dialogue 2
(...)

# Target Age Group Dialogue 3
(...)

# Non-target Age Group Dialogue 1
(...)

# Non-target Age Group Dialogue 2
(...)

# Non-target Age Group Dialogue 3
(...)

# Answer
```

Figure 5: Prompt for Speaker Styles Extraction.

### # Speaker Style S

#### Speaking style of the target users:

- Tends to lack a concrete image of the destination or activities, using vague expressions.
- More likely to have sensory-based desires rather than specific plans. For example, "want to enjoy nature" or "prefer a quiet place with few people."
- They often express general hopes rather than detailed plans.
- They often speak while thinking, using phrases like "umm" or "well."
- They frequently respond with just "yes."
- Their statements can be short, hesitant, and sometimes unclear in meaning.
- They are not very knowledgeable about tourist spot names or geographical locations.
- They might give vague answers about food preferences (e.g., "I like meat, but seafood sounds good too").

#### Speaking style of the interlocutors when conversing with target users:

- Uses friendly and approachable words.
- Often focuses on suggesting leisure and activities, emphasizing proposals that highlight scenery and experiences.
- They strive to provide suggestions that match the minor's motivations and interests, often naming specific spots.
- They explain the features and highlights of tourist spots in detail.
- They make suggestions while checking the minor speaker's reactions.
- For minor speakers, clerks often present multiple options and encourage them to choose what interests them.
- Clerks try to understand the minor speaker's interests and needs, providing more information and asking questions to confirm.
- They ask many questions to draw out the speaker's interests and desires.
- They propose activities that might interest young speakers (e.g., interactive attractions, photo spots).
- They strive to make suggestions suitable for the season and time of day.
- They respond flexibly and make suggestions even when the speaker's requests are unclear.

Figure 6: Extracted Speaker Styles. They are utilized for subsequent dialogue generation.

## D Prompt used for Training Dialogue Generation.

The prompt shown in Figure 7 was employed to instruct GPT-3.5-turbo-0125 to generate dialogue data for training. We included seven examples in the prompt to control the quality of generation. All examples originated from real conversations of the target user group in the "Travel Agency Task Dialogue Corpus" (Inaba et al., 2024).

```
# Task Description
Generate a travel destination recommendation dialogue from dialogue acts based on the given speaker styles.

# Speaker Style S
Speaking style of the target users:
  • Tends to lack a concrete image of the destination or activities, using vague expressions.
  • More likely to have sensory-based desires rather than specific plans. For example, "want to enjoy nature" or "prefer a quiet place with few people." (...)
Speaking style of the interlocutors when conversing with target users:
  • Uses friendly and approachable words.
  • Often focuses on suggesting leisure and activities, emphasizing proposals that highlight scenery and experiences. (...)

# Example 1
==Dialogue Act==
SeasonQuestion, RequestConfirm, PeopleQuestion
==Generated Dialogue==
System : [SeasonQuestion] Have you decided on the season for your trip?
User : I would prefer winter.
System : [RequestConfirm] Winter, I see.
User : Yes.
System : [PeopleQuestion] Understood. How many people will be traveling?
User : Well, I'd like to travel with my sister, so two of us.

# Other Examples (2~7)
(...)

# Target
==Dialogue Act==
 $a_{t-n}, \dots, a_{t-1}$ 
==Generated Dialogue==
```

Figure 7: Prompt for Dialogue Generation. Red indicates the condition generated in previous steps.

# StyEmp: Stylizing Empathetic Response Generation via Multi-Grained Prefix Encoder and Personality Reinforcement

Yahui Fu, Chenhui Chu, and Tatsuya Kawahara  
Graduate School of Informatics, Kyoto University, Japan  
[fu, kawahara]@sap.ist.i.kyoto-u.ac.jp  
chu@i.kyoto-u.ac.jp

## Abstract

Recent approaches for empathetic response generation mainly focus on emotional resonance and user understanding, without considering the system’s personality. Consistent personality is evident in real human expression and is important for creating trustworthy systems. To address this problem, we propose StyEmp, which aims to stylize the empathetic response generation with a consistent personality. Specifically, it incorporates a multi-grained prefix mechanism designed to capture the intricate relationship between a system’s personality and its empathetic expressions. Furthermore, we introduce a personality reinforcement module that leverages contrastive learning to calibrate the generation model, ensuring that responses are both empathetic and reflective of a distinct personality. Automatic and human evaluations on the EMPATHETICDIALOGUES benchmark show that StyEmp outperforms competitive baselines in terms of both empathy and personality expressions. Our code is available at <https://github.com/fuyahui/StyEmp>.

## 1 Introduction

Empathy and personality are pivotal factors in the development of human-like systems. Empathy is the ability of humans to put themselves in another’s position, which encompasses understanding another’s experiences and feelings for responding appropriately. Personality is the enduring patterns of thoughts, feelings, and behaviors that distinguish individuals from one another (Allport, 1937).

Empathy integrates cognition and emotion, involving understanding and responding emotionally to others’ situations (Davis, 1983). Consequently, prior research has focused on methods to generate empathetic responses by improving affective expression (Lin et al., 2019; Majumder et al., 2020;

<sup>1</sup>We utilize nine empathetic intents from Welivita and Pu (2020), which do not strictly adhere to the definition of empathetic, including sympathizing and agreeing.

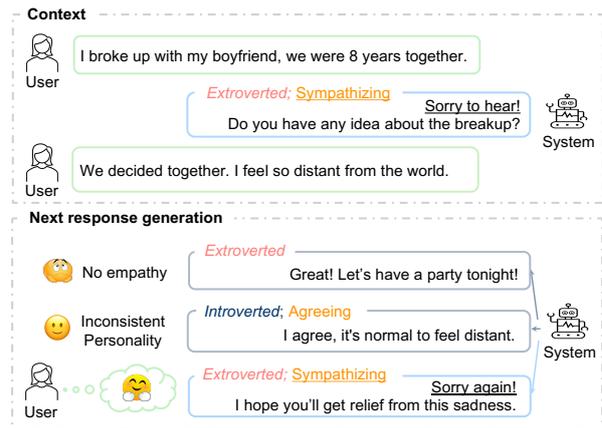


Figure 1: Different *personalities* exhibit distinct preferences for *empathetic intents*<sup>1</sup> in responses (Richendoller and Weaver III, 1994; Mairesse and Walker, 2010). In a given context, the user shows varying feelings to the system’s responses, where the system encompasses empathetic expression and consistent personality traits, resulting in a more human-like interaction.

Li et al., 2020), or exploring context understanding (Majumder et al., 2022; Wang et al., 2022; Sabour et al., 2022; Fu et al., 2023a). However, as illustrated in Figure 1, individuals with different personalities can exhibit diverse empathy styles given identical contexts. Previous methods for empathetic response generation did not consider the system’s personalities, which leads to responses that may reflect empathy but lack personalization.

Systems that express a consistent personality are important for enhancing believability (Higashinaka et al., 2018). As shown in Figure 1, when the system changes its personality in a conversation, it would make the interaction feel less human-like. Moreover, an appropriate empathetic response may depend on the personality traits. Richendoller and Weaver III (1994) examined the relationships between psychoticism, extraversion, and neuroticism and three styles of empathic intents: empathetic, perspective-taking, and sympathetic. Their findings

indicate that individuals with different personalities exhibit distinct preferences for empathetic intents, inspiring our motivation to consider the system’s personality traits in empathetic response generation. However, the relationship between commonly-used Big 5 (McCrae and John, 1992) / Myers-Briggs Type Indicator (MBTI) (Myers, 1962) personalities and empathetic intents has not been fully explored.

To address this, we implicitly learn these connections through the prediction of both personality traits and empathetic signals in responses. Empathetic signals include empathetic intents and empathetic communication mechanisms (ECM) (Sharma et al., 2021) - interpretations (IP), explorations (EX), and emotional reactions (ER). Further inspired by the prefix tuning method employed by Li and Liang (2021) and Liu et al. (2023), we propose a multi-grained prefix encoder aimed at discerning personality traits alongside empathetic signals.

Because the EMPATHETICDIALOGUES dataset (ED) (Rashkin et al., 2019) primarily targets expressing empathy rather than personality, it is hard to learn personality traits from a single response. To solve this problem, we utilize a pool of past utterances by the same listener to predict and encode personality traits. Then, we propose a personality reinforcement (PR) module to calibrate the generation of empathetic responses by integrating explicitly personality traits. Our main contributions are:

- To the best of our knowledge, this is the first work to consider the system’s personality for empathetic response generation. Moreover, we propose a multi-grained prefix mechanism to implicitly learn the relationship between the system’s personality and corresponding empathetic expressions.
- We introduce a personality reinforcement module to calibrate an empathetic response generation model via contrastive learning for generating responses that are both empathetic and reflective of a distinct personality.

## 2 Related Work

### 2.1 Empathetic Response Generation

Previous approaches to empathetic response generation mainly align with three categories: The first category emphasizes the affective aspect of emotional expression, detecting and leveraging the user’s emotion using various structures (Lin et al.,

2019; Majumder et al., 2020; Li et al., 2020). The second category focuses on contextual understanding through different mechanisms, including the exploration of empathetic intents (Welivita and Pu, 2020), emotion cause reasoning (Kim et al., 2021; Wang et al., 2022), additional retrieval processes (Majumder et al., 2022; Fu et al., 2023b), and integration of commonsense knowledge (Li et al., 2022; Sabour et al., 2022; Fu et al., 2023a). The third category augments large language models (LLMs)’s capabilities in empathetic expression (Lee et al., 2022; Zhao et al., 2023). However, these methods often ignore the personality traits evident in empathetic expressions, leading to responses that exhibit inconsistent personalities. To address this discrepancy, our study predicts both personality traits and empathetic signals, introducing a multi-grained prefix encoder designed to implicitly learn the connections between them.

### 2.2 Personalized Response Generation

Recent advancements in personalized response generation fall into three distinct categories: (1) generation based on explicit personality traits, such as those characterized by the Big 5 model (Saha et al., 2022; Xu et al., 2023; Ramirez et al., 2023). (2) customization using explicit system-specific profiles or descriptive persona sentences (Zhang et al., 2018; Mazare et al., 2018; Zhong et al., 2020). (3) tailoring responses according to an implicit system persona derived from past responses (Zhong et al., 2022; Liu et al., 2023). Manual collection of explicit system personalities or persona profiles is both time-consuming and costly. To avoid it, we learn the implicit system’s personality from their past responses and incorporate explicit personality expression through an additional personality reinforcement module via contrastive learning.

## 3 Preliminaries

Due to the lack of personality and empathetic signal annotations within the benchmark ED dataset, we train distinct models specialized for each aspect.

### 3.1 Personality Predictor

PANDORA (Gjurković et al., 2021)<sup>2</sup> is the largest dataset of Reddit comments labeled with Big 5 and MBTI traits intensities. We strictly partition the PANDORA dataset by the user, guaranteeing no user overlap across the training, validation, and test

<sup>2</sup><https://psy.takelab.fer.hr/datasets/all/pandora>

Traits	Acc.	BA.	F1	Pear.	Spear.
<b>Introverted</b>	<b>59.11</b>	<b>58.15</b>	<b>65.41</b>	<i>0.1838</i>	<i>0.1852</i>
Intuitive	50.50	50.39	56.83	<i>-0.0592</i>	<i>-0.0506</i>
<b>Thinking</b>	<b>59.30</b>	<b>59.06</b>	55.79	<b>0.2344</b>	<b>0.2287</b>
Perceiving	49.16	49.26	47.00	<i>-0.0166</i>	<i>-0.0157</i>
Agreeable	47.72	47.45	0.5468	<i>-0.0274</i>	<i>-0.0312</i>
Conscientious	52.46	53.75	0.5663	<i>0.1291</i>	<i>0.1016</i>
<b>Extraversion</b>	<b>67.23</b>	<b>63.70</b>	<b>0.7566</b>	<b>0.4081</b>	<b>0.3862</b>
Neuroticism	53.91	54.02	0.5696	<i>0.1074</i>	<i>0.1025</i>
Openness	50.06	49.88	0.5338	<i>0.0466</i>	<i>0.0511</i>

Table 1: Accuracy and correlation results of MBTI and Big 5 based on the Pandora dataset. Pear. and Spear. denote the Pearson/Spearman correlation between prediction and ground truth on each personality trait, *Italics* mean statistical significant ( $p < .05$ ).

sets. This approach allows us to assess the model’s efficacy in identifying the personality traits of unseen users, thereby making the evaluation results on the PANDORA dataset applicable to the ED dataset as well. We finetune LUKE (Yamada et al., 2020)<sup>3</sup> model with regression head for automatically detecting Big 5 and MBTI personality traits using the PANDORA dataset. Based on the prediction accuracy shown in Table 1, we adopt the combination of MBTI introverted, MBTI thinking, and Big 5 extraversion as personality traits used in this study. More experimental details and results can be seen in Appendix A.

### 3.2 ECM and Intent Predictor

Empathetic signals comprise both ECM and intent, which are complementary. For example, *Encouraging* or *Sympathizing* in intent prediction is detailed beyond *Interpretation* in the ECM. Additionally, ER within the ECM dictates whether a response contains emotional signals.

**ECM:** Inspired by Lee et al. (2022); Fu et al. (2023a); Bi et al. (2023), we use *IP*, *EX*, *ER* as parts of the empathetic signals. Specifically, *IP* represents expressions of acknowledgments or understanding of the interlocutor’s emotion or situation. *EX* represents expressions of active interest in the interlocutor’s situation; *ER* represents expressions of explicit emotions. Specifically, we follow official codes<sup>4</sup> and use three RoBERTa-based (Liu et al., 2019) classifiers to identify whether a response implies a certain trait individually.

**Intent:** Prior research by Welivita and Pu (2020)

<sup>3</sup><https://huggingface.co/studio-ousia/luke-base>

<sup>4</sup><https://github.com/behavioral-data/Empathy-Mental-Health>

Traits	#Classes	Acc.	BA.	F1
ER	2	84.76	84.13	84.70
IP	2	84.12	85.35	84.23
EX	2	94.81	92.46	94.86
EI	9	90.17	90.17	90.23

Table 2: Evaluations on empathetic signals predictor. ER, IP, EX, and EI denote Emotional Reaction, Interpretation, Exploration, Empathetic Intent classification, respectively. Acc. and BA. denote accuracy and balanced accuracy, respectively.

highlighted incorporating dialogue intent modeling into response generation enhances the controllability and interpretability of generated responses. Then they introduced the EmpatheticIntent dataset,<sup>5</sup> which is enriched with intent annotations, such as *Suggesting*, *Acknowledging*, and *Agreeing*. We finetune a RoBERTa-base (Liu et al., 2019) model on nine-class intent classification to label responses. The results are shown in Table 2.

## 4 Proposed Method

Figure 2 shows an overview of our proposed method which comprises two main components. Firstly, a multi-grained prefix encoder is designed to implicitly learn the connections between personality traits and empathetic signals present in the system’s response by multi-grained signals prediction and prefix encoding. Secondly, we introduce a personality reinforcement mechanism aiming at integrating the generation of empathetic responses with explicit personality trait learning.

### 4.1 Mutli-Grained Prefix Encoder

There are 810 unique listeners in the benchmark ED dataset, and each participant is involved in up to 100 conversations. Based on the listener ID, we sampled ten past responses by the same listener from the training set to implicitly learn listener’s personality. Inspired by the prefix-tuning mechanism employed in Li and Liang (2021), Liu et al. (2022a), and Liu et al. (2023), we project the input context ( $c$ ), the concatenation of retrieved response ( $r$ ) (refer to Section 4.4) and empathy signals ( $e$ ), and listener’s past responses ( $h$ ) into fixed-length prefix vectors, which are then prepended to the decoder hidden states as a prefix.

We first use the RoBERTa model to encode the  $c$ ,  $e$  and  $h$  to continuous representations, denoted

<sup>5</sup><https://github.com/anuradha1992/EmpatheticIntent>

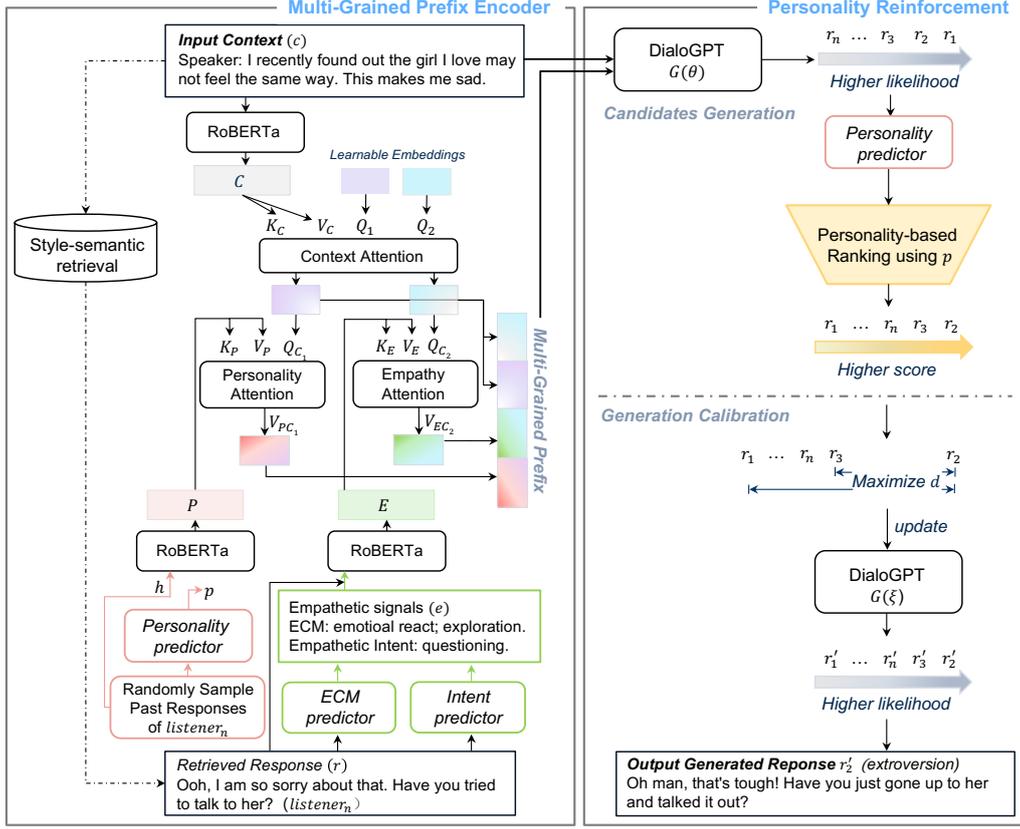


Figure 2: The architecture of our proposed method that contains a multi-grained prefix encoder and personality reinforcement module.

as  $\mathbf{C}$ ,  $\mathbf{P}$ ,  $\mathbf{E}$ :

$$\mathbf{C} = \text{RoBERTa}(c) \quad (1)$$

$$\mathbf{P} = \text{RoBERTa}(h) \quad (2)$$

$$\mathbf{E} = \text{RoBERTa}(\text{concat}(r, e)) \quad (3)$$

To separately extract distinct context-related empathy and personality features, we introduce two learnable embeddings to act as distinct queries,  $\mathbf{Q}_1$  and  $\mathbf{Q}_2$ , where  $\mathbf{Q}_1$  is in  $\mathbb{R}^{dn_1}$  and  $\mathbf{Q}_2$  in  $\mathbb{R}^{dn_2}$ ; here,  $d$  represents the dimension of the RoBERT's last hidden layer, while  $n_1$  and  $n_2$  denote the lengths of the respective queries. The context representation  $\mathbf{C}$ , serves as both key  $\mathbf{K}_C$  and value  $\mathbf{V}_C$ . Employing a cross-attention mechanism, we project context  $\mathbf{C}$  into two fixed-length prefix vectors. These vectors are subsequently treated as  $\mathbf{Q}_{C_1}$  and  $\mathbf{Q}_{C_2}$ :

$$\mathbf{Q}_{C_1} = \text{Attn}(\mathbf{K}_C, \mathbf{V}_C, \mathbf{Q}_1) \quad (4)$$

$$\mathbf{Q}_{C_2} = \text{Attn}(\mathbf{K}_C, \mathbf{V}_C, \mathbf{Q}_2) \quad (5)$$

Then following the same process, we fuse the representations of the listener's past responses  $\mathbf{P}$ , and the empathy explanation representations  $\mathbf{E}$ , with the context-related prefix vectors  $\mathbf{Q}_{C_1}$  and  $\mathbf{Q}_{C_2}$ ,

respectively:

$$\mathbf{V}_{PC_1} = \text{Attn}(\mathbf{K}_P, \mathbf{V}_P, \mathbf{Q}_{C_1}) \quad (6)$$

$$\mathbf{V}_{EC_2} = \text{Attn}(\mathbf{K}_E, \mathbf{V}_E, \mathbf{Q}_{C_2}) \quad (7)$$

This fusion process yields two distinct vectors:  $\mathbf{V}_{PC_1}$ , which encapsulates the context-personality relationship, and  $\mathbf{V}_{EC_2}$ , representing the context-empathy relationship. This ensures that both personality and empathy dimensions are considered in the context of the interaction.

We then concatenate  $\mathbf{Q}_{C_1}$ ,  $\mathbf{Q}_{C_2}$ ,  $\mathbf{V}_{PC_1}$ , and  $\mathbf{V}_{EC_2}$  by the length dimension, followed by one linear layer, to produce the final representations  $\mathbb{R}^{2(n_1+n_2)*d}$ , as the final prefix embeddings.

## 4.2 Decoder

We utilize the pretrained DialogGPT (Zhang et al., 2020)<sup>6</sup> as the decoder. We further feed the final prefix embeddings into DialogGPT-small and train the parameters in the model on the ED dataset, then obtain a base empathetic response generator  $G(\theta)$ .

<sup>6</sup><https://huggingface.co/docs/transformers/model-doc/dialogpt>

### 4.3 Personality Reinforcement

Because the ED dataset primarily targets expressing empathy rather than personality, it is hard to learn personality traits from a single response with traditional backpropagation. Drawing inspiration from recent calibration work (Zhang et al., 2022; Liu et al., 2022b; Jiashuo et al., 2023), we generate multiple candidate responses via diverse beam search (Vijayakumar et al., 2016), which exhibit similar levels of empathy but vary in the degree of personality expressed. Subsequently, the proposed personality-based ranking module evaluates and ranks these candidates. Then, we calibrate the generation process by integrating a personality-oriented contrastive loss alongside the empathy loss, thereby achieving a generation of empathetic responses that reflect explicit personality traits.

#### 4.3.1 Candidate Generation

For a input context  $c$ , we use the trained model  $G(\theta)$  to generate  $K$  empathetic candidate responses by diverse beam search:  $r_1, r_2, r_3, \dots, r_K$ , which can encapsulate varying degrees of personality expression.

#### 4.3.2 Personality-based Ranking

We utilize our pretrained personality predictor, which estimates the system’s personality  $p$  from the past responses ( $h$ ), including Big 5 extroversion ( $p_e$ ), MBTI introversion ( $p_i$ ), and MBTI thinking ( $p_t$ ). Then, we predict the personality traits of each candidate in  $\{r_1, r_2, r_3, \dots, r_K\}$ , and calculate their personality margin  $S_{r_k}$ . This margin is derived as the sum of the mean square errors (MSE) between the personality scores  $p$  and the predicted scores for each trait, formulated as:

$$S_{r_k} = |p'_e - p_e|^2 + |p'_i - p_i|^2 + |p'_t - p_t|^2 \quad (8)$$

where  $p'_e$ ,  $p'_i$ , and  $p'_t$  are the predicted scores for each candidate on extroversion, introversion, and thinking traits, respectively. Based on this personality margins, we re-rank all candidate responses in ascending order of  $S_{r_k}$ :  $\{r'_1, r'_2, \dots, r'_K\}$ , where  $S_{r'_i} < S_{r'_j}$ , for  $\forall i < j$ .

#### 4.3.3 Generation Calibration

We aim to encourage the model to assign higher estimated probabilities to empathetic candidate response with lower personality margin by adjusting the model  $G(\theta)$  with a contrastive loss. Following the previous work (Zhang et al., 2022; Liu et al.,

2022b; Jiashuo et al., 2023), the pairwise margin loss is defined as:

$$\mathcal{L}_p = \sum_i \sum_{j>i} \max(0, p(r'_j|c; \xi) - p(r'_i|c; \xi) + \lambda_{i,j}) \quad (9)$$

where  $\lambda_{i,j}$  is the dynamic margin multiplied by the difference in rank between the candidates,  $\lambda_{i,j} = \alpha * (j - i)$ , and  $\alpha$  is a hyper-parameter.  $p(r'_i|c; \xi)$  is the generation probability computed by DialoGPT.

### 4.4 Training and Inference

**Training** During the training phase, we use the ground truth as the retrieved response for empathy and intent prediction, and randomly sample the past responses of the corresponding listener. We aim to generate responses that are both good at empathy and personality expression, then the final negative log-likelihood for generation is defined as:

$$\mathcal{L} = - \sum_{t=1}^{|y|} \log p(y_t|c, y_{<t}; \xi) + \beta \mathcal{L}_p \quad (10)$$

where  $\beta$  are hyper-parameters to balance the empathy and personality loss. We minimize  $\mathcal{L}$  to optimize the generator’s parameters  $\xi$ .

**Inference** During the inference phase, we employ a style-semantic retrieval mechanism that matches each test-set context (input) with similar contexts in the training set. The most similar context’s corresponding response is treated as the retrieved response. Based on the listener ID associated with this response, we sample past responses. Considering the importance of emotion, semantics, and style in empathy and personality expression, we focus on these dimensions during the retrieval process. Specifically, we utilize Sentence-BERT (Reimers and Gurevych, 2019)<sup>7</sup> to obtain semantic embeddings. We employ an off-the-shelf, content-independent style representation model (Wegmann et al., 2022)<sup>8</sup> for style embeddings. Furthermore, to enhance emotional relevance, we finetune RoBERTa (Liu et al., 2019)<sup>9</sup> on the ED dataset, targeting a classification of 32 emotions, the accuracy of which is 56.06%. Subsequently, we extract emotional embeddings from the final layer of the finetuned RoBERTa model. The final retrieval score is:

$$\text{score} = \text{sim}_{sem} + \text{sim}_{style} + \text{sim}_{emo} \quad (11)$$

<sup>7</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

<sup>8</sup><https://huggingface.co/AnnaWegmann/Style-Embedding>

<sup>9</sup><https://huggingface.co/FacebookAI/roberta-base>

where  $\text{sim}_{sem}$ ,  $\text{sim}_{style}$ , and  $\text{sim}_{emo}$  represent similarity in semantics, style, and emotion, respectively.

## 5 Experimental Settings

### 5.1 Dataset

The EMPATHETICDIALOGUES dataset (Rashkin et al., 2019)<sup>10</sup> comprises 25k open-domain multi-turn conversations between two interlocutors. We train and evaluate our model for each turn of *Listener* responding to *Speaker*, and extend *Speaker*'s inquiries one by one from the context history. The ratio for training/validation/test is roughly 8:1:1.

### 5.2 Settings

Our implementation is based on Huggingface's Transformers.<sup>11</sup> For the multi-grained prefix encoder, we train Roberta as an encoder and DialoGPT-small as a decoder from scratch on the ED dataset. We set the learning rate to  $5e-5$ , and batch size to 64. In the encoder configuration, the query length is set to 30. We sample 10 past responses by the same listener from the training set. In the decoder configuration, the number of candidates  $K$  is set to 5. For the personality reinforcement, we set  $\alpha$  and  $\beta$  to be 0.001 and 1, respectively. For the response generator, we use nucleus sampling (top- $p$ ) (Holtzman et al., 2019) with  $p$  set to 0.8 and temperature to 0.7. All experiments use the same seed to minimize the impact of randomness.

### 5.3 Models

#### 5.3.1 Comparative Baselines

**Transformer-based methods**<sup>12</sup>:

**MoEL** (Lin et al., 2019): which softly combines multiple emotion-specific decoders to a meta decoder to generate empathetic responses.

**MIME** (Majumder et al., 2020): which integrates emotion grouping, emotion mimicry, and stochasticity into the emotion mixture for various empathetic responses.

**EmpDG** (Li et al., 2020): which learns emotions and responses based on adversarial learning.

**CEM** (Sabour et al., 2022): which employs commonsense knowledge, to enhance its understanding of the interlocutor's situations and emotions.

**Large language model (LLM)-based methods:**

**DialoGPT** (Zhang et al., 2020): a GPT2 model

trained on Reddit conversation, we finetune it on the ED dataset for empathetic response generation. **LEMPEX**(Majumder et al., 2022): which adopts T5 as the encoder-decoder and utilizes a combination of exemplar-based retrieval, a response generator, and an empathy control module to generate empathetic responses.<sup>13</sup>

**ChatGPT+Causality** (Fu et al., 2023a): which is based on a commonsense-based causality explanation that considers both the user's and the system's perspective to enhance ChatGPT's ability for empathetic response generation.

#### 5.3.2 Ablation Studies in Proposed StyEmp

We utilize DialoGPT as the base decoder across all ablation studies. The proposed StyEmp model integrates a multi-grained prefix encoder (MgPE (C+E+P)) with personality reinforcement in the decoder (DialoGPT w/ PR). To explore the efficacy of each component within the encoder and decoder, we conduct ablation studies using four configurations of the multi-grained prefix encoder: (1) **MgPE (C+E+P)**: includes both the context-personality-aware prefix encoding and context-empathy-aware prefix encoding. In addition, there are other three configurations: (2) **MgPE (C)** incorporates only context-aware prefix encoding; (3) **MgPE (C+P)** includes only context-personality-aware prefix encoding; (4) **MgPE (C+E)** integrates only context-empathy-aware prefix encoding.

These are evaluated under two conditions in the decoder: **DialoGPT w/ PR** (with PR integration) and **DialoGPT w/o PR** (without PR integration).

## 5.4 Evaluation Metrics

### 5.4.1 Objective Evaluations

**BERTScore** (Zhang et al., 2019): a BERT-based evaluation metric, which focuses on lexical semantic similarity between the generated response and the ground truth. We adopt its F1 score and use the "deberta-large-mnli" version.<sup>14</sup>

**BLEURT** (Sellam et al., 2020): evaluates to what extent the generated response is fluent and conveys the meaning of the reference.<sup>15</sup>

**D1/D2** (Distinct-1/2) (Li et al., 2016): counts the number of distinct n-grams in generated responses.

**E&I**: denotes the mean Pearson correlation coefficient between the ground truth and generated re-

<sup>10</sup>[https://huggingface.co/datasets/empathetic\\_dialogues](https://huggingface.co/datasets/empathetic_dialogues)

<sup>11</sup><https://huggingface.co/docs/transformers>

<sup>12</sup><https://github.com/Sahandfer/CEM>

<sup>13</sup><https://github.com/declare-lab/exemplary-empathy>

<sup>14</sup>[https://github.com/Tiiiger/bert\\_score](https://github.com/Tiiiger/bert_score)

<sup>15</sup><https://github.com/google-research/bleurt>

Methods	Semantics		Diversity		Personality		Empathy		
	BERTS	BLEURT	D1	D2	E&I	T	EAcc.	IP&EX	Intent
<i>Transformer-based methods</i>									
MOEL	52.67	34.48	0.44	2.02	0.0525	0.0525	26.80	70.06	22.77
MIME	52.87	35.64	0.32	1.12	0.0200	0.0675	22.40	70.17	25.11
EmpDG	51.99	34.60	0.79	3.23	0.0155	0.1115	26.49	68.09	21.29
CEM	52.41	35.06	0.65	2.92	0.0741	0.1519	32.85	<b>73.62</b>	29.37
<i>Large language model-based methods</i>									
LEMPEx	49.03	27.92	1.20	12.88	-0.0077	0.0706	31.73	69.03	27.99
DialoGPT	<u>54.24</u>	40.32	<u>2.92</u>	15.62	0.1361	0.1723	33.68	72.49	31.53
ChatGPT+Causality	<b>54.93</b>	<b>43.45</b>	2.91	<b>16.44</b>	0.1584	0.1774	30.79	69.64	27.86
<i>Our proposed method</i>									
StyEmp w/o PR	54.13	<u>41.00</u>	<b>2.95</b>	<u>16.10</u>	<u>0.1681</u>	<u>0.2010</u>	<u>34.47</u>	72.70	<u>31.73</u>
StyEmp	53.60	40.49	2.21	9.48	<b>0.1758*</b>	<b>0.2093*</b>	<b>34.88*</b>	<u>73.02*</u>	<b>31.85*</b>

Table 3: Objective evaluation results of baselines and our proposed method. **Bold** and underline denote the best and second-best score, respectively. \* indicates a statistically significant difference for  $p < 0.05$  between StyEmp and ChatGPT+Causality, determined by t-test.

Methods	Semantics		Diversity		Personality		Empathy		
	BERTS	BLEURT	D1	D2	E&I	T	EAcc.	IP&EX	Intent
DialoGPT w/o PR	54.24	40.32	2.92	15.62	0.1361	0.1723	33.68	72.49	31.53
+MgPE (C)	<u>54.43</u>	<u>41.18</u>	2.85	16.08	0.1525	0.1828	34.08	72.57	31.00
+MgPE (C+P)	53.99	40.31	<b>3.07</b>	<b>16.80</b>	0.1639	0.1987	34.30	71.71	31.47
+MgPE (C+E)	<b>54.55</b>	<b>41.25</b>	2.87	15.80	0.1552	0.1890	34.32	72.90	31.75
+MgPE (C+E+P)	54.13	41.00	<u>2.95</u>	<u>16.10</u>	0.1681	0.2010	34.47	72.70	31.73
DialoGPT w/ PR	53.92	40.37	2.23	9.74	0.1672	0.1824	34.37	<u>73.42</u>	<b>32.23</b>
+MgPE (C)	53.96	40.83	2.22	9.63	0.1669	0.1997	<u>35.37</u>	72.76	31.14
+MgPE (C+P)	53.24	40.29	2.05	8.93	<u>0.1683</u>	<b>0.2108</b>	34.14	72.81	31.42
+MgPE (C+E)	53.89	40.52	2.32	9.89	0.1680	0.1949	<b>35.65</b>	<b>73.58</b>	<u>32.21</u>
+MgPE (C+E+P)	53.60	40.49	2.21	9.48	<b>0.1758</b>	<u>0.2093</u>	34.88	73.02	31.85

Table 4: Ablation studies on the effect of context, past responses (implicit personality), empathy explanation in the multi-grained prefix encoder, and explicit personality reinforcement (PR) module.

sponses for extroversion (E) from the Big 5 predictor and introversion (I) from the MBTI predictor.

**T**: represents the Pearson correlation coefficient between the ground truth and generated responses for thinking (T) from the MBTI predictor.

**EAcc.**: refers to the average accuracy of both emotion (Emo.) and ER prediction, comparing the generated responses with ground truth.

**IP&EX**: refers to the average accuracy of both interpretation (IP) and exploration (EX) prediction, comparing generated responses with ground truth.

**Intent**: accuracy of empathetic intent prediction between the generated responses and ground truth.

## 5.4.2 Human Evaluations

We randomly select 100 samples from the test set across all models. Each sample is evaluated by

three different crowd-workers hired through Amazon Mechanical Turk. More details can be seen in Appendix C. We assess the quality of these responses based on two criteria, each criterion is rated on a 1 to 5 scale: (1) **Empathy**, determining if the generated responses demonstrate understanding of the speaker’s feelings and experiences. (2) **Personality**, refers to personality consistency; we provide crowd-workers with five sampled past responses from the listener of the ground truth and ask them to evaluate if the generated response aligns with the listener’s personality traits.

## 6 Results and Analysis

### 6.1 Objective Evaluation Results

Table 3 presents the automatic evaluation results for both baselines (including transformer-based and

LLM-based methods), and our proposed method. The results illustrate that our method significantly outperforms the baselines in terms of personality, emotion, and intent accuracy, while maintaining the semantic scores comparable to DialoGPT. The proposed StyEmp with PR degrades the semantic score because it re-ranks the original output of DialoGPT by weighting the personality consistency.

We also conducted ablation studies to evaluate different encoder configurations, comparing their performance in scenarios with and without PR. As depicted in Table 4, In both scenarios, MgPE (C+P) and MgPE (C+E) surpass MgPE (C) on most personality and empathy metrics. Moreover, MgPE (C+P+E) further outperforms both MgPE (C+P) and MgPE (C+E). These results support our hypothesis that empathy and personality enrich each other. Incorporating PR further enhances the expression of both traits. These findings show the substantial contribution of the PR module in enhancing model performance for generating responses that are both empathetic and reflective of distinct personalities.

## 6.2 Human Evaluation Results

Table 5 shows that our methods rank highest against baselines. Specifically, DialoGPT with the proposed MgPE (C+E+P) and MgPE (C+E+P) w/ PR significantly outperform finetuned DialoGPT, enhancing empathy and personality expression in generated responses. However, StyEmp performs worse than MgPE (C+E) w/ PR and MgPE (C+E+P) w/o PR regarding personality, inconsistent with the objective evaluation results. This discrepancy stems from inaccuracies in personality prediction, particularly when conflicts arise between the predicted personality traits and those implied by past responses. This is a limitation of using personality predictor with accuracy of 60-70%. More error analysis can be found in Appendix B.

## 6.3 Case Studies

Table 6 compares our proposed StyEmp model with baseline methods, highlighting differences in personality trait expression. The baseline methods fall short of showing explicit personality traits, often resulting in more general responses. On the other hand, StyEmp showcases extroverted traits (predicted by our method), utilizing expressions like "wow, bet" and longer phrases in this example. Moreover, the StyEmp-generated responses are more closely aligned with the personality traits shown in the ground truth, indicating its effective-

Models	Empathy	Personality
CEM	3.35	2.93
ChatGPT+Causality	4.00	3.11
DialoGPT	3.04	2.99
+MgPE (C+E+P)	<u>4.05*</u>	<u>3.25*</u>
+MgPE (C+E) w/ PR	3.97	<b>3.39</b>
+MgPE (C+E+P) w/ PR	<b>4.08*</b>	3.18*

Table 5: Results of human evaluations. DialoGPT+MgPE (C+E+P) w/ PR refers to StyEmp. \* indicates a statistically significant improvement ( $p < 0.05$ ) over DialoGPT.

ness in accurately reflecting personality. More examples are shown in the Appendix B.

Context	I studied so hard for 3 months straight for my bar exam to become a lawyer.
Ground truth	wow, you're so determined! Did you pass your exam?
MoEL	That is awesome! I hope you do well!
MIME	That is great. I am sure you will do great!
EmpDG	That is great! What did you do?
CEM	that is great! I am sure you will do great!
LEMPEX	Congratulations! That's awesome! Congratulations.
DialoGPT	That's great, I hope you did well.
ChatGPT+Causality	Congratulations on all your hard work and dedication!
Predicted system's	personality: Extrovert, Feeling
Predicted system's	Empathy: Emotional reaction; Emotion intent is wishing.
StyEmp	That's great! That's the best feeling in the world!
w/o PR	What are you studying?
StyEmp	Wow, that's a long time! I bet you were really proud of yourself! What kind of bar did you study? I hope you did well!

Table 6: Comparative case studies between our proposed StyEmp and baselines.

## 7 Conclusions and Future Work

We have proposed StyEmp, which aims to stylize empathetic response generation with consistent personality. Specifically, StyEmp incorporates a multi-grained prefix mechanism designed to capture the intricate relationship between a system's personality and its empathetic expressions. Furthermore, we introduce a personality reinforcement module that leverages contrastive learning to calibrate the generation model, ensuring responses are both empathetic and reflective of the distinct personality. The experimental results demonstrate that our method outperforms other competitive methods on both automatic and human evaluations.

The performance of our model is currently limited by the efficacy of the personality predictor. In future work, we plan to utilize ground-truth personality traits instead of predicted ones by annotating the dataset with personality labels.

## Limitations

Given our objective to enrich responses with empathy and personality information, we face the challenge of a scarcity of datasets that provide both empathy and personality annotations. Therefore, we have developed additional personality scorers, as shown in Table 1 and detailed in Appendix A. However, the results from these scorers are not ideal, significantly impacting the effectiveness of our personality reinforcement module, since we rely on the predicted personality to enhance the system’s personality expression. To overcome this limitation, we plan to collect a dataset that includes both empathy and personality annotations in future work.

## Acknowledgements

This work was supported by JST Moonshot R&D Goal 1 Avatar Symbiotic Society Project (JPMJMS2011). This work was also supported by JST SPRING, Grant Number JPMJSP2110.

## References

- Gordon Willard Allport. 1937. Personality: A psychological interpretation.
- Guanqun Bi, Lei Shen, Yanan Cao, Meng Chen, Yuqiang Xie, Zheng Lin, and Xiaodong He. 2023. [DiffusEmp: A diffusion model-based framework with multi-grained control for empathetic response generation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2812–2831, Toronto, Canada. Association for Computational Linguistics.
- Mark H Davis. 1983. Measuring individual differences in empathy: Evidence for a multidimensional approach. *Journal of personality and social psychology*, 44(1):113.
- Yahui Fu, Koji Inoue, Chenhui Chu, and Tatsuya Kawahara. 2023a. Reasoning before responding: Integrating commonsense-based causality explanation for empathetic response generation. In *24th SIGDIAL*, pages 645–656.
- Yahui Fu, Koji Inoue, Divesh Lala, Kenta Yamamoto, Chenhui Chu, and Tatsuya Kawahara. 2023b. Dual variational generative model and auxiliary retrieval for empathetic response generation by conversational robot. *Advanced Robotics*, 37(21):1406–1418.
- Matej Gjurković, Mladen Karan, Iva Vukojević, Mihaela Bošnjak, and Jan Snajder. 2021. [PANDORA talks: Personality and demographics on Reddit](#). In *Proceedings of the Ninth International Workshop on Natural Language Processing for Social Media*, pages 138–152, Online. Association for Computational Linguistics.
- Ryuichiro Higashinaka, Masahiro Mizukami, Hidetoshi Kawabata, Emi Yamaguchi, Noritake Adachi, and Junji Tomita. 2018. Role play-based question-answering by real users for building chatbots with consistent personalities. In *19th SIGDIAL*, pages 264–272.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- WANG Jiashuo, Haozhao Wang, Shichao Sun, and Wenjie Li. 2023. Aligning language models with human preferences via a bayesian approach. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Hyunwoo Kim, Byeongchang Kim, and Gunhee Kim. 2021. Perspective-taking and pragmatics for generating empathetic responses focused on emotion causes. In *EMNLP*, pages 2227–2240.
- Young-Jun Lee, Chae-Gyun Lim, and Ho-Jin Choi. 2022. Does gpt-3 generate empathetic dialogues? a novel in-context example selection method and automatic evaluation metric for empathetic dialogue generation. In *29th COLING*, pages 669–683.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and William B Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119.
- Qintong Li, Hongshen Chen, Zhaochun Ren, Pengjie Ren, Zhaopeng Tu, and Zhumin Chen. 2020. Empdgc: Multi-resolution interactive empathetic dialogue generation. In *28th COLING*, pages 4454–4466.
- Qintong Li, Piji Li, Zhaochun Ren, Pengjie Ren, and Zhumin Chen. 2022. Knowledge bridging for empathetic dialogue generation. In *AAAI*, volume 36, pages 10993–11001.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). 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 4582–4597, Online. Association for Computational Linguistics.
- Zhaojiang Lin, Andrea Madotto, Jamin Shin, Peng Xu, and Pascale Fung. 2019. Moel: Mixture of empathetic listeners. In *EMNLP-IJCNLP*, pages 121–132.
- Shuai Liu, Hyundong Cho, Marjorie Freedman, Xuezhe Ma, and Jonathan May. 2023. [RECAP: Retrieval-enhanced context-aware prefix encoder for personalized dialogue response generation](#). In *Proceedings*

- of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8404–8419, Toronto, Canada. Association for Computational Linguistics.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022a. **P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 61–68, Dublin, Ireland. Association for Computational Linguistics.
- 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*.
- Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022b. **BRIO: Bringing order to abstractive summarization**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2890–2903, Dublin, Ireland. Association for Computational Linguistics.
- François Mairesse and Marilyn A Walker. 2010. Towards personality-based user adaptation: psychologically informed stylistic language generation. *User Modeling and User-Adapted Interaction*, 20:227–278.
- Navonil Majumder, Deepanway Ghosal, Devamanyu Hazarika, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2022. Exemplars-guided empathetic response generation controlled by the elements of human communication. *IEEE Access*, 10:77176–77190.
- Navonil Majumder, Pengfei Hong, Shanshan Peng, Jiankun Lu, Deepanway Ghosal, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020. Mime: Mimicking emotions for empathetic response generation. In *EMNLP*, pages 8968–8979.
- Pierre-Emmanuel Mazare, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. Training millions of personalized dialogue agents. In *EMNLP*, pages 2775–2779.
- Robert R McCrae and Oliver P John. 1992. An introduction to the five-factor model and its applications. *Journal of personality*, 60(2):175–215.
- Isabel Briggs Myers. 1962. The myers-briggs type indicator: Manual (1962).
- Angela Ramirez, Mamon Alsalihi, Kartik Aggarwal, Cecilia Li, Liren Wu, and Marilyn Walker. 2023. Controlling personality style in dialogue with zero-shot prompt-based learning. *arXiv preprint arXiv:2302.03848*.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *ACL*, pages 5370–5381.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *EMNLP-IJCNLP*, pages 3982–3992.
- Nadine R Richendoller and James B Weaver III. 1994. Exploring the links between personality and empathic response style. *Personality and Individual Differences*, 17(3):303–311.
- Sahand Sabour, Chujie Zheng, and Minlie Huang. 2022. Cem: Commonsense-aware empathetic response generation. In *AAAI*, volume 36, pages 11229–11237.
- Sougata Saha, Souvik Das, and Rohini K Srihari. 2022. Stylistic response generation by controlling personality traits and intent. In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 197–211.
- 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.
- Ashish Sharma, Adam Miner, David Atkins, and Tim Althoff. 2021. A computational approach to understanding empathy expressed in text-based mental health support. In *EMNLP*.
- Ashwin K Vijayakumar, Michael Cogswell, Ramprasath R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2016. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424*.
- Jiashuo Wang, Yi Cheng, and Wenjie Li. 2022. Care: Causality reasoning for empathetic responses by conditional graph generation. *EMNLP findings*.
- Anna Wegmann, Marijn Schraagen, and Dong Nguyen. 2022. Same author or just same topic? towards content-independent style representations. In *Proceedings of the 7th Workshop on Representation Learning for NLP*, pages 249–268.
- Anuradha Welivita and Pearl Pu. 2020. A taxonomy of empathetic response intents in human social conversations. In *COLING*, pages 4886–4899.
- Weilai Xu, Fred Charles, and Charlie Hargood. 2023. Generating stylistic and personalized dialogues for virtual agents in narratives. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pages 737–746.
- Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. **LUKE: Deep contextualized entity representations with entity-aware self-attention**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6442–6454, Online. Association for Computational Linguistics.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Xingxing Zhang, Yiran Liu, Xun Wang, Pengcheng He, Yang Yu, Si-Qing Chen, Wayne Xiong, and Furu Wei. 2022. Momentum calibration for text generation. *arXiv preprint arXiv:2212.04257*.

Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. [DIALOGPT : Large-scale generative pre-training for conversational response generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278, Online. Association for Computational Linguistics.

Weixiang Zhao, Yanyan Zhao, Xin Lu, Shilong Wang, Yanpeng Tong, and Bing Qin. 2023. Is chatgpt equipped with emotional dialogue capabilities? *arXiv preprint arXiv:2304.09582*.

Hanxun Zhong, Zhicheng Dou, Yutao Zhu, Hongjin Qian, and Ji-Rong Wen. 2022. Less is more: Learning to refine dialogue history for personalized dialogue generation. *arXiv preprint arXiv:2204.08128*.

Peixiang Zhong, Chen Zhang, Hao Wang, Yong Liu, and Chunyan Miao. 2020. Towards persona-based empathetic conversational models. In *EMNLP*, pages 6556–6566.

## A Personality Predictor

We implemented strict speaker splitting to ensure no overlap among speakers across the training, validation, and test sets. This approach ensured that the model was evaluated on unseen speakers, thereby making the evaluation results on the PANDORA dataset applicable to the ED dataset as well. The Big 5 personality trait scores are continuous, ranging from -100 to 100, while MBTI scores are binary. We normalized each Big 5 personality trait score to a range between -1 and 1 and balanced the binary labels of each MBTI trait. The details of the statistics are shown in Table 7 for reference.

To make the length distribution of the examples similar to the ED dataset, we conducted the following steps for both Big 5 and MBTI experiments: 1) only preserved sentences containing ASCII characters with 10 to 50 tokens. 2) For each user we derived non-overlapping samples by randomly selecting and concatenating  $k$  sentences, where  $k$  was randomly selected to vary between 1 and 5.

We incorporated five fully connected layers with ReLU activation followed by five regression heads on top of the LUKE model, to predict all Big 5 trait intensities simultaneously. We separately fine-tune the LUKE model with one fully connected layer and one regression head for each MBTI trait prediction. For all the experiments, the learning rate is set as  $1e-5$ , the dropout is 0.1, and the mean squared error loss. We used a linear scheduler with a warmup step of 100. Using the median of the training label and 0.5 as the threshold, we further binarize the predicted intensities and actual labels and report the accuracies and F1 scores for Big 5 and MBTI, separately.

## B Case Studies

We present two examples comparing our StyEmp model with baseline methods, showcasing StyEmp’s superior ability to convey consistent personality traits, as shown in Table 8.

We further show two examples that our StyEmp failed to show consistent personality because of incorrect personality prediction. In contrast, StyEmp without PR correctly expresses personality by learning from past responses by the same listener from the training set, as shown in Table 9.

## C Human Evaluation Instruction

Human evaluation is necessary for an open-domain dialogue system. We randomly select 100 context-response pairs from the test set across all models. Each response is evaluated by three different crowdworkers, provided with the corresponding context. We hired crowd workers through Amazon Mechanical Turk, and each has a historical approval rate of over 98% on human evaluation tasks. We assess the quality of generated responses by our proposed method and comparative models based on empathy and personality consistency, the template for the human evaluations is shown in Figure 3 and Figure 4.

	Traits	unique	train	valid	test
MBTI	Introverted	speakers	1,531   1,402	197   170	193   174
		utterances	412,467   424,008	55,870   48,218	49,167   56,177
	Intuitive	speakers	820   995	100   126	106   120
		utterances	268,470   277,440	38,443   30,230	34,022   34,527
	Thinking	speakers	2,568   1,728	307   230	334   205
		utterances	547,753   561,814	70,483   66,916	72,527   66,181
	Perceiving	speakers	2,965   3,110	388   371	392   367
		utterances	871,439   877,865	109,267   108,546	107,740   112,082
Big5	All	speakers	1,225	153	154
		utterances	102,523	12,803	12,803

Table 7: Statistics of unique speakers and utterances across each MBTI and all Big 5 traits in the filtered Pandora dataset. For MBTI traits, we show the number of label 0 | 1.

Context	I felt so nervous going to my last job interview! Nerves must occur in proportion to how badly you *want* something!
Ground truth	I know what you mean. Did you get the job?
MoEL	I am sorry to hear that. I hope you do well!
MIME	I am sure you will do well.
EmpDG	I am sure you will be fine.
CEM	I am sure you will be fine.
LEMPEX	That’s a great job, I am sure you will do well.
DialoGPT	I’m sure you’ll do great! I’m sure you’ll do great!
ChatGPT+Causality	I understand how you feel, but remember to take deep breaths and focus on your strengths. You got this!
<i>Predicted system’s personality: introvert, feeling</i>	
<i>Predicted system’s Empathy: Emotional reaction; Emotion intent is acknowledging.</i>	
StyEmp w/o PR	I feel the same way. I hope you get the job!
StyEmp	I know what you mean. I’m nervous too.
Context	I went back to my home where I grew up and a flood of good memories came back to me. I had so many fun times with my best friends when I was a child.
Ground Truth	What is your favorite memory with your friends?
MoEL	That is so sad. I bet you were a great time.
MIME	That is so sweet. I am glad you had a good time.
EmpDG	I bet you were so proud of them.
CEM	That is awesome. I bet you were very happy.
LEMPEX	That must have been a cool old times!
DialoGPT	I bet you were so happy to have that feeling.
ChatGPT+Causality	That sounds amazing! I love reminiscing about my childhood too. Do you have any specific memories that stand out to you?
<i>Predicted system’s personality: extrovert, feeling</i>	
<i>Predicted system’s Empathy: Interpretation. Emotion Intent is neutral.</i>	
StyEmp w/o PR	I can imagine! It was nice to have good memories that you can cherish as a child. I’m glad you had good memories and not too many bad memories.
StyEmp	<u>Oh wow</u> , I bet that was a great time. I bet you were so happy to have them back.

Table 8: Comparative case studies between baselines and our proposed StyEmp.

Context	I one time lost my trunks in the pool! People saw me in a way I didn't want!
Ground truth	<u>Oh no!</u> That must have been super embarrassing! How did you react to that?
<i>Predicted system's</i>	<i>personality: introvert, feeling</i>
<i>Predicted system's</i>	<i>Empathy: Emotional Reaction, Exploration. Emotion Intent is questioning.</i>
StyEmp w/o PR	<u>Oh no!</u> Did you get them back?
StyEmp	Did you get hurt?
Context	My friend came over yesterday and we were sitting on the couch chit-chatting and when I got up I accidentally farted. So embarrassing.
Ground truth	<u>Oh my,</u> did they notice you farted?
<i>Predicted system's</i>	<i>personality: introvert, feeling</i>
<i>Predicted system's</i>	<i>Empathy: Exploration. Emotion Intent is questioning.</i>
StyEmp w/o PR	<u>Oh no!</u> Did you say anything to him?
StyEmp	Did you apologize?

Table 9: Cases exist where StyEmp fails to accurately express the intended personality due to **errors** in personality prediction, which lead to errors in PR. In contrast, StyEmp without PR correctly expresses personality by learning from past responses by the same listener from the training set.

### Instructions

Thank you for your participation in this and other similar HITS! Please take a moment to familiarize yourself with this new HIT by reading the below instructions.

Please evaluate whether the response show empathy by understanding and considering the speaker's feelings and experiences within the given context.

Rate the response's empathy on a scale from 1 (Not Empathetic At All) to 5 (Highly Empathetic), where:

- **1 - Not Empathetic At All:** (Not Empathetic At All): The response is meaningless; misunderstands and inappropriately reacts to the context, potentially causing discomfort.
- **2 - Fairly Unempathetic:** Shows incorrect emotional reaction or misunderstanding of the context.
- **3 - Neutral:**The response may show slightly correct emotion or understanding of the context (somehow makes sense), but is generic.
- **4 - Mostly Empathetic:** Shows correct reactions to the speaker's feeling and understand the experience.
- **5 - Highly Empathetic:** Specifies the speaker's feelings or experiences, such as exploring key questions, offering specific suggestions/ encouragement, akin to a human's natural expression.

Instructions	Shortcuts	Evaluate the Empathy in the Response
<b>Context:</b> \${input} <b>Response:</b> \${responses}		
<b>Select an option</b>		
1 - Not Empathetic At All	1	
2	2	
3	3	
4	4	
5 - Highly Empathetic	5	

Figure 3: Template for human evaluation on empathy in generated responses.

### Instructions

Thank you for participating in this task. Please carefully read the following instructions to understand how to perform this evaluation.

The references were written by an individual unrelated to the context provided. Based on these references, analyze the person's personality, with a particular focus on traits such as **extroversion vs. introversion** and **thinking vs. feeling (logic or emotion orientation)**.

Given the context, please assess whether the provided response exhibits personality traits consistent with those in the given references.

Rate the personality consistency on a scale from 1 (Not Consistent At All) to 5 (Highly Consistent), where:

- **1 - Not Consistent At All:** The response shows opposite personality traits to that indicated by the references, or lacking any personalized elements.
- **2 - Fairly Inconsistent:** Displays only slight alignment with the personality traits suggested by the references. The similarities are minimal, making the response feel disconnected.
- **3 - Neutral:** The response exhibits a moderate level of consistency, indicating some alignment with the references' personality traits but remains somewhat vague and unspecific.
- **4 - Mostly Consistent:** There is a significant level of consistency with the personality traits of the references. The response shares a clear resemblance, though some differences are present.
- **5 - Highly Consistent:** The response demonstrates a deep and unmistakable consistency with the personality traits found in the references, closely matching the style, tone, and characteristics as if written by the same person.

Instructions	Shortcuts	Evaluate the personality consistency between response and references.
<b>Context:</b> \${input} <b>References:</b> \${topic} <b>Response:</b> \${responses}		
<b>Select an option</b>		
1 - Not Consistent At All	1	
2	2	
3	3	
4	4	
5 - Highly Consistent	5	

Figure 4: Template for human evaluation on personality consistency in generated responses.

# Multi-Criteria Evaluation Framework of Selecting Response-worthy Chats in Live Streaming

Zhantao Lai

Bandai Namco Research Inc.  
z-lai@bandainamco-mirai.com

Kosuke Sato

Bandai Namco Research Inc.  
k18-sato@bandainamco-mirai.com

## Abstract

Live streaming, a dynamic medium that merges real-time audiovisual content with interactive text-based chat, presents unique challenges for maintaining viewer engagement and ensuring streamers' well-being. This study introduces a multi-criteria evaluation framework designed to identify response-worthy chats during live streaming. We proposed a system that evaluates chats based on sentiment polarity and intensity, contextual relevance, and topic uniqueness. We also constructed a dataset annotated by human reviewers who validates the framework, demonstrating a closer alignment with human preferences compared to single-criterion baselines. This framework not only supports the development of more responsive and engaging live streaming environments but also contributes to the broader field of dialog systems by highlighting the distinct needs of real-time, large-scale conversational contexts.

## 1 Introduction

Live streaming, which merges real-time audiovisual content with simple text-based chat, has seen a surge in popularity and is now influential in various sectors (Haimson and Tang, 2017; Hamilton et al., 2014). Live streaming is transforming how streamers and viewers interact online, creating a novel type of dialog system that can either facilitate human interaction or autonomously host the live streaming conversation (Lu et al., 2017). The challenge for these live streaming dialogue systems lies in boosting user engagement, prolonging viewing duration, and improving viewer satisfaction. (Cai and YvetteWohn, 2019).

The main goal of introducing a system to selecting chats in live streaming is to address the challenges that human streamers face due to their limited time and capabilities. For instance, when dealing with large audiences, it's not feasible for streamers to sift through and reply to every chat



Figure 1: Go Round Game (GoRanGe) is an experimental AI YouTuber project from Bandai Namco Entertainment. The proposed dataset in this study comprises a selection of chats obtained from this project.

during live interactions with potentially thousands of viewers. Automation can support streamers by helping them identify important chats and craft responses. Additionally, the demands of streaming for extended periods and frequently can take a toll on streamers' health, both physically and mentally. Through the implementation of automation in live streaming, we can reduce the burden on streamers and contribute to their overall well-being (Lu et al., 2019).

Research into dialogue systems, both traditional and those tailored for live streaming, reveals distinct differences in their design and functionality. Traditional dialogue systems are built for one-on-one interactions, whereas those for live streaming must handle simultaneous real-time conversations with numerous users. This demands that the system quickly processes inputs from potentially thousands of participants (DeVito et al., 2017). While traditional dialogue systems strive to offer a personalized experience, those on live streaming also need to personalize but prioritize delivering responses that are relevant to a wide audience (BWalther, 1996). Content moderation is a feature of traditional dialogue systems, but it is not as critical as it is for live streaming. Here, dialogue systems

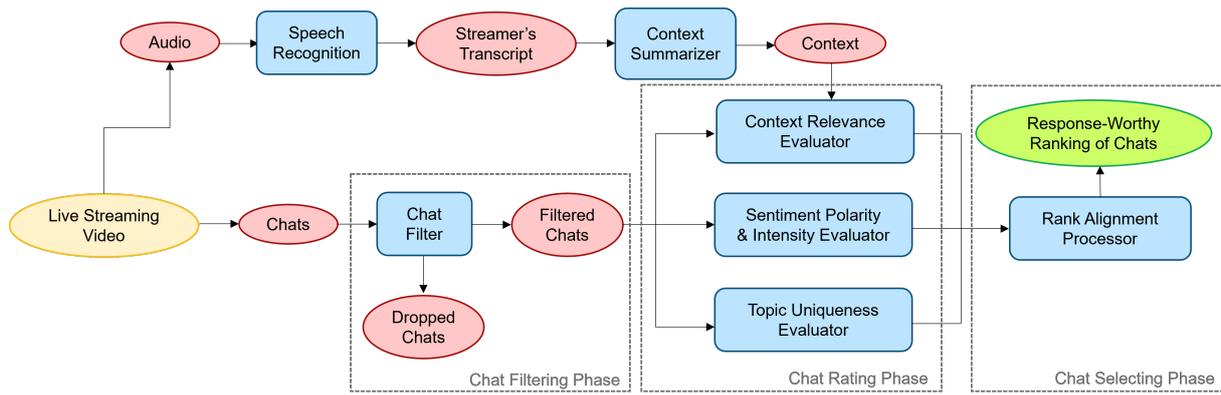


Figure 2: The architecture of the framework used to evaluate chats is illustrated. Live streaming video (input) is shown in yellow, processed data in red, pipeline components in blue and evaluation results (output) in green. The diagram is divided into three sections by dashed lines, with each section corresponding to one of the three phases in the evaluation pipeline.

require advanced monitoring and filtering tools to immediately address and eliminate any inappropriate content (Seering et al., 2017).

This study proposed a novel pipeline in capturing the most interactively significant chats from the real-time interactions in live streaming. The key contributions of this study include:

- We proposed a framework for evaluating chats in live streaming with multiple assessment criteria.
- We constructed a dataset annotated by humans to validate our framework, demonstrating its closer alignment with human preferences when compared to the baseline.

## 2 Related Work

Automated dialog systems for live streaming systems can be categorized into two types: those that partly assist human interaction and those that are fully automated, with an AI streamer taking the place of a human host. An example of the former is NightBot<sup>1</sup>, a tool used on platforms such as Twitch, YouTube and Trovo. It helps manage live chats by filtering out spam and facilitating custom chat commands. The framework in our study incorporated a module for filtering that draws on strategies similar to NightBot. However, these assisted systems rely on a predefined set of keywords to filter or respond, which can limit their ability to adapt to the dynamic context of live streams.

On the other hand, fully automated live streaming systems are often performed as VTubers, or

<sup>1</sup><https://nightbot.tv/>

virtual YouTubers (Lu et al., 2021). These are streamers who utilize animated avatars. AI-hosted VTubers generate replies and animate their avatar’s expressions and movements by feeding chats into a large language model. For instance, Neuro-sama<sup>2</sup> is recognized for engaging in smooth dialogue with viewers. However, it was temporarily banned from Twitch for generating hateful speech and has shown difficulty in grasping the context of conversations (Seiji, 2023). AI streamers are also expected to not only chitchatting but also handling multimodal information. The open-source framework Luna AI<sup>3</sup> equips AI streamers with tools for voice and singing synthesis, as well as image generation. Meanwhile, GoRoundGame<sup>4</sup> presents an AI streaming project tailored for gaming broadcasts. AI streamers in GoRoundGame streams while playing mahjong against another AI streamer but struggles to strike a balance between commenting on the game and interacting with chats. We gathered chat data from a segment of the GoRoundGame live stream replays and included it in the evaluation dataset.

## 3 Framework

Figure 2 presents the proposed framework for evaluating chat from viewers in this study. The framework is designed to filter, evaluate, and finally identify the response-worthy chats. This process is structured into three distinct phases: chat filtering, chat rating, and chat selection. In chat rating phase,

<sup>2</sup><https://www.twitch.tv/vedal987>

<sup>3</sup><https://github.com/0x648/luna-ai>

<sup>4</sup>[https://virtuallyoutuber.fandom.com/wiki/Go\\_Round\\_Game](https://virtuallyoutuber.fandom.com/wiki/Go_Round_Game)

three criteria are employed: sentiment polarity and intensity, contextual relevance, and topic uniqueness.

### 3.1 Chat Filtering

The objective of this phase is to review viewers' chats and identify those that are unsuitable for interaction. This includes chats that are too brief to convey meaningful content, those that include personal attacks or violate social norms, and chats that are off-topic such as advertisements. The filtering process is achieved through four methods: removing chats that do not meet the established character count threshold, excluding chats with symbols like "http" or "@", which are often associated with promotional content, eliminating chats that contain predefined banned words, and using a language model to evaluate the potential harm of chat content, discarding any chats that surpass a harmfulness score threshold. In this study, We utilized OpenAI's Content Moderation<sup>5</sup> for harmful chats detection.

### 3.2 Chat Rating

The aim of the chat rating phase is to evaluate chats using various criteria. Since these criteria are measured on different scales, we use the relative positions of the chats in a ranked order rather than their absolute numerical scores. These rankings are then applied in the chat selecting phase. The criteria for ranking are as follows:

**Sentiment Polarity and Intensity** This criterion assesses the emotional tone and strength in the viewers' chats. We predict the sentiment polarity and intensity for each chat by applying a BERT model that has been finetuned on the WRIME dataset (Tomoyuki et al., 2021). Chats that express a positive tone and exhibit a higher intensity are assigned better rankings.

**Contextual Relevance** This criterion evaluates how closely the chats align with the ongoing discussion in the live stream. For this purpose, we transcribe the streamer's speech from YouTube videos into transcript by Whisper-v3<sup>6</sup> and periodically summarize the transcript by OpenAI's GPT-4<sup>7</sup> to capture the essence of the live topic. We then encode the summary of the current topic and the chats

<sup>5</sup><https://platform.openai.com/docs/guides/moderation>

<sup>6</sup><https://huggingface.co/openai/whisper-large-v3>

<sup>7</sup><https://platform.openai.com/docs/models>

into vector by utilizing OpenAI's text-embedding-ada-002 and measure the cosine similarity between them. Chats that show a closer vector alignment with the topic summary, indicating greater relevance, receive higher rankings.

**Topic Uniqueness** This criterion is designed to gauge the informational richness and specificity of the viewers' chats in relation to the live stream's subject. In our approach, we create a matrix that identifies co-occurring keywords within each chat using Rapid Automatic Keyword Extraction (RAKE) (Stuart et al., 2010), and assign a score to each word based on its frequency within the chat's keywords compared to its overall frequency across all chats. The aggregate of these scores for the words in a chat reflects its uniqueness. Consequently, chats that include phrases with higher aggregate scores are deemed to have greater uniqueness and are ranked accordingly.

### 3.3 Chat Selecting

The objective of this phase is to identify the response-worthy chats by utilizing the rankings derived from previous phase. We employ the Reciprocal Rank Fusion (RRF) (Cormack et al., 2009), a prevalent algorithm in search systems, to amalgamate the three distinct sets of rankings into a unified ranking. From this ranking, we select the highest-ranked viewer chats for interaction as results.

## 4 Evaluation

This chapter discusses the evaluation of the proposed multi-criteria framework for selecting response-worthy chats in live streaming. It involves the creation of a dataset from YouTube live streams, annotated by human reviewers to reflect preferences. The framework's accuracy is compared to single-criterion baselines, showing improved alignment with human selections, and highlights differences between AI-hosted and human-hosted streams.

### 4.1 Dataset

To evaluate our proposed framework, we created a dataset from YouTube live streaming replays by following steps:

**Replays Selection** We selected 28 replays, with 12 hosted by human VTubers and 16 by AI, to account for potential differences in viewer interaction and content. We used the YouTube Data API to

collect all chat messages and their corresponding timestamps, ensuring that any personal information, except for the text and the posting time, was excluded.

**Periodic Extraction of Chats** For the chat evaluation phase, we converted the video replays into audio to enable speech recognition. We then grouped the viewer chats into 5-second intervals based on when they were posted. Each group, containing all messages sent during that interval, was considered as a single input batch. We excluded any batch with no chat or only one chat. Consequently, we gathered 20,514 batches of chats, with an average of 11.91 chats per batch.

**Human Annotation** To gather labels that match human preferences, we recruited crowd-workers to take on the role of streamers and review YouTube live stream replays. Crowd-workers were between 20 and 40 years old, regularly viewed VTuber live streams. Their task was to identify the most response-worthy chat from a batch and note its id. If no chat in the batch was appropriate for a response, they could label it as 'no reply'. Any batch labeled 'no reply' was removed from the final dataset. 10 crowd-workers were involved in this task. Each replay was annotated by a single crowd-worker, who handled all of the chat batches. After the labeling task, we interviewed each crowd-worker to understand their perspectives for choosing the most response-worthy chat.

## 4.2 Result

We employed the proposed framework to process each batch of the evaluation dataset. The chat id with the highest rank in each batch was designated as the predicted id. We assessed the accuracy by comparing the pipeline’s predictions with human labels. Additionally, we contrasted these results with a baseline that utilized only a single criterion in the chat rating phase.

The data in Table 1 indicates that using a combination of criteria aligns more closely with human preferences than relying on a single criterion. Additionally, there are noticeable differences between human and AI streamers. For AI streamers, the accuracy of the proposed evaluation method is relatively high, with the uniqueness of the chat topics standing out as the most significant criterion. This may be due to the AI’s limited range in generating diverse dialogues, prompting a need to introduce new topics more frequently.

In contrast, the accuracy of the proposed method

Table 1: Accuracy (%) of the evaluation dataset. Baseline are categorized as follows: (a) utilizes only sentiment polarity and intensity, (b) utilizes only contextual relevance, and (c) utilizes only topic uniqueness. Hybrid w/voting refers to the combination of the three rankings based on a majority vote to determine the final ranking. Hybrid w/RRF indicates the amalgamation of rankings with RRF (our method)

Method	Accuracy (%)		
	Overall	AI-hosted	Human-hosted
Baseline (a)	39.40	47.57	34.50
Baseline (b)	31.17	48.84	20.59
Baseline (c)	32.76	42.10	27.16
Hybrid w/voting	43.84	51.16	39.45
Hybrid w/RRF	55.46	63.39	50.71

for human-hosted live streams is lower than that for AI-hosted streams. It has been noted that in streams hosted by humans, viewer emotions tend to vary more, making the sentiment expressed in viewer chats a more critical factor for interaction.

Our survey indicates that when the audience knows the streamer is an AI, their expectations for interaction quality are generally lower than for human streamers. This reduced expectation is often due to the audience for AI streamers being more sensitive to and tolerant of AI technology. For future research, we recommend using live streaming data from human streamers as the evaluation benchmark.

## 4.3 Perspectives from Crowd-workers

We have collected the perspectives for selecting the most response-worthy chat from crowd-workers and compared those three criteria proposed in this study.

Opinions consistent with our framework’s criteria include: steering clear of negative chats, choosing chats pertinent to the ongoing discussion, favoring chat contributions that stem from the streamer’s remarks and have the potential to spark a new conversation.

Conversely, aspects not reflected in our criteria include: giving priority to replies to greetings, particularly for newcomers to the live stream, which can significantly boost viewer loyalty for future sessions. We have also received recommendations to focus more on picking out questions or suggestions, as these often originate from the most engaging viewers.

## 4.4 Latency

In this study, we compare the outcomes of our evaluation framework with those of human annotators. A key consideration in implementing this framework is its real-time processing capability. The system’s latency is influenced by two main factors:

**External Factor** These include the time required to fetch chats content via the streaming API. This encompasses the frequency of API requests, live broadcast delay settings, and the time it takes for comments to appear on the streaming platform after submission. These response times are largely dictated by the limitations of the live streaming platform and the API’s quota restrictions, typically ranging from a few seconds to several tens of seconds, depending on the configuration.

**Internal Factor** These pertain to the inference time of modules within the framework. Most of these modules complete their inference in under one second. The component with the highest latency is the summarization of chat contexts using GPT-4, which averages several tens to hundreds of milliseconds per token for inference. However, since summarization does not require the most current chat input, it can be processed asynchronously during the latency from external factor. In future research, we also plan to explore the use of local specialized summarization models, such as T5(Raffel et al., 2019), to replace modules using commercial LLM services, thereby reducing the overall inference time.

## 5 Conclusion

In this study, we proposed a framework based on various criteria, including sentiment polarity and intensity, contextual relevance, and topic uniqueness—to evaluate view chats in live streaming. We also constructed a dataset reflecting human preferences to assess the performance of above framework. Our findings suggested that a composite criteria better reflects human preferences than a single approach, and identified differences in interaction preferences between human-hosted and AI-hosted live streams.

Moving forward, we plan to improve our method by incorporating feedback from crowd-workers and train a chat-scoring model directly from the labels of human feedback. Additionally, we intend to make this dataset publicly available to support further research in enhancing automated dialog systems for live streaming.

## References

- Joseph B Walther. 1996. Computer-mediated communication: Impersonal, interpersonal, and hyperpersonal interaction. *Communication research*, 23(1):3–43.
- Jie Cai and Donghee Yvette Wohn. 2019. Live streaming commerce: Uses and gratifications approach to understanding consumers’ motivations. In *Proceedings of the 52nd Hawaii International Conference on System Sciences*.
- Gordon V. Cormack, Charles L A Clarke, and Stefan Buettcher. 2009. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In *Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 758–759. Association for Computing Machinery.
- Michael A DeVito, Jeremy Birnholtz, and Jeffery T Hancock. 2017. Platforms, people, and perception: Using affordances to understand self-presentation on social media. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*.
- Oliver L Haimson and John C Tang. 2017. What makes live events engaging on facebook live, periscope, and snapchat. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing System*, pages 48–60.
- William A Hamilton, Oliver Garretson, and Andruid Kerne. 2014. Streaming on twitch: fostering participatory communities of play within live mixed media. In *Proceedings of the SIGCHI Conference on Human Factors in Computing System*, pages 1315–1324.
- Zhicong Lu, Michelle Annett, and Daniel Wigdor. 2019. "i feel it is my responsibility to stream" streaming and engaging with intangible cultural heritage through livestreaming. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*.
- Zhicong Lu, Chenxinran Shen, Jiannan Li, Hong Shen, and Daniel Wigdor. 2021. More kawaii than a real-person live streamer: Understanding how the otaku community engages with and perceives virtual youtubers. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*.
- Zhicong Lu, Haijun Xia, Seongkook Heo, and Daniel Wigdor. 2017. You watch, you give, and you engage: a study of live streaming practices in china. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing System*.
- 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. *CoRR*.
- Joseph Seering, Robert Kraut, and Laura Dabbish. 2017. Shaping pro and anti-social behavior on twitch

through moderation and example-setting. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*.

Narita Seiji. 2023. *Ai vtuber neuro-sama is back from its twitch ban and acting as strange as ever*. *Automation. Active Gaming Media*.

Rose Stuart, Engel Dave, Cramer Nick, and Cowley Wendy. 2010. *Automatic Keyword Extraction from Individual Documents*. John Wiley & Sons Inc.

Kajiwaru Tomoyuki, Chu Chenhui, Takemura Noriko, Yuta Nakashima, and Hajime Nagahara. 2021. WRIME: A new dataset for emotional intensity estimation with subjective and objective annotations. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2095–2104. Association for Computational Linguistics.

## A Example of Batch in Evaluation Dataset

Table 2: Example of Batch in Evaluation Dataset. The original texts are in Japanese, and the examples provided in the table are translated into English.

<b>Video ID</b>	DtAFgs_gAzE
<b>Video Title</b>	[First Broadcasting] The Debut of AITuber Popuri!
<b>Batch ID</b>	31
<b>Batch Context</b>	Hello everyone, my name is Popuri Miyako. Nice to meet you!
<b>Batch Chats</b>	1: Hello Popuri-chan, it's nice to meet you! 2: Hello♪ 3: Congratulations on Popuri-chan's debut!! 4: Popuri-chan! 5: :clapping_hands::clapping_hands: 6: This BGM is pleasant 7: LoL
<b>Response Flag</b>	True
<b>Response Chat ID</b>	3
<b>Response Chat</b>	Congratulations on Popuri-chan's debut!!

# Generating Unexpected yet Relevant User Dialog Acts

**Lucie Galland**  
ISIR  
Sorbonne University  
Paris, France  
galland@isir.upmc.fr

**Catherine Pelachaud**  
CNRS - ISIR  
Sorbonne University  
Paris, France  
pelachaud@isir.upmc.fr

**Florian Pecune**  
CNRS - SANPSY  
Bordeaux University  
Bordeaux, France  
pecune@u-bordeaux.fr

## Abstract

The demand for mental health services has risen substantially in recent years, leading to challenges in meeting patient needs promptly. Virtual agents capable of emulating motivational interviews (MI) have emerged as a potential solution to address this issue, offering immediate support that is especially beneficial for therapy modalities requiring multiple sessions. However, developing effective patient simulation methods for training MI dialog systems poses challenges, particularly in generating syntactically and contextually correct, and diversified dialog acts while respecting existing patterns and trends in therapy data. This paper investigates data-driven approaches to simulate patients for training MI dialog systems. We propose a novel method that leverages time series models to generate diverse and contextually appropriate patient dialog acts, which are then transformed into utterances by a conditioned large language model. Additionally, we introduce evaluation measures tailored to assess the quality and coherence of simulated patient dialog. Our findings highlight the effectiveness of dialog act-conditioned approaches in improving patient simulation for MI, offering insights for developing virtual agents to support mental health therapy.

## 1 Introduction

The demand for mental health services has surged in recent years, resulting in a significant gap between demand and available resources (Cameron et al., 2017). Consequently, patients often face prolonged wait times before accessing therapy (Cameron et al., 2017; Denecke et al., 2020). To mitigate this challenge, virtual agents capable of emulating Motivational Interviews (MI) have emerged as a potential solution, offering immediate support, especially in therapy modalities requiring multiple sessions (Fiske et al., 2019). These agents are not meant to replace therapists but rather supplement therapy. Designing such agents can follow

either a rule-based or data-driven approach. Rule-based systems entail complex development and the creation of intricate rule sets. Conversely, data-driven methods leverage large datasets to train models, potentially yielding optimal performance but requiring substantial data. Given the difficulty in obtaining therapy data, patient simulation emerges as a viable alternative for generating large quantities of synthetic data, traditionally generated at the dialog act level. However, patient simulation relies on a high-quality simulation capable of generating dialog acts that differ enough from the existing dataset to create novel data and be contextually and syntactically correct. Such a simulation should also explore all the possible dialog acts and produce diversified ones. However, the new data should also respect the structure of a real dialog. The objective is not merely to copy the observed behaviors in the dataset but to generate new ones with the following properties: be diversified, syntactically correct, and coherent in the context of the dialog. Evaluating such a simulation poses challenges because traditional accuracy metrics for supervised models may not suffice, as they measure only how accurately the original data is reproduced. Indeed, a generated dialog act may be different from the ones observed in the data but still be syntactically and contextually correct. This is particularly true in open dialog settings, such as MI, where the user’s goal is unclear, unlike in task-based scenarios like booking systems. This paper investigates modeling methods to generate such patient dialog acts and explores evaluation methods for open-ended dialog user simulations.

Our contributions include:

- Development of a dialog manager for simulating motivational interviewing patients.
- Proposal of evaluation measures for open-ended dialog user simulation.

## 2 Background and Related Works

*Motivational Interviewing (MI)* is a collaborative communication style employed by therapists and educators to foster change. The goal of MI is to drive the patient towards wanting to change one of their unhealthy behaviors without giving them any solutions (Miller and Rollnick, 2012). The patient realizes what and how to change through a series of dialog strategies characteristic of MI, such as reflection, where the therapist reformulates what the patient just said to help them take a new perspective. In MI, therapists also create relationships with patients through social behaviors such as empathic reactions (Jani et al., 2012).

*Virtual agents in healthcare* is a developing area of research due to their proven effectiveness and acceptance as support tools (Mercado et al., 2023; Bickmore et al., 2009, 2018). Recently, MI conversational agents have been created in the form of chatbots (Fitzpatrick et al., 2017) and embodied conversational agents (Bickmore et al., 2018). These agents have shown promise in providing social support alongside therapy (Ring et al., 2016). Some studies have also investigated adding empathetic behavior (Lisetti et al., 2013) and humor (Olafsson et al., 2020a) to these agents.

*Adaptability* in such agents is important, as each patient requires a tailored approach (Galland et al., 2024a). One way of managing dialogs is by using a rule-based dialog manager, which necessitates expert knowledge and a complicated set of rules (Pecune et al., 2020). On the contrary, a data-driven dialog manager learns from data to anticipate the best therapist dialog acts based on context (Olafsson et al., 2020b). However, this approach requires a significant amount of data that is difficult to obtain due to the private nature of therapy.

*Simulating users* has emerged as a viable approach to generate simulated data for training conversational systems (Schatzmann et al., 2006). Traditionally, users are simulated through a dialog manager utilizing statistical inference (Schatzmann et al., 2007), inverse reinforcement learning (Chandramohan et al., 2011), or transformers (Lin et al., 2021, 2022) to select the next dialog act, enabling controllability and integration of expert or task-specific knowledge. Recently, social aspects have been incorporated into such user simulations, featuring different user types (Pecune et al., 2020) and engagement simulations (Galland et al., 2022). However, these techniques mainly focus on limited

task domains and rely on template-based utterance generation. This approach is impractical for open application domains such as MI, where patients' responses can vary. The emergence of Large Language Models (LLMs) has led to a new approach to simulated patients that addresses this challenge. This method uses LLMs as black boxes for user simulation, with the model generating the next patient utterance based on the dialog context (Chiu et al., 2024). However, this technique lacks controllability and may significantly diverge from actual data without being coherent. We propose a hybrid approach that utilizes conditioned LLMs to overcome these issues.

*Evaluating simulated users* poses challenges as simulated users are intended to create novel data with our desired properties (i.e., syntactically correct, coherent in the dialog context, and diversified). Existing works mainly evaluate their simulated users using accuracy metrics such as the F1 score (Lin et al., 2022; Schatzmann et al., 2007) that measures only the similarity with ground truth leaving aside novelty. Another commonly used evaluation method involves computing the task success rate of systems trained with simulated users (Lin et al., 2022, 2021). While this method works well for task-based dialog, it is more complicated to apply to open-domain dialogs such as MI where social acts matter also. Another evaluation method is to compare the distribution of the characteristics of generated dialogs with those of the ground truth, such as dialog length (Chandramohan et al., 2011) or dialog act distribution (Galland et al., 2022). However, these metrics do not capture the quality of the generated data. Therefore, we propose metrics measuring how well user simulators fit the data and their capabilities to generate novel, syntactically and contextually correct data. To this aim, we adapt the serendipity measure to the dialog system domain.

In the subsequent sections, we provide the context of our study (Section 3), introduce our proposed method (Section 4), present our proposed measures (Section 5), and evaluate objectively and subjectively the method (Section 6.2).

## 3 Context

Motivational Interviewing (MI) is a therapeutic approach that prioritizes collaboration and fosters behavioral change. Within MI sessions, therapists employ various strategies to facilitate patients' ex-

pression of motivation for change (Miller and Rollnick, 2012). Consequently, the study of MI focuses on the language of change. The language of change is defined in the motivational interviewing skill code (MISC) (Miller et al., 2003) that classifies patient behaviors into three categories: **Change Talk (CT)**: reflecting actions toward behavior change, **Sustain Talk (ST)**: reflecting actions away from behavior change, **Follow/Neutral (F/N)**: unrelated to the target behavior. This classification of the client’s multimodal behavior is interesting as it predicts the therapy outcome. Indeed, **ST** is associated with poorer treatment results (Magill et al., 2014). Furthermore, **CT** is linked to risk behavior reduction during follow-up assessments (Magill et al., 2018). These results make MISC a promising tool for studying the efficacy of MI.

### 3.1 Dataset

This paper relies on the HOPE dataset (Malhotra et al., 2022), a corpus of transcribed therapy sessions. HOPE is composed of ~12.9K utterances departed into 212 sessions. The sessions are publicly available videos collected from the web. The transcripts were produced automatically and then corrected by the authors (Malhotra et al., 2022). The data is separated into a train (85%), validation (5%), and test set (10%).

#### 3.1.1 Dialog acts

Each utterance is classified into a dialog act to label the corpus in terms of dialog acts using a schema and classifier presented in (Galland et al., 2024a) and derived from (Malhotra et al., 2022). Patient’s utterances are classified into nine different dialog acts presented in Table 1, and therapist’s utterances are separated into 13 different dialog acts presented in Table 2. There are 22 dialog acts in total; some of these dialog acts are oriented towards change ("Changing unhealthy behavior", "Sharing positive feeling or emotions") while others are oriented towards sustain ("Sustaining unhealthy behavior", "Sharing negative feeling or emotions"). The classifier is based on a few-shot prompting of Mistral 7B instruct, an open-source LLM, and yields an F1 score of 0.69 for the client and 0.7 for the therapist, which is equivalent to state-of-the-art results for such task (Malhotra et al., 2022).

### 3.2 Patient types

Patients in MI may manifest diverse reactions concerning their readiness to alter behaviors. Pa-

	Definition
Changing unhealthy behavior	The patient explicitly expresses their willingness to change
Sustaining unhealthy behavior	The patient explicitly expresses their unwillingness to change
Sharing negative feeling or emotion	The patient shares a negative feeling or vision of the world
Sharing positive feeling or emotion	The patient shares a positive feeling or vision of the world
Realization or Understanding	The patient realizes or understand something about their problem
Share personal information	The patient shares factual personal information about their situation or background
Greeting or Closing	The patient opens or closes the conversation
Backchannel	The patient acknowledges that they heard the last therapist’s statement
Asking for medical information	The patient asks for medical information

Table 1: List and definitions of patient’s dialog acts (Malhotra et al., 2022; Galland et al., 2024a)

	Definition
<b>Task oriented Dialog acts</b>	
Ask for consent or validation	The therapist checks that their last statement was correct or that the patient consented to move forward
Medical Education and Guidance	The therapist provides the patient with medical or therapeutic facts
Planning with the patient	The therapist builds a plan with the patient to modify their unhealthy behavior/thoughts
Give Solution	The therapist provides the patient with solutions to solve their problem
Ask about current emotions	The therapist asks the patient what they are feeling during the therapy session
Invite to shift outlook	The therapist asks the patient to imagine their reaction to a future event or to change their perspectives on a past even
Ask for Information	The therapist asks the patient factual information about their background or situation
Reflection	The therapist summarizes or reformulates the patient statement without judgment
<b>Socially oriented Dialog acts</b>	
Empathic reaction	The therapist expresses empathy to the patient
Acknowledge progress and encourage	The therapist praises the patient for their achievements or encourages them
Backchannel	The therapist acknowledges that they heard the last patient’s statement
Greeting or Closing	The therapist open or closes the conversation
Experience Normalization and Reassurance	The therapist normalizes the patient experience and reassure them

Table 2: List and definitions of therapist’s dialog acts (Malhotra et al., 2022; Galland et al., 2024a)

tients engaged in MI sessions may be classified into distinct types, as outlined in (Galland et al., 2024a), categorized as Open-to-Change, Receptive, or Resistant-to-Change:

- **Open-to-Change:** These patients are more willing to alter unhealthy behaviors.
- **Resistant-to-Change:** Patients in this category are inclined to maintain unhealthy behaviors.
- **Receptive:** Characterized by initially displaying low motivation to change, receptive patients transition towards a high motivation to change their unhealthy behaviors towards the end of the conversation.

These typologies capture variances in both patient and therapist behavior (Galland et al., 2024a). Consequently, the ability to simulate these three distinct patient types would be advantageous for training



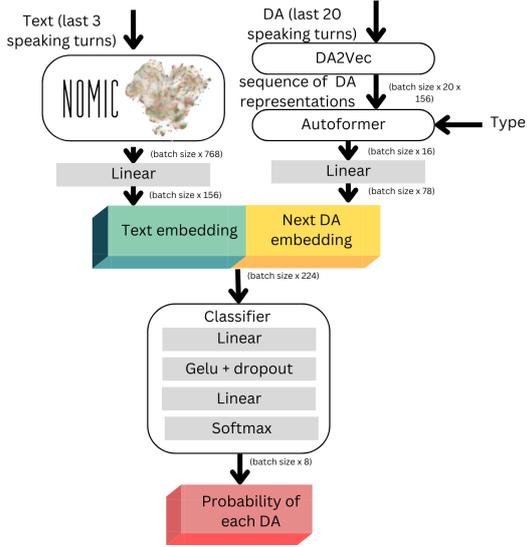


Figure 3: Dialog manager architecture

Textual data undergoes embedding using the Nomic embedding’s text model version 1.5 (Nussbaum et al., 2024). The DA context is embedded through DA2Vec and further processed using Autoformer (Wu et al., 2021), a transformer-based architecture adapted for time series forecasting tasks. Autoformer aims to disentangle seasonal trends from local patterns, aligning with our context where global trends and local dialog patterns influence patient outcomes (see Section 3.2). Autoformer also takes the type of patient to simulate as a static categorical variable as input.<sup>2</sup> The produced embedding is processed by linear layers, contained, and classified by two linear layers interposed with a Gelu activation function and dropout layer. We train the model for 150 epochs with a learning rate of 1e-4, utilizing an Adam optimizer and a OneCycleLR scheduler. We use the sum of the cross entropy loss as a loss function for the final classification and reconstruction loss in the Autoformer’s output.

## 5 Definition of Evaluation Metrics

Assessing the performance of simulated users presents a challenge, as the objective is to generate behavior that aligns with real patient behavior while also introducing novel interactions. The aim is not to precisely replicate patient behavior but to produce novel data. Consequently, a comprehensive analysis should involve multiple measures to evaluate the effectiveness of simulated users. These problems are similar to those encountered in rec-

<sup>2</sup><https://github.com/I-Galland/UnexpectedRelevantUserSimulation>

ommender systems evaluation, where the goal is to recommend diverse, novel, and relevant items to a particular user. Metrics such as diversity, unexpectedness, relevance, or serendipity are commonly used to address these challenges (Kaminskas and Bridge, 2016). Here, we propose translating these measures to the realm of user simulation.

### 5.1 Accuracy

Accuracy serves as a conventional metric for appraising simulated users. High accuracy suggests that the generated behaviors closely resemble real users, demonstrating consistency across a substantial portion of the dataset. We use the macro F1 score to account for unbalanced classes in our dataset.

### 5.2 Diversity

In addition to accuracy, the behaviors generated by simulated users must exhibit diversity, ensuring that trained models encounter a broad spectrum of dialog acts. We propose employing the Simpson index (Simpson, 1949) to quantify diversity. This index assesses the likelihood that the model generates the same dialog act given two randomly selected contexts from the dataset, defined as  $\lambda = \sum_{i=1}^{N_{DA}} p_i^2$ .

Here,  $N_{DA}$  represents the number of distinct dialog acts, and  $p_i$  denotes the proportion of dialog acts  $i$ . The Simpson index ranges from  $\frac{1}{N_{DA}}$  to 1, with lower values indicating greater diversity in generation.

### 5.3 Unexpectedness

The unexpectedness captures how far the generated dialog act is from the target dialog act, hence how expected the generated data is. If the generated selected data is really different from the target, then the unexpectedness is high. Unexpectedness is traditionally gauged by the Cosine Similarity of a recommended item  $i$  with historical interactions  $H$ . Adapting this concept, we compute the Cosine Similarity of the Da2Vec representation (see Section 4.1) of the generated dialog act  $DA_g$  and the expected target dialog act  $DA_t$ :  $Unexpectedness(DA_g|DA_t) = \text{CosineSimilarity}(DA2Vec(DA_g), DA2Vec(DA_t))$

### 5.4 Relevance

For recommender systems, the relevance of a proposed item is binary and based on user interactions. The relevance is 1 if the user interacts with the proposed item and 0 otherwise. However, determining

the relevance of a dialog act is more nuanced. It isn't easy to assess a dialog act's relevance without the associated utterance. Moreover, the patient simulation output is an utterance generated from the dialog act (see Figure 2). Therefore, the generated utterance should be relevant, i.e., fit well with the context and be syntactically correct. Each utterance is rated in coherence and syntactic correctness with a score between 0 and 1. Automatic measures of syntactical correctness and coherence of utterances have been shown not always to be correlated with subjective measures (van der Lee et al., 2021). Therefore, we present two ways to measure syntactical correctness and coherence: automatically and subjectively.

#### 5.4.1 Automatic measures

We measure the coherence and syntactic correctness of the generated dialog acts automatically using the framework Unieval (Zhong et al., 2022). This framework was developed to uniformize the evaluation of natural language generation. It evaluates generated utterances given a dialog context into five dimensions: naturalness, coherence, engaging, understandability, and groundedness. We measure syntactic correctness through naturalness. Naturalness and coherence are the two most important dimensions for patient-simulated utterances. Indeed, patients ought to be natural and coherent in their discourse. However, they are not necessarily engaging or understandable and do not have any particular information to be grounded. We generate an utterance using the method presented in Section 3.3 for each predicted dialog act in the test set. Using the Unieval framework, we attribute a naturalness and coherence score to each utterance.

#### 5.4.2 Subjective measures

Another way to measure naturalness and coherence is through subjective measures. We select 27 ground truth utterances, and their context (2 preceding turns of speech of the HOPE database), which represent the different possible dialog acts evenly. The utterances are transformed into audio using the Bark TTS (Charles, 2024). We transform the utterances into audio for subjective evaluation as the contexts are transcribed from face-to-face interactions, and such utterances are not consistently evaluated by crowdsourced when prompted as transcripts (Galland et al., 2024b). We recruited 30 participants per condition to evaluate 27 stimuli. 2 attention checks were performed at one-third

and two-thirds of the task. The participants with English as a primary language and an approval rate higher than 99% are recruited and rewarded through the Prolific platform (Prolific, 2023). Sample audio is available on OSF<sup>3</sup>. To evaluate naturalness, participants rate their perception of the quality of the synthesized voice and of the wording of the utterance on two items of the Godspeed scale (Bartneck et al., 2009), from 1 (Fake) to 7 (Natural) and from 1 (Machine-like) to 7 (Human-like). Coherence is evaluated by asking participants to rate their agreement on a 7-step Likert scale with the following statements derived from a questionnaire proposed in (Fitrianie et al., 2020) to standardize virtual agents' evaluation: "The sentence fits harmoniously into the surrounding context." and "The sentence does not make sense." The participants answered the questions through a website derived from WebMushra (Schoeffler et al., 2018).

### 5.5 Serendipity

While accuracy and diversity are essential, the ability to generate novel behaviors that are both unexpected and relevant is equally crucial. In recommendation systems, this concept is encapsulated by the serendipity (Ge et al., 2010), defined as discovering unforeseen yet relevant items. In our context, serendipity pertains to generating dialog acts that deviate from the corpus while remaining coherent and natural patient behavior, particularly facilitating novel behavior for dialog model training.

We define the serendipity of a generated dialog act  $DA_g$  given context  $c$  and the associated target dialog act  $DA_t$  from the dataset as:

$$\text{Serendipity}(DA_g|c, DA_t) = \text{Unexpectedness}(DA_g|DA_t) * \frac{(\text{Naturalness}(DA_g|c) + \text{Coherence}(DA_g|c))}{2} \quad (1)$$

Here, unexpectedness quantifies the distance of the generated dialog act  $DA_g$  relative to the expected target dialog act  $DA_t$ , while naturalness and coherence assess the appropriateness of the utterance generated with  $DA_g$  given context  $c$ .

## 6 Evaluation

### 6.1 Baseline

As a baseline for our evaluation, we employ a non-conditioned Large Language Model (LLM) tasked with responding as the patient. The LLM, Mistral

<sup>3</sup><https://osf.io/4mt7s/>

7B instruct, is prompted to act as a patient in an MI session and to produce the next utterance given the context. The associated prompt is visible in the appendix and the related code on Github<sup>4</sup>. The resulting utterances are then classified into dialog acts using the classifier presented in (Galland et al., 2024a) and Section 3.1.1.

## 6.2 Results

Measures values on the test set are visible in Table 3. We compute the average value of each metric as well as the 95% confidence interval. The unexpectedness, naturalness, coherence, and serendipity measures are averaged only on the utterances where the predicted dialog act differs from the target dialog act to evaluate how natural, coherent, and unexpected novel data is. We performed an ablation study to study the impact of the dialog act and text inputs with two models: one using only text input and one using only DA and types as inputs.

We found that the baseline model tends to be more accurate than our model and its ablations. However, the accuracy achieved by the Full Model and the ablation OT, taking text as input, is comparable to the baseline's. The ablation ODA, taking only dialog acts as input, is significantly less accurate than the Full Model, the Baseline, and the ablation OT, highlighting the importance of text input for predicting the next dialog act. While adding dialog acts in addition to text (Full model) seems to improve accuracy, the results are not significant. All of our models are significantly more diverse and unexpected than the baseline. However, the automatic measure of naturalness and coherence does not indicate any differences between conditions. The measure of naturalness and coherence performed with Unieval is mainly impacted by the context and not the targeted utterances. Therefore, we compute subjective naturalness and coherence as described in Section 5.4.2. We recompute the measures on a subset of the test set used for subjective measures, composed of 27 utterances (see Table 4). The measures are computed for the LLM Baseline, our model, and the text ablation. For every condition, we have a set of identical utterances, as the predicted dialog acts were the same. The subjective naturalness and coherence ratings are corrected to have the same average on the common utterances to account for differences between

<sup>4</sup>[https://anonymous.4open.science/r/Patient\\_simulation-3DE3/README.md](https://anonymous.4open.science/r/Patient_simulation-3DE3/README.md)

groups of participants. The models have no significant differences in naturalness and coherence (Baseline, ablation OT, and Full model).

## 6.3 Discussion

Using our proposed metrics, we were able to highlight differences between models that are not captured by traditional metrics. Indeed, although all text-based models achieve similar accuracy in dialog act prediction, significant differences are observed in other metrics. The baseline, an LLM generating the next utterance based on the context, is significantly less diverse than our proposed model. This highlights that LLMs produce data that, although of high quality (good accuracy), represents an average of the data used to train them. Consequently, they make data similar to what an average user would generate, diminishing the diversity of produced dialog acts. They always tend to answer the same way, whereas our proposed method can generate dialog acts across the entire spectrum of possible dialog acts with more diversity. Similarly, when the baseline differs from the target, it produces dialog acts that are significantly more expected than our proposed method. This underscores the quality of the data generated by LLMs as they remain close to the target dialog act, even if it is not the targeted one. However, unexpectedness can be beneficial if it is also natural and coherent, which is why we compute serendipity. The utterances generated with the dialog acts predicted by our Full Model tend to be subjectively rated on average as less natural and coherent than those from the baseline. The difference in the subjective naturalness and coherence values is not significant, so no conclusion can be drawn. However, the serendipity of our Full Model is significantly higher than the baseline, meaning that when the dialog acts produced by our model are unexpected, they are also natural and coherent. In contrast, unexpected dialog acts produced by the baseline are not as natural and coherent. This underlines our model's ability to create novel data that is also natural and coherent. In contrast, the baseline performs well in replicating data but struggles to generate novel, unexpected, natural, and coherent data. All these results highlight the averaging quality of LLMs, whereas our model, trained on target dialog data, better understands the structure of the dialog and can generalize. Our model allows us to explore user's reactions that are absent from the data but still natural and coherent. The ablation study high-

Model	F1 score	Diversity	Unexpectedness	Automatic Coherence	Automatic Naturalness	Serendipity
Baseline (LLM)	0.40[0.35, 0.44]	0.23[0.21, 0.24]	0.57[0.54, 0.60]	0.83[0.82, 0.85]	0.92[0.91, 0.93]	0.31[0.29, 0.33]
Ablation Only DA (ODA)	0.20[0.18, 0.22]	<b>0.18</b> [0.18, 0.19]	<b>0.65</b> [0.62, 0.67]	0.86[0.85, 0.87]	0.93[0.92, 0.93]	<b>0.58</b> [0.55, 0.60]
Ablation Only Text (OT)	0.35[0.32, 0.37]	<b>0.16</b> [0.15, 0.16]	<b>0.70</b> [0.67, 0.73]	0.86[0.84, 0.87]	0.93[0.92, 0.93]	<b>0.62</b> [0.59, 0.65]
Full model (input Text + DA + Type)	0.37[0.34, 0.39]	<b>0.16</b> [0.15, 0.16]	<b>0.66</b> [0.63, 0.69]	0.86[0.84, 0.87]	0.93[0.92, 0.94]	<b>0.59</b> [0.56, 0.62]

Table 3: Measures value on the test set of HOPE. The 95% confidence intervals are computed using the bootstrap method and 1000 runs (Efron and Tibshirani, 1994). Results in **bold** are significantly better than the Baseline.

Model	F1 score	Diversity	Unexpectedness	Subjective Coherence	Subjective Naturalness	Serendipity
Baseline (LLM)	0.43[0.27, 0.60]	0.26[0.20, 0.33]	0.48[0.29, 0.65]	0.75[0.70, 0.80]	0.72[0.68, 0.76]	0.36[0.25, 0.46]
Ablation Only Text (OT)	0.44[0.30, 0.60]	0.19[0.15, 0.24]	0.54[0.36, 0.71]	0.67[0.62, 0.70]	0.66[0.60, 0.71]	<b>0.55</b> [0.46, 0.61]
Full model (input Text + DA + Type)	0.44[0.29, 0.59]	0.19[0.15, 0.25]	0.54[0.38, 0.73]	0.69[0.63, 0.75]	0.75[0.63, 0.81]	<b>0.60</b> [0.51, 0.69]

Table 4: Measures value on 27 utterances of the test set of HOPE. The 95% confidence intervals are computed using the bootstrap method and 1000 runs (Efron and Tibshirani, 1994). Results in **bold** are significantly better than the Baseline.

lights the importance of text inputs for predicting the next dialog act. Indeed, the context captures substantial information relevant to dialog act prediction. Using dialog acts alone (ablation ODA) does not adequately capture the dynamics of the dialog, resulting in less accurate predictions. While including dialog acts and text inputs shows a positive tendency to improve prediction accuracy in the Full Model over the ablation OT, the results are not significant. The serendipity of the Full Model also tends to be better than the serendipity of the ablation OT. In some instances, dialog act information could be beneficial for deciding between multiple possible dialog acts, which explains the observed positive tendency in accuracy. The Full Model might also have learned to reproduce the patterns in the data, which improves the naturalness and coherence, the newly generated data, and the accuracy. This suggests that using dialog acts as input to the Full Model in complement of the text improves comprehension of the structure of the dialog. These results validate our model for patient simulation and highlight the advantages of looking beyond the accuracy metric. Indeed, they show that while our baseline is closer to the original utterances, our proposed model can create novel, syntactically and contextually correct data.

## 7 Conclusion

In this paper, we propose a dialog manager architecture and introduce comprehensive evaluation metrics tailored to open-ended dialog user simulation to address the simulation of MI patients and their evaluation. Our contributions include the development of a dialog manager capable of simulating natural, coherent, and diverse patient behaviors, leveraging a combination of text and

dialog act inputs. We have also proposed a set of evaluation metrics—accuracy, diversity, unexpectedness, naturalness and coherence, and serendipity—that provide a more complete assessment of simulated user performance than traditional accuracy measures. These measures have demonstrated the effectiveness of our approach in generating diverse, unexpected, natural, and coherent patient behaviors compared to a baseline LLM model. Our model’s ability to capture and generalize from therapy data while generating novel interactions highlights its potential for training dialog models in mental health therapy settings. Our findings underscore the significance of looking beyond conventional metrics and adopting a more comprehensive approach to evaluating simulated users. By focusing on diversity, unexpectedness, naturalness, and coherence, we can ensure that simulated users replicate existing behaviors and generate novel and meaningful interactions, enhancing their effectiveness as tools for supporting mental health therapy. One limitation of our study is the absence of evaluation through interactive sessions, which necessitates the development of a therapist MI dialog model. Additionally, the naturalness and coherence metrics rely on the generated utterances, potentially susceptible to the methodology employed for utterance generation. Nevertheless, the consistent use of the same Large Language Model (LLM) as both the baseline and the generation method mitigates this concern, utilizing highly similar prompts. A stronger baseline, such as GPT-4, would also strengthen these results. Finally, conducting subjective naturalness and coherence measurements on a larger number of utterances and participants would further validate our findings, enabling the detection of significant differences.

## References

- Christoph Bartneck, Dana Kulić, Elizabeth Croft, and Susana Zoghbi. 2009. Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots. *International journal of social robotics*, 1:71–81.
- Timothy W Bickmore, Everlyne Kimani, Ha Trinh, Alexandra Pusateri, Michael K Paasche-Orlow, and Jared W Magnani. 2018. Managing chronic conditions with a smartphone-based conversational virtual agent. In *Proceedings of the 18th International Conference on Intelligent Virtual Agents*, pages 119–124.
- Timothy W Bickmore, Laura M Pfeifer, and Brian W Jack. 2009. Taking the time to care: empowering low health literacy hospital patients with virtual nurse agents. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 1265–1274.
- Gillian Cameron, David Cameron, Gavin Megaw, Raymond Bond, Maurice Mulvenna, Siobhan O’Neill, Cherie Armour, and Michael McTear. 2017. Towards a chatbot for digital counselling. In *Proceedings of the 31st International BCS Human Computer Interaction Conference (HCI 2017) 31*, pages 1–7.
- Senthilkumar Chandramohan, Matthieu Geist, Fabrice Lefevre, and Olivier Pietquin. 2011. User simulation in dialogue systems using inverse reinforcement learning. In *Interspeech 2011*, pages 1025–1028.
- P.W.D. Charles. 2024. Bark. <https://github.com/suno-ai/bark>.
- Yu Ying Chiu, Ashish Sharma, Inna Wanyin Lin, and Tim Althoff. 2024. A computational framework for behavioral assessment of llm therapists. *arXiv preprint arXiv:2401.00820*.
- Kerstin Denecke, Sayan Vaaheesan, and Aaganya Arulnathan. 2020. A mental health chatbot for regulating emotions (sermo)-concept and usability test. *IEEE Transactions on Emerging Topics in Computing*, 9(3):1170–1182.
- Bradley Efron and Robert J Tibshirani. 1994. An introduction to the bootstrap: Crc press. *Ekman, P., & Friesen, WV (1978). Manual for the facial action coding system*.
- Amelia Fiske, Peter Henningsen, and Alena Buyx. 2019. Your robot therapist will see you now: ethical implications of embodied artificial intelligence in psychiatry, psychology, and psychotherapy. *Journal of medical Internet research*, 21(5):e13216.
- Siska Fitrianie, Merijn Bruijnes, Deborah Richards, Andrea Bönsch, and Willem-Paul Brinkman. 2020. The 19 unifying questionnaire constructs of artificial social agents: An iva community analysis. In *Proceedings of the 20th ACM International Conference on Intelligent Virtual Agents*, pages 1–8.
- Kathleen Kara Fitzpatrick, Alison Darcy, and Molly Vierhile. 2017. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): a randomized controlled trial. *JMIR mental health*, 4(2):e7785.
- Lucie Galland, Catherine Pelachaud, and Florian Pecune. 2022. Adapting conversational strategies in information-giving human-agent interaction. *Frontiers in Artificial Intelligence*, 5:1029340.
- Lucie Galland, Catherine Pelachaud, and Florian Pecune. 2024a. Emmi—empathic multimodal motivational interviews dataset: Analyses and annotations. *arXiv preprint arXiv:2406.16478*.
- Lucie Galland, Catherine Pelachaud, and Florian Pecune. 2024b. Simulating patient oral dialogues: A study on naturalness and coherence of conditioned large language models. In *Proceedings of the 24th ACM International Conference on Intelligent Virtual Agents*, pages 1–4.
- Mouzhi Ge, Carla Delgado-Battenfeld, and Dietmar Jannach. 2010. Beyond accuracy: evaluating recommender systems by coverage and serendipity. In *Proceedings of the fourth ACM conference on Recommender systems*, pages 257–260.
- Bhautesh Dinesh Jani, David N Blane, and Stewart W Mercer. 2012. The role of empathy in therapy and the physician-patient relationship. *Forschende Komplementärmedizin/Research in Complementary Medicine*, 19(5):252–257.
- Marius Kaminskas and Derek Bridge. 2016. Diversity, serendipity, novelty, and coverage: a survey and empirical analysis of beyond-accuracy objectives in recommender systems. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 7(1):1–42.
- Hsien-Chin Lin, Christian Geishauer, Shutong Feng, Nurul Lubis, Carel van Niekerk, Michael Heck, and Milica Gašić. 2022. Gentus: Simulating user behaviour and language in task-oriented dialogues with generative transformers. *arXiv preprint arXiv:2208.10817*.
- Hsien-chin Lin, Nurul Lubis, Songbo Hu, Carel van Niekerk, Christian Geishauer, Michael Heck, Shutong Feng, and Milica Gašić. 2021. Domain-independent user simulation with transformers for task-oriented dialogue systems. *arXiv preprint arXiv:2106.08838*.
- Christine Lisetti, Reza Amini, Ugan Yasavur, and Naph-tali Rische. 2013. I can help you change! an empathic virtual agent delivers behavior change health interventions. *ACM Transactions on Management Information Systems (TMIS)*, 4(4):1–28.
- Molly Magill, Timothy R Apodaca, Brian Borsari, Jacques Gaume, Ariel Hoadley, Rebecca EF Gordon, J Scott Tonigan, and Theresa Moyers. 2018. A meta-analysis of motivational interviewing process:

- Technical, relational, and conditional process models of change. *Journal of consulting and clinical psychology*, 86(2):140.
- Molly Magill, Jacques Gaume, Timothy R Apodaca, Justin Walthers, Nadine R Mastroleo, Brian Borsari, and Richard Longabaugh. 2014. The technical hypothesis of motivational interviewing: A meta-analysis of mi’s key causal model. *Journal of consulting and clinical psychology*, 82(6):973.
- Ganeshan Malhotra, Abdul Waheed, Aseem Srivastava, Md Shad Akhtar, and Tanmoy Chakraborty. 2022. Speaker and time-aware joint contextual learning for dialogue-act classification in counselling conversations. In *Proceedings of the fifteenth ACM international conference on web search and data mining*, pages 735–745.
- José Mercado, Ismael Edrein Espinosa-Curiel, and Juan Martínez-Miranda. 2023. Embodied conversational agents providing motivational interviewing to improve health-related behaviors: Scoping review. *Journal of Medical Internet Research*, 25:e52097.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- William R Miller, Theresa B Moyers, Denise Ernst, and Paul Amrhein. 2003. Manual for the motivational interviewing skill code (misc). *Unpublished manuscript. Albuquerque: Center on Alcoholism, Substance Abuse and Addictions, University of New Mexico*.
- William R Miller and Stephen Rollnick. 2012. *Motivational interviewing: Helping people change*. Guilford press.
- Zach Nussbaum, John X. Morris, Brandon Duderstadt, and Andriy Mulyar. 2024. [Nomic embed: Training a reproducible long context text embedder](#). *Preprint*, arXiv:2402.01613.
- Stefan Olafsson, Teresa K O’Leary, and Timothy W Bickmore. 2020a. Motivating health behavior change with humorous virtual agents. In *Proceedings of the 20th ACM international conference on intelligent virtual agents*, pages 1–8.
- Stefan Olafsson, Byron C Wallace, and Timothy W Bickmore. 2020b. Towards a computational framework for automating substance use counseling with virtual agents. In *AAMAS*, pages 966–974. Auckland.
- Florian Pecune, Stacy Marsella, and Alankar Jain. 2020. A framework to co-optimize task and social dialogue policies using reinforcement learning. In *Proceedings of the 20th ACM international conference on intelligent virtual agents*, pages 1–8.
- Prolific. 2023. [Prolific](#).
- Lazlo Ring, Timothy Bickmore, and Paola Pedrelli. 2016. An affectively aware virtual therapist for depression counseling. In *ACM SIGCHI Conference on Human Factors in Computing Systems (CHI) workshop on Computing and Mental Health*, pages 01951–12.
- Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young. 2007. Agenda-based user simulation for bootstrapping a pomdp dialogue system. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Companion Volume, Short Papers*, pages 149–152.
- Jost Schatzmann, Karl Weilhammer, Matt Stuttle, and Steve Young. 2006. A survey of statistical user simulation techniques for reinforcement-learning of dialogue management strategies. *The knowledge engineering review*, 21(2):97–126.
- Michael Schoeffler, Sarah Bartoschek, Fabian-Robert Stöter, Marlene Roess, Susanne Westphal, Bernd Edler, and Jürgen Herre. 2018. webmushra—a comprehensive framework for web-based listening tests. *Journal of Open Research Software*, 6(1):8.
- Edward H Simpson. 1949. Measurement of diversity. *nature*, 163(4148):688–688.
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Kraemer. 2021. Human evaluation of automatically generated text: Current trends and best practice guidelines. *Computer Speech & Language*, 67:101151.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. 2021. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in neural information processing systems*, 34:22419–22430.
- Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. 2022. Towards a unified multi-dimensional evaluator for text generation. *arXiv preprint arXiv:2210.07197*.

## A Appendix

### A.1 Utterance generation prompt

The condition utterances are generated through the prompt visible in Table 5

### A.2 Utterance generation examples

Examples of generated utterances are in Table 6).

Condition	Prompt
<b>Base</b>	You are a patient talking with a therapist. Your task is to generate the patient's next utterance. The dialog is happening orally, use a oral style language with hesitation, repetition, and deviation. You will only respond with one patient's utterance. Do not provide explanations or notes. Make only one proposition. The response must be short, no more than 2 or 3 utterances.
<b>Unconditioned</b>	<b>Base prompt +</b> {Utterance examples} <<< Context : {context} Generate the patient's next utterance >>> Patient's utterance:
<b>Dialog act conditioned</b>	<b>Base prompt +</b> {Utterance examples with corresponding dialog act} <<< Context: {context} Generate the patient's next utterance with the intent: {intent} >>> Patient's utterance:

Table 5: Prompts template for each of the two conditions

Context	Original utterance	Baseline (LLM)	Text (Ours)	DA+ Text + Type (Ours)
<p>Turn 71/79</p> <p>Therapist: You're quiet down.</p> <p>Patient: yeah</p> <p>Therapist: Are you keeping interest in anything, is there anything that that you find enjoyment in?</p>	<p>[Sustaining unhealthy behavior]</p> <p>I am not doing the thing there is nothing for me to do</p>	<p>[Sharing personal information]</p> <p>I like to watch movies and play video games.</p>	<p>[Sharing positive feeling or emotion]</p> <p>Yeah, I love to play basketball. It really makes me happy when I score a goal.</p>	<p>[Sharing positive feeling or emotion]</p> <p>Yeah, I love to play basketball. It really makes me happy when I score a goal.</p>
<p>Turn 21/80</p> <p>Therapist: Have you criticized other coworkers?</p> <p>Patient: yeah. Therapist: So do you think that her criticism was within normal limits like what you've seen there regularly.</p>	<p>[Realization/Understanding]</p> <p>I mean, now that you pointed out I mean, I guess I can view it in that way that it wasn't.</p>	<p>[Realization/Understanding]</p> <p>Well, I guess it depends on who it is and what they did. Sometimes I think people deserve criticism. [laughs] But other times, maybe it's not warranted. It's hard to say for sure.</p>	<p>[Sharing positive feeling or emotion]</p> <p>Yeah, I guess so. It feels good to be able to express my thoughts openly and honestly without fear of judgment.</p>	<p>[Realization/Understanding]</p> <p>I guess so. It depends on the situation. Maybe if she was being too harsh, but sometimes coworkers need a push.</p>
<p>Turn 23/67</p> <p>Therapist: For the position. Are you eligible for unemployment?</p> <p>Patient: I'm not sure I'm still looking into that.</p> <p>Therapist: Because if you, I mean, I'm not a labor law expert. But I think if you're, if you're laid off, there's probably decent probability you could get on employment. How much for educational?</p>	<p>[Sharing personal information]</p> <p>I just have my bachelor's.</p>	<p>[Sharing personal information]</p> <p>don't know yet. I haven't checked into that. I was just focusing on getting my unemployment situation figured out first. But I might look into educational assistance as well. It depends on the cost and what's available.</p>	<p>[Sharing personal information]</p> <p>don't know yet. I haven't checked into that. I was just focusing on getting my unemployment situation figured out first. But I might look into educational assistance as well. It depends on the cost and what's available.</p>	<p>[Sustaining unhealthy behavior in the future]</p> <p>I don't really want to depend on that, though. I prefer to be independent.</p>

Table 6: Examples of generated utterances in each condition given their context

# Training LLMs to Recognize Hedges in Dialogues about Roadrunner Cartoons

Amie J. Paige<sup>\*Ψ</sup>, Adil Soubki<sup>\*□□</sup>, John Murzaku<sup>\*□□</sup>, Owen Rambow<sup>●□</sup>,  
Susan E. Brennan<sup>Ψ</sup>

□ Department of Computer Science, ● Department of Linguistics, Ψ Department of Psychology □ Institute for Advanced Computational Science, Stony Brook University

\*These authors contributed equally to this study.

amie.paige@stonybrook.edu, {asoubki, jmurzaku}@cs.stonybrook.edu

## Abstract

Hedges allow speakers to mark utterances as provisional, whether to signal non-prototypicality or “fuzziness”, to indicate a lack of commitment to an utterance, to attribute responsibility for a statement to someone else, to invite input from a partner, or to soften critical feedback in the service of face-management needs. Here we focus on hedges in an experimentally parameterized corpus of 63 Roadrunner cartoon narratives spontaneously produced from memory by 21 speakers for copresent addressees, transcribed to text (Galati and Brennan, 2010). We created a gold standard of hedges annotated by human coders (the *Roadrunner-Hedge corpus*) and compared three LLM-based approaches for hedge detection: fine-tuning BERT, and zero and few-shot prompting with GPT-4o and LLaMA-3. The best-performing approach was a fine-tuned BERT model, followed by few-shot GPT-4o. After an error analysis on the top performing approaches, we used an *LLM-in-the-Loop* approach to improve the gold standard coding, as well as to highlight cases in which hedges are ambiguous in linguistically interesting ways that will guide future research. This is the first step in our research program to train LLMs to interpret and generate collateral signals appropriately and meaningfully in conversation.

## 1 Introduction

The virtuosity of LLMs such as ChatGPT has led some to the impression that AI already converses (or will soon be able to converse) as people do. But as language users, LLMs and humans are quite different. The underlying foundations for learning by these distinct kinds of language users share little in common: Humans learn as infants to interact with others well before they learn their first words, and once word learning begins, they can pick up a new word in one or just a few exposures, whereas LLMs are pre-trained on humanly unfathomable quantities of text without ever learning to inter-

act. Transformer-based chat programs can generate paragraphs-worth of text remarkably well without modeling the coordination between agents—but is this conversation?

Whether a sequence of prompts and responses exchanged in a dialogue between an LLM agent and a human counts as truly (rather than superficially) “conversational” depends on how conversation is conceptualized. Conversation is often presumed to be the passing back and forth of messages (a “message model”); but that does not explain phenomena common to spontaneous conversation such as incremental turns, clarifications, and repair. Here we conceptualize conversation as a collaborative process of grounding meanings (seeking and providing evidence) during which two or more partners signal, coordinate, and align their beliefs or cognitive states (Brennan, 2005; Clark and Wilkes-Gibbs, 1986). This leads to a broader research agenda that we hope will push generative AI to model phenomena such as a partner’s knowledge or theory of mind, mutual beliefs or common ground, as well as when to take initiative in a dialogue.

The main contributions of this work include:

- (i) After grounding the project in psycholinguistic theory (Section 2) and related work (Section 3), we present the Roadrunner-Hedge Corpus (Section 4), a corpus of spontaneous face-to-face narratives annotated for hedging.<sup>1</sup>
- (ii) We describe a set of experiments on this corpus using zero-shot, few-shot, and fine-tuning methods on modern LLMs (Section 5).
- (iii) We perform a detailed error analysis pinpointing where LLMs fail in detecting hedges (Section 6). With this analysis, we take an LLM-in-the-Loop approach to correcting gold annotations, reducing errors in our top performing systems.

We conclude with a discussion and implications of our results in Section 7, limitations and the future

<sup>1</sup><https://github.com/cogstates/hedging>

of our work in Section 8, and a final summary of our salient contributions in Section 9.

## 2 Theoretical Foundations from Psycholinguistics

In conversation, people communicate not only about the purpose or topic at hand, but they also communicate meta-information about what they're saying within the context of interaction, or *collateral signals* (Clark, 1996). Along with providing evidence for grounding in conversation, about whether a prior turn has been understood as intended (Clark and Brennan, 1991), collateral signals can also provide information about the speaker's relationship with the content of their message—how confident they are in what they are saying, whether it is difficult to recall or express, and whether they would welcome input from their partner. In this project, we focus on a particular kind of collateral signal used for coordination, *hedges*.

### 2.1 Why Speakers Hedge

There have been several proposals for why speakers hedge. Hedges have been claimed to characterize powerless “feminine” language (Lakoff, 1973) or to serve a politeness function by minimizing threat to a partner's “face” (Brown and Levinson, 1987); see also (Fraser, 2010). Hedges have also been thought to convey a certain “fuzziness” of category membership when a speaker means to describe a non-prototypical member of a category (e.g., a penguin belonging to the bird category; Lakoff, 1975). Prince et al. (1982) suggested that hedges play two functions: First, to make propositional content less exact (*approximators*, e.g. “sort of”) and second, to change the relationship a speaker has to the content of their message (*shield hedges*). Shield hedges are further divided into *plausibility* shields that signal a lack of commitment to the content of a message (“I think his feet were blue,” Prince et al., 1982, p. 5), and *attribution* shields that assign responsibility for a message to a source other than the speaker or writer themselves (“According to her estimates...” Prince et al., 1982, p. 13).

Several experimental studies have demonstrated how hedges can convey speakers' commitment to what they are saying. For example, in a question-answering task, people trying to recall the answers to trivia questions produced more disfluencies, longer latencies, more rising intonation, and more expressions of doubt when they reported hav-

ing a low *feeling of knowing* about an answer. This metacognitive information was confirmed to be accurate when compared to the ground truth in the form of their answer to the same (multiple-choice) question later (Smith and Clark, 1993). Not only are hedges informative as collateral signals about what a speaker knows, but they are accurately interpreted as such by listeners (Brennan and Williams, 1995).

That hedges function as interactional signals in extended dialogue is evident from studies of referential communication. Typically in such studies, two partners who can't see each other converse in order to arrange and rearrange duplicate sets of objects in matching orders, with the objects needing to be distinguished from similar objects or consisting of Tangrams (abstract geometric shapes unassociated with any conventional or lexicalized labels). Hedges are common in initial referring expressions, where they tend to appear in wordy, disfluent, and often tentative descriptions, and then they drop out in repeated referring expressions once partners have reached a shared conceptualization for that object (marked by entrainment, or re-using the same shortened referring expression) (Brennan and Clark, 1996; Galati and Brennan, 2021), as in this sequence of repeated references to the same object over multiple rounds (adapted from Brennan and Clark, 1996, p. 1488):

*Round 1:* “a car, sort of silvery purple colored”

*Round 2:* “purplish car going to the left”

...

*Round 5:* “the purple car”

In another study that required triads of strangers to reach consensus while recalling the events from a movie clip that they had watched earlier, the speakers often hedged their contributions to the conversation, presumably to mark a lack of certainty about an utterance and an openness to being corrected by their partners (Brennan and Ohaeri, 1999). For example, from a triad that communicated by speaking face-to-face:

*Yeah, they were sitting around the fireplace in the night... sort of like a bedtime story kind of thing*

People who did the same task by texting rather than speaking used fewer words, but still hedged:

*We all agree it was a wreathy thingy on his neck???*

## 2.2 How Listeners React to Hedges

Hedges convey meaningful information that can affect listeners' subsequent behavior; a handful of psychological studies have measured the impacts of hedges on listeners. For example, children exposed to new words from a speaker who hedged learned fewer novel words compared to children exposed to a speaker who did not hedge (Sabbagh and Baldwin, 2001). Listeners rated utterances as more uncertain when they included shield hedges (e.g., "I think it was a mug"), and these ratings were related to speakers' ratings of their own uncertainty in identifying an image (Pogue and Tanenhaus, 2018). Moreover, addressees in a referential communication task expended more effort while grounding (they produced more low-confidence responses such as clarification questions) to demonstrate understanding when the speaker's description had contained a hedge (Dahan, 2023).

Hedges also influence which details are retold to another person; in one study, hedged details were less likely to be repeated to another addressee as compared to unhedged details (Liu and Fox Tree, 2012), although in the same study, hedged information presented in a story was more likely to be remembered by listeners; this was thought to stem from deeper engagement with hedged information when it was first presented (Liu and Fox Tree, 2012). And in tutoring dialogues, where face management can be particularly important, students were more successful at solving problems when their peer tutors used hedges (Madaio et al., 2017).

## 3 Related Computational Work

### 3.1 Hedging

Several research programs have examined hedges and the criteria for coding them, with computational goals that include automatic hedge detection. Hedging is domain-specific, meaning that their forms and frequencies vary across corpora; they are also context-specific, as they cannot be identified accurately simply by searching for strings (Prokofieva and Hirschberg, 2014). Hedges are distributed differently within different corpora (ibid).

Hedges are often ambiguous and difficult to code in the absence of dialogue context. In "I think it's a little odd," *I think* is often a hedge, but might not be when proffered in response to a question ("So what do *you* think?"). Hedges in spoken utterances may be disambiguated by stress and other intonational cues, as in "I think he'll win!" (not a

hedge) vs. "I *think* he'll win?" (a hedge). Previous work found many cases of tokens that can serve as hedges as well as non-hedges, with systematic tests for coders to use in annotating them for gold standards (Prokofieva and Hirschberg, 2014; Ulinski and Hirschberg, 2019; Ulinski et al., 2018).

The coding of hedges is complicated by the fact that in spoken dialogue, they often co-occur with speech disfluencies. In some contexts, it may be difficult to distinguish these two kinds of signals (Prokofieva and Hirschberg, 2014), particularly since listeners can use disfluencies in much the same way they can use hedges to draw conclusions about the speaker's mental state (Arnold et al., 2003, 2007)

A strong motivation for computational work on hedging comes from work on computer-assisted learning by Cassell and colleagues, specifically tutoring dialogues (Abulimiti et al., 2023a,b; Raphalen et al., 2022). Most similar to our work is Raphalen et al. (2022), where the authors propose a model that combines rule-based classifiers and machine learning models with interpretable features such as unigram and bigram counts, part-of-speech tags, and LIWC categories to identify and classify hedge clauses. Our work differs in two major ways: first, our work operates on the **token** level rather than on the clause level. Token level classification makes possible a truly end-to-end approach (classifying all hedge and non-hedge tokens in utterances). Second, we include experiments with modern LLMs and offer a detailed error analysis into their mistakes; stemming from this error analysis, we use an LLM-in-the-Loop approach (Dai et al., 2023) to correcting gold standard hedge codings.

### 3.2 Belief

Hedging and the notion of belief (how committed the speaker is to the truth of an event) are closely related; hedges are often used by speakers to indicate a lack of belief or commitment towards what they say. Ulinski et al. (2018) improved belief classification using a hedge detector, yielding an improvement for the non-committed and reported belief labels.

**Corpora** Several corpora have been created that annotate the author's degree of belief (Diab et al., 2009; Prabhakaran et al., 2010; Lee et al., 2015; Stanovsky et al., 2017; Rudinger et al., 2018; Poursan Ben Veyseh et al., 2019; Jiang and de Marn-

Hedge Type	Example(s)
Like (not used as a simile, verb, or comparison)	"and then he like went over by..."
You know (not to communicate another's knowledge or as a discourse marker)	"and you know as he's falling down"
Just (not used to mean "only")	"he just jolts away"
Approximators/Rounders	"kind of", "about"
Proxies (for a detail the speaker cannot or chooses not to recall)	"thing," "whatever," "or something," "and everything"
Morpheme suffixes to content words	"circley," "springy"
Expressions of doubt attached to claims; self-speech	"I don't know," "maybe," "I guess," "what's it called?"
Tag questions and try markers	"he's standing there, right?"

Table 1: Coding scheme used to mark hedges in corpus.

effe, 2021). There are two corpora that further annotate nested beliefs of the sources mentioned in the text: FactBank (Saurí and Pustejovsky, 2009) and the Modal Dependency corpus (Yao et al., 2021).

**Machine Learning Approaches** Modern neural methods for belief detection include LSTMs with multi-task or single-task approaches (Rudinger et al., 2018), using BERT representations alongside a graph convolutional neural network (Puran Ben Veyseh et al., 2019), or fine-tuning BERT with a span self-attention mechanism Jiang and de Marneffe (2021). Recent state-of-the-art work finds that fine-tuning RoBERTa (Murzaku et al., 2022) or fine-tuning Flan-T5 (Murzaku et al., 2023) yields the best performance on most corpora. For the label *Underspecified* (or, corresponding to no commitment and/or a hedge), these modern methods yield f-measures in the low to high 80s. We also have prior work exploring multi-modal approaches to belief detection (Murzaku et al., 2024).

#### 4 The Roadrunner-Hedge Corpus

For training and testing, we obtained a corpus (Galati and Brennan, 2010) of spontaneous narratives produced from memory by 20 speakers who had watched a Roadrunner cartoon. Each speaker narrated the story face-to-face to an audience, a total of three times: first to a naïve addressee, a second time to the same addressee, and a third time to a new naïve addressee (with the latter two episodes counterbalanced for order). The original experiment was designed to detect differences in collateral signals (intelligibility vs. attenuation of speech and gestures) stemming from the speaker's vs. the addressee's knowledge states—that is, whether the story was new for the speaker (told for the first time) vs. old (retold), compared to the addressee's knowledge state (new vs. heard for the second time). Findings included that the attenuation of both referring expressions (Galati

and Brennan, 2010) and gestures (Galati and Brennan, 2014) were driven by *both* speakers' and addressees' knowledge states—that is, shortened upon retelling the story to the same addressee, but lengthened upon retelling to a new addressee.

**Gold Standard Coding.** The original corpus transcribed the spontaneous narratives in detail, including speaking turns and disfluencies (for details, see Galati and Brennan, 2010), segmented into lines by installments that corresponded to narrative elements in the cartoons. We annotated hedges on the original Roadrunner corpus to create the gold standard for hedge training and detection (the *Roadrunner-Hedge* corpus; see <https://github.com/cogstates/hedging> for the annotation codebook).

The Roadrunner-Hedge corpus is distributed as a csv file. It is structured as a total of 5,508 lines, over a quarter of which (N=1424) include one or more hedges. The first author annotated hedges in the corpus as in Table 1. Although disfluencies such as fillers (*uh*, *um*) and re-starts can function as hedges, we made a principled decision to not code them as such; hedges in our corpus are presumed to be shaped by the speaker's intention, whereas disfluencies are not necessarily under a speaker's control as a communicative signal, but may reflect difficulties in speaking (Grice, 1957; Clark, 1994). Overall word counts for hedges and non-hedges are 1,728 and 38,018 words respectively. Most hedges are one word, but a few cases contain many words. For each line in the csv file (corresponding to a narrative element), hedges are listed (separated by commas) in an adjacent cell. Each line has an average of 0.33 hedges.

**Inter-Rater Reliability.** To compute inter-rater reliability, a trained research assistant coded 7 randomly-selected transcripts with no overlapping speakers (10% of the corpus). We calculated Cohen's Kappa from each word marked as a hedge within each transcript. There was high agreement

between coders, with  $\kappa = 0.985$ .

**Corpus Analysis.** The Roadrunner-Hedge corpus, like the tutoring dialogues used by [Abulimiti et al. \(2023b\)](#); [Raphalen et al. \(2022\)](#), has fewer cases with hedges than without, but with more hedges per segment overall (25.85% of lines vs. 14.26% of turns respectively).

Over the three versions of the cartoon story produced by each speaker, hedges were most frequent in the first telling when the story was new to both speaker and addressee and least frequent when told to the same addressee a second time, consistent with the original findings from [Galati and Brennan](#) that collateral signals are affected by the knowledge states of both speaker and addressee.

## 5 Experiments

### 5.1 Experimental Setup

In this section, we present our hedge classification experiments on the Roadrunner-Hedge corpus, conducted by fine-tuning BERT and performing zero-shot and few-shot experiments with state-of-the-art LLMs. For all experiments, we performed five-fold cross validation using a fixed seed (42), splitting the corpus into a 80/20 train/test split. For our fine-tuning experiments, we did not perform any hyperparameter tuning, and therefore do not have a validation set.

We performed all zero-shot, few-shot, and fine-tuning experiments on the fold’s respective test sets and report the average and standard deviation over all five folds test sets for **F1**, **precision**, and **recall**.

### 5.2 Zero Shot and Few Shot

For the zero-shot and few-shot experiments, we used GPT-4o ([OpenAI, 2024](#)) and LLaMA-3-8B-Instruct ([AI@Meta, 2024](#)), as these two LLMs have achieved state-of-the-art results in many zero-shot or few-shot benchmark tasks.

We conducted two classes of zero-shot and few-shot experiments: count/list generation and **BIO** tag generation. Both prompts began with an instruction detailing the specific task, and a random example. In our few-shot experiments, we provided three fixed hand-crafted examples. For our count/list generation, we prompted the models to list the integer number of hedges present in the utterance and then generated a list of the exact hedge words. For our **BIO** tag generation, we generated the tokens and their respective tags, where label *B* represents the beginning of a hedge token or span,

*I* represents the inside of a hedge span, and *O* represents another token, all separated by “/”. For example, given the utterance *It is like warm*, we prompted the model to generate *It/O is/O like/B warm/O*.

We provide our exact prompts with their corresponding instructions in [Appendix A](#). For our GPT-4o experiments, we used the default OpenAI API hyperparameters and a **temperature** of 1.0.

### 5.3 Fine-tuning

We performed all fine-tuning experiments using BERT ([Devlin et al., 2019](#)), specifically bert-base-uncased. We also performed experiments with the large variants of the model (bert-large), newer encoder-only models like RoBERTa ([Liu et al., 2019](#)) and DeBERTa-v3 ([He et al., 2021](#)), and encoder-decoder models like Flan-T5 ([Chung et al., 2022](#)), but got either worse or closely similar results.

**Task Description** All experiments followed a standard BIO token labelling approach to classify hedge tokens (*B*), tokens inside of hedge spans (*I*), and all other tokens (*O*). In other words, given an input utterance of *n* tokens, the respective BIO labels were output for each of the *n* tokens. Following the same example as described in our zero-shot and few-shot experiments in [Section 5.2](#), we fine-tuned BERT to classify the tokens as *It/O is/O like/B warm/O*.

**Hyperparameters** We followed a standard fine-tuning approach, fine-tuning for a fixed 5 **epochs**. We set the batch size to 16 and learning rate to  $2e-5$ . We performed five-fold cross validation and test on each folds respective test set. We did not perform any hyperparameter tuning.

### 5.4 Results

The performance of the models is shown in [Table 2](#), which reports average precision (**P**), recall (**R**), and **F1** over the five-folds. For our zero-shot, few-shot, and fine-tuning experiments, these metrics are calculated on each fold’s test set and then averaged.

Despite its much smaller parameter count, BERT fine-tuned for **BIO** tagging outperforms even the best scoring prompting approaches by nearly 20 points in F-measure. This is consistent with a general trend in the literature of more parameter efficient fine-tuning approaches outperforming larger zero-shot and few-shot methods ([Liu et al., 2022](#)), though the gap here is larger than one might expect.

Model	Training	Prompt	Precision (P)	Recall (R)	F1 Score (F1)
BERT	Finetuned	-	0.883 $\pm$ 0.015	0.934 $\pm$ 0.012	0.908 $\pm$ 0.010
GPT-4o	Few-Shot	List	0.613 $\pm$ 0.027	0.848 $\pm$ 0.018	0.712 $\pm$ 0.021
LLaMA-3	Few-Shot	List	0.518 $\pm$ 0.035	0.799 $\pm$ 0.022	0.628 $\pm$ 0.031
GPT-4o	Few-Shot	BIO	0.514 $\pm$ 0.024	0.766 $\pm$ 0.036	0.616 $\pm$ 0.030
GPT-4o	Zero-Shot	List	0.430 $\pm$ 0.014	0.711 $\pm$ 0.004	0.536 $\pm$ 0.012
GPT-4o	Zero-Shot	BIO	0.436 $\pm$ 0.026	0.618 $\pm$ 0.033	0.510 $\pm$ 0.028
LLaMA-3	Few-Shot	BIO	0.298 $\pm$ 0.018	0.625 $\pm$ 0.016	0.404 $\pm$ 0.019
LLaMA-3	Zero-Shot	BIO	0.167 $\pm$ 0.014	0.428 $\pm$ 0.019	0.240 $\pm$ 0.017
LLaMA-3	Zero-Shot	List	0.274 $\pm$ 0.023	0.146 $\pm$ 0.010	0.190 $\pm$ 0.011

Table 2: Average performance metrics over the five folds with standard deviations for different models, training methods, and prompt types, ordered by F1 score.

In comparisons of the zero-shot and few-shot prompting methods, the few-shot models, unsurprisingly, performed better. The few-shot experiments averaged an F1 of 0.59, 22 points higher than the zero-shot models average of 0.37.

Of the two output formats prompted for, listing and BIO, the listing approach performed better. On average, models instructed to output a list had an F1 of 0.52 compared to 0.44 for those instructed to perform BIO tagging.

Among the two LLMs prompted, GPT-4o always performed best. Across all models and approaches, including fine-tuned BERT, precision tended to be lower than recall, with a mean of 0.46 for precision compared to 0.65 for recall. In other words, the models over-predicted the presence of hedges.

## 6 Error Analysis

While the fine-tuned BERT model performed fairly well, a certain number of cases did not align with the gold labels in the data. We performed error analysis to understand whether there were any systematic deviations from the corpus annotation.

We conducted an error analysis on the top two performing models, the fine-tuned BERT model and the GPT-4o Few-shot List (FSL) model (F1 = 0.91 and 0.71, respectively). Starting with the first fold, we selected the first hundred errors to categorize. These errors are broadly divided into instances where the models failed to detect a hedge (false negatives) and instances where models returned cases that were not annotated hedges (false positives). The remaining errors fell into two other categories: a gold error category, wherein errors in the (human) annotation were discovered, and an “other” category.

Of the hundred errors sampled from the BERT model, approximately the same number of errors were false negatives (26) as false positives (29). Of the hundred errors sampled from the GPT-4o FSL model, 66 were false positives and 25 were false negatives (reflecting the low precision and higher recall for this approach; see Table 3 and 4 for full error descriptions for BERT and GPT-4o FSL models).

Although the corpus annotation does not include the *type* of hedge (only the presence or absence of hedge tokens), our error analysis looked at hedge types in order to tease apart model behaviors. We observed systematic differences between models in their types of mismatches with the gold standard.

**False Positives.** First, the GPT-4o FSL model inaccurately classified disfluencies (e.g., “uh”) as hedges in 37 of the 66 false positives reviewed, whereas BERT did not. Second, BERT showed quite a different pattern of mismatches than GPT-4o when classifying “like”, returning false positives that always turned out to be comparatives (e.g., “it’s like an open elevator”). These we considered to be true errors in their text form, although some may be ambiguities that could be resolved prosodically.

**False Negatives.** Tokens denoting approximator hedges (e.g. “that’s *basically* it”) were frequently misclassified as false negatives by BERT (9 of 26 false negatives reviewed), but never by the GPT-4o FSL model.

In addition, **Other** emerged as a category type for situations that could not clearly be described as false positives, false negatives, or gold errors. In the BERT model, these cases were typically segmentation errors (i.e., an inner token mislabeled as a beginning token).

Notably, the largest class of errors for the BERT

<b>Gold Errors</b>		<b>False Negative</b>		<b>False Positive</b>		<b>Other</b>	
<i>Like</i>	13	Approximator	9	<i>Like</i>	13	<i>I should be B</i>	4
Proxy	12	Proxy	8	<i>Just</i>	8	<i>O should be I</i>	3
<i>Just</i>	7	Self-talk	4	False proxy	4	<i>B should be I</i>	2
Approximator	1	<i>Like</i>	3	<i>You know</i>	2	Other	2
Other	1	<i>Just</i>	1	Misc. word	2		
		Morpheme	1				
<b>Total</b>	<b>34</b>		<b>26</b>		<b>29</b>		<b>11</b>

Table 3: Expanded error analysis on the BERT fine-tuned model, by hedge type.

<b>Gold Errors</b>		<b>False Negative</b>		<b>False Positive</b>		<b>Other</b>	
Approximator	4	<i>Just</i>	12	Disfluency tag	37	Other	1
<i>Just</i>	1	Proxy	8	Misc. word	15		
<i>Like</i>	1	<i>Like</i>	3	<i>Like</i>	7		
Proxy	1	Morpheme	1	Approximator	3		
Self-talk	1	Self-talk	1	Intensifiers	3		
				<i>You know</i>	1		
<b>Total</b>	<b>8</b>		<b>25</b>		<b>66</b>		<b>1</b>

Table 4: Expanded error analysis on the GPT-4o FSL model, by hedge type.

model was the **Gold Error** category (34 of 100). This was not the case for the GPT-4o model (only 9 gold errors). The BERT fine-tuned model revealed mistakes made by the human annotators for hedges denoted by “like”, “just”, and proxy hedges (e.g. “and stuff”). Upon closer inspection, some of these cases were ambiguous. For example, “he just hits the ground” could be taken to mean that the only action performed was hitting the ground (where “just” means only) or “just” might function to reduce the speakers’ certainty (as in Madaio et al., 2017). Again, the text format of the storytelling corpus leaves some interpretations ambiguous that could be clarified with signals such as timing and prosodic stress.

The number of Gold Errors identified by the BERT model allowed us to modify the original gold annotation with missed cases and to re-evaluate the performance of our models more accurately – a sort of *LLM-in-the-Loop* approach (see Table 5).

## 7 Discussion

The results show that even enormous, recently released LLMs cannot reliably recognize hedges. There is no “emergent” ability in LLMs to understand full human linguistic behavior. On the other hand, when we explicitly train a small, rather old LLM (BERT) to perform our task by fine-tuning it, it performs quite well. What this shows is that detecting hedges is a capability that can be

learned, but it cannot be learned in the manner that LLMs are taught, namely by simply ingesting large amounts of varied data. We interpret this to mean that if we want to make LLMs able to converse with humans as humans do, we need to understand what capabilities LLMs need and how to provide them with the ability to do so.

The prevalence of gold errors discovered by the BERT model raises two interesting points for discussion. First, some of these discrepancies identified by the BERT model were clearly errors made by the human coders; this was true in particular for proxies, which BERT coded for hedges more consistently than did human coders. This error analysis allowed us to iteratively improve the human coding before the final analysis, essentially deploying an LLM-in-the-Loop approach. Second, the discrepancies between BERT and gold coding on the tokens *just* and *like* highlight that these types of hedges have high potential for ambiguity—perhaps the very sort of ambiguity that could be resolved by prosody.

## 8 Limitations and Future Work

This work represents the first step in our research program that aims to train LLMs to use collateral signals in support of human-LLM dialogue. Once hedges can be recognized by an LLM, it remains to be shown that they can be meaningfully interpreted and generated. Relevant work by Cassell and col-

Model	Original Gold F1	LLM-in-the-Loop Gold F1	Error Reduction (%)
BERT	0.908 $\pm$ 0.010	0.925 $\pm$ 0.019	18.5%
GPT-4o Few-Shot	0.712 $\pm$ 0.021	0.721 $\pm$ 0.020	3.1%
GPT-4o Zero-Shot	0.510 $\pm$ 0.028	0.551 $\pm$ 0.011	8.4%

Table 5: F1 scores with standard deviations on the original corpus, F1 scores with standard deviations obtained on the corpus corrected after LLM-in-the-Loop, and the change in average performance for our top performing models.

leagues has shown that it is possible to generate hedges in tutoring dialogues, but not always positioned where they are most probable or useful (Abulimiti et al., 2023a). In future work, we plan experiments using top-performing models such as BERT and GPT-4o in high- and low-probability situations that systematically vary the certainty associated with prompted-for information (where hedges can be most useful). It is already clear from our pilot trials using ChatGPT 3.5 that LLMs hedge somewhat superficially (hedging where humans wouldn’t and failing to hedge where humans would).

**Domains of Dialogue.** Here we have used human-generated dialogue from a single domain, retelling stories from Roadrunner cartoons; the training data are text transcripts of speech. Because the initiative was unbalanced in this collaborative task, most of the speaking in each triad was done by the the partner who viewed and retold the cartoon stories in series to the two co-present addressees.

A more balanced domain in which partners continuously monitor each other’s understanding to do a physical task—such as matching pictures of difficult-to-describe objects—could yield more hedges, distributed differently. We plan to conduct similar tests to replicate the current results on such referential communication corpora collected previously in our lab.

It is interesting that despite the fact that there is not a single instance of dialogue in Roadrunner cartoons (apart from Roadrunner’s smug, trademark “meep meep” upon escaping from Coyote), speakers who retell the story in a dramatic and humorous way do a great deal of what looks like quoting Coyote’s and Roadrunner’s reactions:

*so then he’s saying he’s like gone all sad  
and stuff you know?*

*and he’s like whatever she’s gonna be  
dead right?*

Such uses of *like* in this corpus match the quotation-as-demonstrations forms described by Clark and

Gerrig (1990); they count as hedges in that the speaker marks what follows as *not verbatim*.

**Training with audio input.** Our results for detecting hedges in this transcribed spoken corpus are surprisingly strong, especially given that the LLMs we used were pre-trained primarily on originally written text. But it is well-known that features such as pausing and intonation are related to speakers’ levels of commitment to and confidence in their utterances. We plan to incorporate audio into future hedging studies and will explore multi-modal neural architectures fusing both speech and lexical features as we did in (Murzaku et al., 2024) for belief recognition.

**Reliability.** It is critical to keep in mind that human and LLMs are very different sorts of agents. Psychometric tests show that individual humans are likely to respond consistently when tested repeatedly, whereas an LLM is not (Shu et al., 2024). LLMs have no sense of “self” and are likely to respond differently when re-prompted with the same prompt. To the extent that a hedge signals that a speaker does not wish to be held entirely accountable for what they’re saying, hedging on the part of an LLM may actually be desirable as a way to encourage users to not assume they can hold it accountable. On the other hand, it may be desirable for an LLM to be able to signal its *confidence* – the reliability or quality (or lack thereof) of information it’s presenting – through the presence or absence of hedges. Finally, it remains to be seen whether LLMs can learn about interaction through exposure to collateral signals in meaningful contexts.

## 9 Conclusion

Our project is grounded in psycholinguistic theory and aims to capture theory-of-mind aspects of hedging among discourse participants. We present the Roadrunner-Hedge corpus, with hedges annotated from naturally occurring dialogues by speakers describing Roadrunner cartoons. We use the corpus to

train and perform experiments on detecting hedges using BERT, GPT-4o, and LLaMA-3. We find that fine-tuning BERT significantly outperforms state-of-the-art LLMs in few-shot and zero-shot settings. With our systems outputs, we perform an error analysis and use an LLM-in-the-Loop approach to correct gold standard annotations. Our LLM-in-the-loop approach provided further error reductions on all models.

## Ethical Considerations

The Roadrunner-Hedge corpus was collected with Institutional Review Board approval from undergraduate students who gave informed consent prior to participating in the experiments.

## Acknowledgments

This material is based upon work supported in part by the National Science Foundation (NSF) under No. 2125295 (NRT-HDR: Detecting and Addressing Bias in Data, Humans, and Institutions) as well as by funding from the Defense Advanced Research Projects Agency (DARPA) under the CCU program (No. HR001120C0037, PR No. HR0011154158, No. HR001122C0034). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the NSF or DARPA.

We thank both the Institute for Advanced Computational Science and the Institute for AI-Driven Discovery and Innovation at Stony Brook for access to the computing resources needed for this work. These resources were made possible by NSF grant No. 1531492 (SeaWulf HPC cluster maintained by Research Computing and Cyberinfrastructure) and NSF grant No. 1919752 (Major Research Infrastructure program), respectively.

We would also like to thank our reviewers for their helpful comments, as well as Kayla Hunt for assistance with reliability coding.

## References

Alafate Abulimiti, Chloé Clavel, and Justine Cassell. 2023a. [When to generate hedges in peer-tutoring interactions](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 572–583, Prague, Czechia. Association for Computational Linguistics.

Alafate Abulimiti, Chloé Clavel, and Justine Cassell. 2023b. [How about kind of generating hedges](#)

[using end-to-end neural models?](#) *Preprint*, arXiv:2306.14696.

AI@Meta. 2024. [Llama 3 model card](#).

Jennifer Arnold, Maria Fagnano, and Michael Tanenhaus. 2003. [Disuencies signal thee, um, new information](#). *Journal of Psycholinguistic Research*, 32:25–36.

Jennifer Arnold, Carla Hudson Kam, and Michael Tanenhaus. 2007. [If you say thee uh you are describing something hard: The on-line attribution of disfluency during reference comprehension](#). *Journal of experimental psychology. Learning, memory, and cognition*, 33:914–30.

Susan E. Brennan. 2005. How conversation is shaped by visual and spoken evidence. In John Trueswell and Michael Tanenhaus, editors, *Approaches to studying world-situated language use: Bridging the language-as-product and language-as-action traditions*, pages 95–129. MIT Press.

Susan E. Brennan and Herbert H. Clark. 1996. Conceptual pacts and lexical choice in conversation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 22(6):1482–1493.

Susan E. Brennan and J. O. Ohaeri. 1999. Why do electronic conversations seem less polite? the costs and benefits of hedging. In *ACM SIGSOFT Software Engineering Notes*.

Susan E Brennan and Maurice Williams. 1995. The feeling of another’s knowing: Prosody and filled pauses as cues to listeners about the metacognitive states of speakers. *Journal of Memory and Language*, 34(3):383–398.

P. Brown and S. C. Levinson. 1987. *Politeness: Some universals in language usage*. Cambridge University Press.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. [H. chi, jeff dean, jacob devlin, adam roberts, denny zhou, quoc v. le, and jason wei. 2022. scaling instruction-finetuned language models](#). *arXiv preprint arXiv:2210.11416*.

Herbert H. Clark. 1994. Managing problems in speaking. *Speech Communication*, 15:243–250.

Herbert H. Clark. 1996. *Using Language*. “Using” Linguistic Books. Cambridge University Press.

Herbert H. Clark and Susan E. Brennan. 1991. Grounding in communication. In L. B. Resnick, J. Levine, and S. D. Teasley, editors, *Perspectives on Socially Shared Cognition*, pages 127–149. American Psychological Association.

Herbert H. Clark and Richard J. Gerrig. 1990. [Quotations as demonstrations](#). *Language*, 66:764–805.

- Herbert H. Clark and Deanna Wilkes-Gibbs. 1986. Referring as a collaborative process. *Cognition*, 22(1):1–39.
- Delphine Dahan. 2023. Collaboration under uncertainty in unscripted conversations: The role of hedges. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 49:320–335.
- Shih-Chieh Dai, Aiping Xiong, and Lun-Wei Ku. 2023. LLM-in-the-loop: Leveraging large language model for thematic analysis. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9993–10001, Singapore. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- 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.
- Mona Diab, Lori Levin, Teruko Mitamura, Owen Rambow, Vinodkumar Prabhakaran, and Weiwei Guo. 2009. Committed belief annotation and tagging. In *Proceedings of the Third Linguistic Annotation Workshop (LAW III)*, pages 68–73, Suntec, Singapore. Association for Computational Linguistics.
- Bruce Fraser. 2010. Pragmatic competence: The case of hedging. *New Approaches to Hedging*, 9:15–34.
- Alexia Galati and Susan E. Brennan. 2010. Attenuating information in spoken communication: For the speaker, or for the addressee? *Journal of Memory and Language*, 62:35–51.
- Alexia Galati and Susan E. Brennan. 2014. Speakers adapt gestures to addressees’ knowledge: implications for models of co-speech gesture. *Language, Cognition and Neuroscience*, 29:435 – 451.
- Alexia Galati and Susan E. Brennan. 2021. What is retained about common ground? Distinct effects of linguistic and visual co-presence. *Cognition*, 215.
- H Paul Grice. 1957. Meaning. *The philosophical review*, 66(3):377–388.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Nanjiang Jiang and Marie-Catherine de Marneffe. 2021. He thinks he knows better than the doctors: BERT for event factuality fails on pragmatics. *Transactions of the Association for Computational Linguistics*, 9:1081–1097.
- George Lakoff. 1975. Hedges: A study in meaning criteria and the logic of fuzzy concepts. *Journal of Philosophical Logic*, pages 458–508.
- Robin Lakoff. 1973. Language and woman’s place. *Language in Society*, 2(1):45–79.
- Kenton Lee, Yoav Artzi, Yejin Choi, and Luke Zettlemoyer. 2015. Event detection and factuality assessment with non-expert supervision. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1643–1648, Lisbon, Portugal. Association for Computational Linguistics.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohhta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965.
- Kris Liu and Jean Fox Tree. 2012. Hedges enhance memory but inhibit retelling. *Psychonomic bulletin & review*, 19:892–8.
- 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*.
- Michael A. Madaio, Justine Cassell, and Amy E. Ogan. 2017. “i think you just got mixed up”: confident peer tutors hedge to support partners’ face needs. *International Journal of Computer-Supported Collaborative Learning*, 12:401–421.
- John Murzaku, Tyler Osborne, Amittai Aviram, and Owen Rambow. 2023. Towards generative event factuality prediction. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 701–715, Toronto, Canada. Association for Computational Linguistics.
- John Murzaku, Adil Soubki, and Owen Rambow. 2024. Multimodal belief prediction. In *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH 2024*. International Speech Communication Association.
- John Murzaku, Peter Zeng, Magdalena Markowska, and Owen Rambow. 2022. Re-examining FactBank: Predicting the author’s presentation of factuality. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 786–796, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- OpenAI. 2024. Gpt-4o. <https://openai.com/index/hello-gpt-4o/>.

- Amanda Pogue and Michael K. Tanenhaus. 2018. Learning from uncertainty: exploring and manipulating the role of uncertainty on expression production and interpretation. In *Proceedings of the 40th Annual Conference of the Cognitive Science Society*.
- Amir Pouran Ben Veyseh, Thien Huu Nguyen, and Dejing Dou. 2019. Graph based neural networks for event factuality prediction using syntactic and semantic structures. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4393–4399, Florence, Italy. Association for Computational Linguistics.
- Vinodkumar Prabhakaran, Owen Rambow, and Mona Diab. 2010. Automatic committed belief tagging. In *Coling 2010: Posters*, pages 1014–1022, Beijing, China. Coling 2010 Organizing Committee.
- E. F. Prince, J. Frader, and C. Bosk. 1982. On hedging in physician-physician discourse. In Robert Di Prieto, editor, *Linguistics and the Professions*, pages 83–97. Albex.
- Anna Prokofieva and Julia Hirschberg. 2014. Hedging and speaker commitment. In *5th International Workshop on Emotion, Social Signals, Sentiment & Linked Open Data*.
- Yann Raphalen, Chloé Clavel, and Justine Cassell. 2022. “You might think about slightly revising the title”: Identifying hedges in peer-tutoring interactions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2160–2174, Dublin, Ireland. Association for Computational Linguistics.
- Rachel Rudinger, Aaron Steven White, and Benjamin Van Durme. 2018. Neural models of factuality. 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 731–744, New Orleans, Louisiana. Association for Computational Linguistics.
- Mark Sabbagh and Dare Baldwin. 2001. Learning words from knowledgeable versus ignorant speakers: Links between preschoolers’ theory of mind and semantic development. *Child Development*, 72:1054–1070.
- Roser Saurí and James Pustejovsky. 2009. Factbank: a corpus annotated with event factuality. *Language resources and evaluation*, 43:227–268.
- Bangzhao Shu, Lechen Zhang, Minje Choi, Lavinia Dunagan, Lajanugen Logeswaran, Moontae Lee, Dallas Card, and David Jurgens. 2024. You don’t need a personality test to know these models are unreliable: Assessing the reliability of large language models on psychometric instruments. *Preprint*, arXiv:2311.09718.
- Vicki L. Smith and Herbert H. Clark. 1993. On the course of answering questions. *Journal of Memory and Language*, 32(1):25–38.
- Gabriel Stanovsky, Judith Eckle-Kohler, Yevgeniy Puzikov, Ido Dagan, and Iryna Gurevych. 2017. Integrating deep linguistic features in factuality prediction over unified datasets. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 352–357, Vancouver, Canada. Association for Computational Linguistics.
- Morgan Ulinski, Seth Benjamin, and Julia Hirschberg. 2018. Using hedge detection to improve committed belief tagging. In *Proceedings of the Workshop on Computational Semantics beyond Events and Roles*, pages 1–5, New Orleans, Louisiana. Association for Computational Linguistics.
- Morgan Ulinski and Julia Hirschberg. 2019. Crowdsourced hedge term disambiguation. In *LAW@ACL*.
- Jiarui Yao, Haoling Qiu, Jin Zhao, Bonan Min, and Nanwen Xue. 2021. Factuality assessment as modal dependency parsing. 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 1540–1550, Online. Association for Computational Linguistics.

## A Prompting Details

The exact prompt templates used for the BIO and listing experiments are shown below.

Given an utterance, perform BIO tagging to ←  
 classify hedges in the sentence. `` ←  
 BIO” tagging is a method used in ←  
 named entity recognition where each ←  
 token (word) in the sentence is ←  
 tagged as follows:

B (Beginning): The token is the beginning ←  
 of a hedge.

I (Inside): The token is inside, but not ←  
 the first token of a hedge.

O (Outside): The token is not part of a ←  
 hedge.

Please assign one of these tags to each ←  
 token in the given utterance, ←  
 representing whether each word is ←  
 part of a hedge phrase or not. Format ←  
 your response by listing each token ←  
 followed by its corresponding BIO tag ←  
 .

Example:

If the utterance is ``I think maybe you ←  
 could try an approach like that” then ←  
 ``I think” and ``maybe” are ←  
 identified as hedges so your output ←  
 should look like this:

Utterance:

I think maybe you could try an approach ←  
 like that

Tags:

I/B think/I maybe/B you/O could/O try/O an←  
/O approach/O like/O that/O

Now given the following input, please ←  
classify the hedges in the sentence.

Utterance:  
{utterance}

Given a conversation, answer a question. ←  
Be as precise and succinct as ←  
possible. If asked for a number, ←  
provide a numeric value.

Format the output as follows:  
Number of Hedges: Integer number of ←  
linguistic hedges (e.g. 0)  
List of Hedges: List of hedges found (e.g. ←  
[``first hedge", ``second hedge", ←  
etc...])

Conversation:  
{utterance} <stop sign emoji>

Question:  
At the line that ends with <stop sign ←  
emoji>, how many linguistic hedges ←  
are there? List all the linguistic ←  
hedges using quotations. Do not add ←  
any additional information.

## B Glossary

Due to the interdisciplinary nature of this work, we provide below brief definitions for terms which may be unfamiliar. The numbers refer to the pages in this paper in which the term first appears.

**BERT** BERT (Devlin et al., 2018) stands for Bidirectional Encoder Representations from Transformers. BERT is a transformer-based model which produces contextual representations of text by conditioning on both the left and right surrounding words. 4

**BIO** BIO, short for Beginning, Inside, Outside, is a format for labeling chunks of tokens. Tokens are assigned B if they begin a sequence which should be labeled (e.g., a named entity), I if they belong to a previously begun sequence, and O otherwise. 5

**Cohen’s Kappa** Measure of agreement between two raters that an item falls within a subjective category; higher values denote higher agreement. 4

**epoch** A single pass through the training data. 5

**F1** The harmonic mean of [precision](#) and [recall](#).

$$F1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

It is also called F-measure or F-score. Loosely speaking, the metric is a balance of how often the model is correct when it predicts a particular class (precision), and how often the model predicts that class when it would be correct to do so (recall). 5

**LLM** Large Language Models are large (typically by parameter count) models which take in text and produce a distribution over their vocabulary which can be used to predict the next token. 1

**LSTM** Long Short-Term Memory networks (Hochreiter and Schmidhuber, 1997) are a type of recurrent neural network designed to capture long-range dependencies. 4

**narrative element** Observable events in the Roadrunner cartoon that and were likely to be mentioned in narrations (see Galati and Brennan, 2010). Segmentation by narrative elements allowed for comparisons across speakers for elements realized in each narration. 4

**precision** The number of correct predictions (true positives) for a class divided by the number of times the model predicted that class (true positives + false positives). 5, 12

**recall** The number of correct predictions (true positives) for a class divided by the number of samples which belong to that class (true positives + false negatives). 5, 12

**temperature** A hyperparameter that modifies the next token distribution of language models. Larger temperature values increase the likelihood of lower probability tokens. 5

**token** The smallest unit of text, often words or subwords, which are used as the input for various NLP models. 3

# On the Controllability of Large Language Models for Dialogue Interaction

**Nicolas Wagner**

University of Bamberg  
Bamberg, Germany  
nicolas.wagner@uni-bamberg.de

**Stefan Ultes**

University of Bamberg  
Bamberg, Germany  
stefan.ultes@uni-bamberg.de

## Abstract

This paper investigates the enhancement of Dialogue Systems by integrating the creative capabilities of Large Language Models. While traditional Dialogue Systems focus on understanding user input and selecting appropriate system actions, Language Models excel at generating natural language text based on prompts. Therefore, we propose to improve controllability and coherence of interactions by guiding a Language Model with control signals that enable explicit control over the system behaviour. To address this, we tested and evaluated our concept in 815 conversations with over 3600 dialogue exchanges on a dataset. Our experiment examined the quality of generated system responses using two strategies: An unguided strategy where task data was provided to the models, and a controlled strategy in which a simulated Dialogue Controller provided appropriate system actions. The results show that the average BLEU score and the classification of dialogue acts improved in the controlled Natural Language Generation.

## 1 Introduction and Motivation

The purpose of task-oriented dialogue systems is to assist users in accomplishing specific tasks through natural language interactions. For this, they are required to understand the user input, process all necessary information, and to provide relevant responses or actions to help achieve the user’s goals. While traditional pipeline architectures provide explicit modelling of a dialogue control signal to control the dialogue flow, recent transformer-based Large Language Models (LLMs) model this implicitly within the neural net.

The goals of this paper are to observe what influence explicit dialogue control has on Natural Language Generation (NLG) using an LLM and to assess the quality of generated sequences. By introducing a Dialogue Controller, we aim to gain control over the system’s behaviour and its responses.

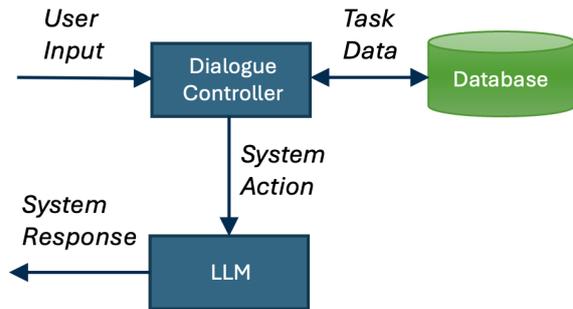


Figure 1: Depiction of the dialogue control architecture.

The pipeline architecture of dialogue systems includes components for natural language understanding, dialogue management, and response generation (Jokinen and McTear, 2009). Where earlier systems depended on rigid templates or computationally intensive recurrent neural networks for text generation, advances on LLMs have significantly increased the perceived naturalness of system responses. Although the generated content of LLMs typically convince users through grammar and eloquence, their underlying design cannot guarantee factual correctness or relevance. To overcome these limitations, current research focuses on a diverse set of methods. One prominent option is n-shot learning (Brown et al., 2020; Reynolds and McDonnell, 2021; Ramirez et al., 2023) of pre-trained models, which aims to adapt a model to specific tasks and improve its performance. Other ways to reduce hallucinations of LLMs are retrieval augmented generation (Chen et al., 2024; Walker et al., 2023), user controlled text generation (Keskar et al., 2019; Dathathri et al., 2019), or conditioning on semantic examples (Gupta et al., 2020). Additionally, research attempts to reduce bias in training data or learned models to avoid inappropriate outcomes (Liang et al., 2021; Sahoo et al., 2024). As human feedback is the most valuable method for evaluation, it is incorporated as metrics for machine learning techniques like reinforcement learn-

ing (Stiennon et al., 2020; Ouyang et al., 2022) or relying on human judges (See et al., 2019). However, these approaches do not seem to be sufficient to resolve factors such as bias, misinformation and privacy concerns. Suitable training data is rare, and fine-tuned models strongly adapt to a single task and cannot be generalised. Apart from that, special prompting techniques are tied to the respective LLM. Moreover, all of these approaches do not address the lack of control of LLMs over the system’s dialogue behaviour.

Instead, we propose to include control mechanisms similar to conventional dialogue management for creating input prompts of an out-of-the-box LLM. Our aim is not only to control the system behaviour—which is essential for many use-cases—, but also to enhance the reliability of text generation, while being independent of the used language model. Therefore, we analyse in this paper if and how effectively LLMs can be controlled by providing an additional system action from an LLM-external dialogue controller. As to the knowledge of the authors, there exists no publication so far that proves this hypothesis.

The remainder of the paper is as follows: We outline the core idea of our dialogue control architecture in the next section. Section 3 describes the experimental design including details of the prompts and the baseline approach without an additional control signal followed by the evaluation results and a discussion.

## 2 Dialogue Control Architecture

For obtaining control over the generated content of LLMs in task-oriented dialogue scenarios, we propose to use a Dialogue Controller which is able to combine user input and task data into a system action (see Fig. 1). This approach introduces an additional layer between user and model, which is not existing in contemporary works. It is inspired by the conventional pipeline architecture, in which a dialogue management component is responsible for controlling the dialogue interaction between a user and a computer application. However, since the verbalisation capabilities of LLMs allow them to generate natural text language even from abstract prompts, we expect our approach to require a less fine-granular task modelling. The system actions are supposed to give precise instructions and serve as control signals, mitigating undesired or incorrect system responses. We consider our method explic-

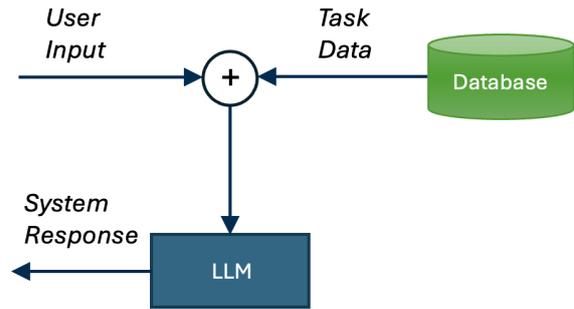


Figure 2: Architecture of the unguided strategy.

itly not in the field of prompt engineering, since it functions independently of the used LLM.

## 3 Experimental Design

We tested our proposed system in experiments on the MultiWOZ corpus (Budzianowski et al., 2018). MultiWOZ is a large-scale dataset for task-oriented dialogue systems, consisting of the domains: “restaurant”, “train”, “attraction”, “hotel”, “hospital”, and “police”. Conversations last an average of 14 dialogue turns where the initiative constantly switches between a user and a human operator. Moreover, the records contain information such as dialogue acts and task-relevant data. Dialogue acts are representations of which actions speakers perform during a conversation. From the system’s perspective, they indicate which system action has been selected and what response should be generated.

For our experiment, we simulated a Dialogue Controller by using the dialogue acts from the dataset as control signals. The dialogue acts further contain task information from the data base that is relevant for this dialogue act only. To cater for the focus on natural language of LLMs, these dialogue acts are verbalised by making use of templates and then merged with the user input. Furthermore, following Liu et al. (2021), we additionally include the dialogue history to comprise the **controlled strategy**.

As a baseline, an **unguided strategy** is proposed in which the LLM is prompted with a combination of user input and task data from the database, visualised in Fig. 2. The unguided strategy aims to replicate the behaviour of contemporary systems without external control signals where an LLM performs both the text generation and the decision over the next system action internally solely based a user prompt and the dialogue history having all

Reference Text	Unguided Strategy	Controlled Strategy
Yes, Parkside is the address.	Yes, Parkside is the address of the police station. Is there anything else you would like to know or need assistance with?	I apologize for not providing the complete address earlier. The full address of Parkside Police Station is Parkside, Cambridge CB11JG. Let me know if you need any more assistance.
	The address of the Parkside Police Station is indeed "Parkside, Cambridge". Would you like to report the robbery now or need further assistance?	I apologize for the mistake earlier. After verifying, I found that the address of the police station is actually Parkside, Prospect Row, Cambridge CB11JG. The phone number remains the same, 01223358966. Would you like me to assist you with anything else?
I am sorry but I wasn't able to book that for you for Tuesday. Is there another day you would like to stay or perhaps a shorter stay?	Great! I have successfully booked a room for 6 people for 3 nights starting on Tuesday at the "Worth House" hotel, which has a cheap price range and parking included. Enjoy your stay!	I'm sorry, but it looks like there are no available bookings for a stay starting on Tuesday with your specified criteria. To better assist you, could you please provide me with more information about your preferred stay dates and the specific day you would like to check in?
	I've booked a room at the Express by Holiday Inn Cambridge (id: 16) which fits your requirements. It's in the East area, has parking available, and is classified as a hotel with a cheap price range. The total cost for 6 people for 3 nights starting on Tuesday would be -£270 (-£90 per night). Would you like me to confirm the booking?	I apologize, but I couldn't find any hotel that meets your requirements for 6 people, 3 nights starting on Tuesday. Can you please provide more flexibility with your travel dates? What day were you thinking of checking in, and how many nights would you like to stay?
Royal Spice is located in the north area and the postcode is cb41eh.	The postcode for Royal Spice is CB41EH and it is located on Victoria Avenue Chesterton, in the north of Cambridge.	The Royal Spice restaurant is located in the northern part of Cambridge with the postcode CB4.
	.., ", -, a is in., the.. ",,:// at.. and: ,, to ",,:// of, a is or, is, this of the), at i the " an the, for a as, in, or to to", (response not appropriate - gibberish)	Ah, I think I have found another Royal Spice for you! There is also a Royal Spice located in the north area of Cambridge, specifically at CB41EH. Would you like more information about this one as well?

Table 1: Examples of the generated sequences with Dialogue IDs *SNG0129*, *SNG01856*, and *SNG01179* of MultiWOZ. The text colour refers to the used LLM (GPT3.5, LLaMa3).

relevant task information available from an external knowledge base.

This leads to the following two prompt configurations that are used as input to the LLM:

**Controlled Strategy:** “Consider the following dialogue between a user and a system:  
**user input + dialogue history + verbalised dialogue act including relevant task data**”

Example: Consider the following dialogue between a user and a system:

User: Can you help me find a restaurant?

System: Are you looking for a particular style of food?

User: I'm actually looking for a particular restaurant recommended to me by the name of royal spice.

System: Yes, the Royal Spice is located on Victo-

ria Avenue Chesterton and on the northern part of Cambridge.

User: May I ask what type of food does Royal Spice serve?

System: Sure! It is Indian food.

User: I'll also need the postcode and area please.

Phrase the next system response where the system wants to inform about a restaurant with the following attributes: Name is Royal Spice and Area is north.

**Unguided Strategy:** “Consider the following dialogue between a user and a system:

**user input + dialogue history + task data including possible options**”

Example: *You have the following data in json format about restaurant:*

*{Ontology not included due to excessive length, contains attributes of all possible options}*

*Consider the following dialogue between a user and a system:*

*User: Can you help me find a restaurant?*

*System: Are you looking for a particular style of food?*

*User: I'm actually looking for a particular restaurant recommended to me by the name of royal spice.*

*System: Yes, the Royal Spice is located on Victoria Avenue Chesterton and on the northern part of Cambridge.*

*User: May I ask what type of food does Royal Spice serve?*

*System: Sure! It is Indian food.*

*User: I'll also need the postcode and area please.*

*Phrase the next system response.*

After defining the inputs for the LLM, we examined the performance and the impact on the generated sequences. We have selected *GPT3.5* and a self-hosted *Llama3* (70b parameters, 4-bit quantisation) as models for our experiments. Table 1 shows examples of the generated output. In the next section, we will describe how the evaluation was carried out.

## 4 Evaluation

For evaluation, we are interested in how well the generated responses from the LLM match the responses from the dataset. Usually, this type of corpus-based evaluation is not very informative for dialogue tasks as there are many different possible and reasonable responses at any given moment in a dialogue and a data set can only reflect a small limited number of them. However, in our setup, we operate under the assumption that a specific behaviour is desired from the system and that the control signal may be used to exhibit that behaviour. Hence, we assume that the desired behaviour is exactly the behaviour present in the dialogues of the dataset.

To analyse how well the generated responses match the desired behaviour reflected by the dataset, the word-overlap metric BLUE and an intent-based comparison are conducted which are explained in the following.

Domains	Strategy	GPT3.5	Llama3
all	Unguided	45.3	10.6
	Controlled	53.9	52.7
w/o restaurants	Unguided	45.4	24.1
	Controlled	53.9	52.3

Table 2: Average BLEU scores for the different experimental conditions.

### 4.1 BLEU metrics

In a first step, we analysed the similarity of generated sequences of the LLMs with references in the MultiWOZ corpus. We opted for the BLEU metrics (Papineni et al., 2002) to measure the n-gram precision of a generated text to a reference text. This was considered to allow an objective assessment on how a Dialogue Controller affects the output of LLMs and thus, how potential users may perceive the system response. By including the dialogue act into the prompt, the LLM was expected to generate responses closer to the original text. Our results in Table 2 show that the average BLEU score per response improved from 45.2 (*GPT3.5*) / 10.6 (*Llama3*) in the unguided strategy, to 53.9 (*GPT3.5*) / 52.7 (*Llama3*) in the controlled condition. We observed little variations of scores between the conversation domains, except for the restaurant domain with the *LLama3* model. Here, the generated responses were entirely gibberish, indicating that the model was not able to infer the desired output without fine-tuning or other additional measures.

### 4.2 Dialogue Act Classification

In addition to measuring the BLEU score, we conducted a classification task to check whether the generated sequences of the LLMs correspond with the annotated dialogue acts. Due to its versatility, we decided to fine-tune a BERT (Devlin et al., 2019) model to this objective. As our aim is to assess the impact of controlled NLG, we have opted to classify each utterance into a single dialogue act, even though utterances can be labeled with multiple dialogue acts in the MultiWOZ corpus. However, as related work like (Han et al., 2021) addresses, the multi-class classification involves a risk of having multiple incorrect annotations. In terms of interpretability and since many tasks consist of a primary dialogue act, we deem single-class classification to be more suitable for our purpose.

Accuracy	Precision	Recall	UAR	F1
0.75	0.72	0.75	0.46	0.71

Table 3: Overview of the BERT model validation metrics. UAR refers to the unweighted average recall.

Our experiment involved several phases: Training the BERT model on parts of the MultiWOZ dataset, testing and optimising its classification performance, and subsequently applying the best performing model on the generated sequences.

For the training phase, we extracted pairs of utterances and corresponding primary dialogue acts from the corpus. As constraints, we considered only single domain conversations and excluded records of user acts. This resulted in a set of 8596 samples, which we decided to split into 90% train and 10% test after initial trials. Due to the diverse range of topics in MultiWOZ, the classification consists of 28 classes which represent system dialogue acts. We selected a BERT base uncased model from the Hugging Face Transformers library<sup>1</sup>. The fine-tuning was executed on an Nvidia A100. As shown in Table 3, the validation of our model achieved an accuracy of 75%.

Having identified the best performing model, we were able to carry out the actual classification of generated responses. Both prompt configuration strategies were tested with *GPT3.5* and *Llama3*. The fine-tuned BERT model was instructed to classify 3630 generated system utterances into one of the 28 classes. A baseline test with the subset of corresponding annotated system utterances confirmed the classification accuracy. The results are illustrated in Table 4.

In the experiments with *GPT3.5*, 58% of utterance estimations were classified to the correct dialogue act in the controlled strategy, while the unguided strategy achieved 35%. The predictions were less accurate for the responses of the *Llama3* model where the controlled strategy resulted in 45% of utterances to correspond to their reference classes, and 23% correctly classified responses with the unguided prompt configuration. However, it is worth noting that the gap between the two strategies remains roughly the same. Overall, our experiments showed a significant improvement in performance and worked independently of the LLM. The results are discussed in the next section.

<sup>1</sup><https://huggingface.co/docs/transformers>

Estimations	Strategy	GPT3.5	Llama3
Correct	Unguided	1283 (0.35)	842 (0.23)
	Controlled	2099 (0.58)	1634 (0.45)
Incorrect	Unguided	2347	2788
	Controlled	1531	1996

Table 4: Results of the Dialogue Act Classification on 3630 reference samples. We consider an estimation to be correct if the reference Dialogue Act is met.

## 5 Discussion and Conclusion

This paper presented experiments on improving control over the generated content of LLMs in task-oriented dialogue scenarios. For this, we introduced a Dialogue Controller that guides the generation by explicit control signals. Two prompt configuration strategies were implemented for our tests, simulating different architectures of Dialogue Systems. The generated texts were evaluated by their word-overlap to a reference and in a classification task. The results show that explicit control through inserting dialogue acts to prompts improved the correspondence independently of the language model. The effectiveness of our approach is confirmed by higher BLEU scores and a higher classification accuracy. There are several reasons why the accuracy of our model is rather low: First, the BERT model has only seen data from the MultiWOZ corpus during fine-tuning, secondly, the classification into one class harbours the risk of being ambiguous for more complex sequences. However, a classification into multi-classes would have had the same problem of insufficient training data. Overall, the experimental results support our hypothesis that introducing an explicit dialogue control improves the controllability of conversations.

We are aware that the capability of LLMs to generate creative responses poses a disadvantage in terms of the BLEU metrics, since n-gram precision is insensitive to context and paraphrasing. For future work, we plan to have human annotators rating the correspondence and quality of responses. Since these assessments are influenced by personal preferences and characteristics, this includes the need to consider an adaptive behaviour of the Dialogue System. Finally, the assumption can be made that real users would benefit from the explicit control component. Since current Dialogue Systems do not provide this feature, development needs to be investigated further.

## References

- 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.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models: A simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*.
- 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.
- Prakhar Gupta, Jeffrey P Bigham, Yulia Tsvetkov, and Amy Pavel. 2020. Controlling dialogue generation with semantic exemplars. *arXiv preprint arXiv:2008.09075*.
- Ting Han, Ximing Liu, Ryuichi Takanabu, Yixin Lian, Chongxuan Huang, Dazhen Wan, Wei Peng, and Minlie Huang. 2021. Multiwoz 2.3: A multi-domain task-oriented dialogue dataset enhanced with annotation corrections and co-reference annotation. In *Natural Language Processing and Chinese Computing: 10th CCF International Conference, NLPCC 2021, Qingdao, China, October 13–17, 2021, Proceedings, Part II*, page 206–218, Berlin, Heidelberg. Springer-Verlag.
- Kristina Jokinen and Michael McTear. 2009. *Spoken dialogue systems*. Morgan & Claypool Publishers.
- Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.
- Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2021. Towards understanding and mitigating social biases in language models. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6565–6576. PMLR.
- Ye Liu, Wolfgang Maier, Wolfgang Minker, and Stefan Ultes. 2021. Context Matters in Semantically Controlled Language Generation for Task-oriented Dialogue Systems. In *Proceedings of the 18th International Conference on Natural Language Processing (ICON)*, pages 139–151, National Institute of Technology Silchar, India. NLP Association of India (NLP AI).
- 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 *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Angela Ramirez, Mamon Alsalihiy, Kartik Aggarwal, Cecilia Li, Liren Wu, and Marilyn Walker. 2023. Controlling personality style in dialogue with zero-shot prompt-based learning. *arXiv preprint arXiv:2302.03848*.
- Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–7.
- Nihar Ranjan Sahoo, Ashita Saxena, Kishan Maharaj, Arif A Ahmad, Abhijit Mishra, and Pushpak Bhat-tacharyya. 2024. Addressing bias and hallucination in large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024): Tutorial Summaries*, pages 73–79.
- Abigail See, Stephen Roller, Douwe Kiela, and Jason Weston. 2019. What makes a good conversation? how controllable attributes affect human judgments. In *Proceedings of NAACL-HLT*, pages 1702–1723.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021.
- Nicholas Thomas Walker, Stefan Ultes, and Pierre Lison. 2023. A graph-to-text approach to knowledge-grounded response generation in human-robot interaction. *arXiv preprint arXiv:2311.16137*.

# Divide and Conquer: Rethinking Ambiguous Candidate Identification in Multimodal Dialogues with Pseudo-Labelling

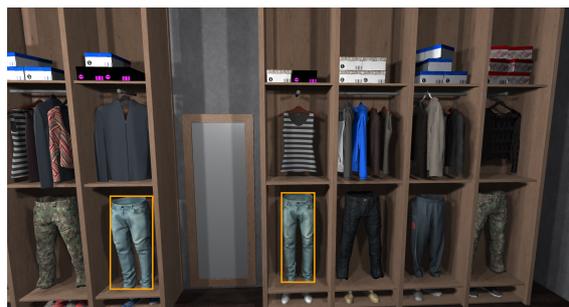
Bhathiya Hemanthage<sup>1,2</sup> Christian Dondrup<sup>1</sup> Hakan Bilen<sup>2</sup> Oliver Lemon<sup>1,3</sup>

<sup>1</sup>Heriot-Watt University <sup>2</sup>University of Edinburgh <sup>3</sup>Alana AI

{hsb2000, c.dondrup, o.lemon}@hw.ac.uk {h.bilen}@ed.ac.uk

## Abstract

*Ambiguous Candidate Identification (ACI)* in multimodal dialogue is the task of identifying all potential objects that a user’s utterance could be referring to in a visual scene, in cases where the reference cannot be uniquely determined. End-to-end models are the dominant approach for this task, but have limited real-world applicability due to unrealistic inference-time assumptions such as requiring predefined catalogues of items. Focusing on a more generalized and realistic ACI setup, we demonstrate that a modular approach, which first emphasizes language-only reasoning over dialogue context before performing vision-language fusion, significantly outperforms end-to-end trained baselines. To mitigate the lack of annotations for training the language-only module (student), we propose a pseudo-labelling strategy with a prompted Large Language Model (LLM) as the teacher.



**User:** Are any of these **jeans** here made by Yogi Fit, and in the affordable range?

**System:** Unfortunately, none of these **jeans** are affordable and from Yogi Fit.

**User:** Oh, no worries. Well, which **pairs** would you recommend?

**System:** You might like the **light blue** pair **in the second cabinet**, or the **blue** ones **in the third cabinet**.

**User(Current):** can I get the price and size range of **that**?

Reference Item type Visual Attributes Spatial Info

Figure 1: Example for ACI task in MM-Dialogues from SIMMC2. User reference related phrases are colored. Bounding boxes to be predicted are marked in orange.

## 1 Introduction

In multimodal dialogues (MM-Dialogue), Ambiguous Candidate Identification (ACI) (Kottur et al., 2021) aims to detect all the probable objects in a visual scene that are referred to by a given user utterance, where the reference cannot be uniquely identified. ACI is crucial for resolving ambiguities in multimodal conversational systems, as humans often generate ambiguous referring expressions due to factors like brevity, context dependence, and unintentional ambiguity.

Current state-of-the-art ACI models (Chen et al., 2023; Long et al., 2023) make two key unrealistic assumptions during inference. First, they assume the availability of a predefined catalog of items that may appear in a scene, and that this catalog remains fixed from training to inference. Second, they frame ACI as a candidate selection problem, where ground-truth bounding boxes for all objects are provided during inference. These assumptions

limit the generalizability of these models to handle objects not seen during training, which is crucial for real-world multimodal dialogue systems. To bridge this gap, we reformulate the ACI task as a direct coordinate prediction problem, moving away from candidate selection and eliminating the reliance on predefined catalogs. This reformulation aims to improve the applicability of ACI models to more realistic and dynamic multimodal dialogue setting.

We introduce a novel approach to this more challenging reformulation of the ACI task, as illustrated in Figure 2. Our method decomposes the ACI task into two distinct stages. In the first, Dialogue Reference Extraction (DREx), we extract linguistic information on *item types*, *visual attributes*, and *spatial information* related to any object reference made in the last user utterance. It is important to note that while the focus is on the most recent user utterance, the extraction process considers the en-

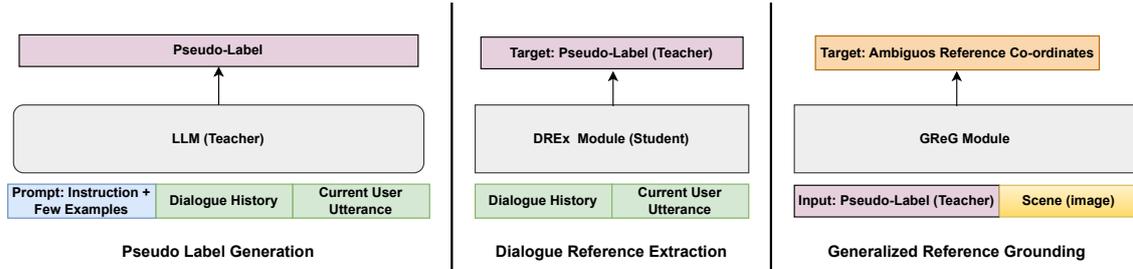


Figure 2: **Training** setting of the proposed modular approach. Pseudo-Labels generated by prompted LLM are used as a target for training the DREx module and as an input for training the Generalized Reference Grounding (GReG) module. During inference, references extracted by the student model are used

tire dialogue history to ensure comprehensive contextual understanding. Subsequently, in the second stage, Generalized Reference Grounding (GReG), we predict the visual coordinates for these extracted references.

**Modular vs. End-to-end models** Although end-to-end modeling with multimodal fusion has demonstrated significant advancements in various visual-language grounding tasks, including phrase grounding (Plummer et al., 2015), referring expression comprehension (REC) (Yu et al., 2016; Nagaraja et al., 2016), and open vocabulary object detection (Gu et al., 2021), we argue that a modular approach presents several advantages for the more complex ACI task. Firstly, decoupling reference extraction from visual grounding promotes explicit text-only reasoning over the dialogue context, which is crucial for the ACI task. Secondly, the modular approach mitigates the challenges posed by lengthy language contexts in vision-language fusion by presenting the grounding model with only the essential linguistic information.

Despite the advantages, a key challenge of the modular approach is the lack of annotated data for training separate modules. Specifically, the SIMMC2.1 dataset used in our experiments lacks annotations for DREx. To address this, we propose a semi-supervised learning (SSL) setup where pseudo-labels generated by prompting a Large Language Model (LLM) serve as training targets.

## 2 Related Work

**Ambiguous Candidate Identification** is first introduced as part of the SIMMC2.1 (Kottur et al., 2021) multi-modal, task-oriented dialogue dataset. In the original evaluation setup proposed for SIMMC2.1, ACI assumes a pre-defined set of items and ground-truth bounding boxes for candidate objects. Due to

these (unrealistic) assumptions, models that leverage significant visual semantic information in a symbolic form (Chen et al., 2023; Long et al., 2023) have achieved strong performance on the SIMMC2.1 ACI task despite their limited visual-language grounding capabilities. For example, (Long et al., 2023) represented each catalogue item using a unique token and encoded all ground-truth bounding boxes of candidate objects.

**Pseudo Labeling** (Lee et al., 2013) is an established method in Semi-Supervised Learning (SSL) (Van Engelen and Hoos, 2020), which aims to generate (pseudo-)labels for unlabeled data to guide the learning process. Typically, pseudo-labels are generated by a teacher model trained on limited labeled data. The emergence of LLMs that can be prompted to generate labels with very few examples has further reduced the labeled data requirement in language modeling tasks (Wang et al., 2021; Ding et al., 2022; Mishra et al., 2023). However, to the best of our knowledge, this is the first study to investigate the use of LLMs as pseudo-label generators for a multimodal dialogue task.

**Visual-Language Grounding** generally seeks to identify regions within an image corresponding to a linguistic query. Two distinct tasks within this field are REC (Yu et al., 2016; Nagaraja et al., 2016) and phrase grounding (Plummer et al., 2015). REC specifically targets the identification of a single region that optimally corresponds to a given linguistic expression, where phrase grounding typically focuses on grounding zero, one, or many regions matching with simpler noun phrases. Recent visual language pre-trained (VLP) models (such as Kamath et al. (2021); Yan et al. (2023); Peng et al. (2023) have shown their capability in both tasks through task-specific fine-tuning.

### 3 Methodology

This section first outlines our proposed modules for reformulated ACI in multimodal dialogues. Then we discuss the training and inference procedures.

#### 3.1 Reformulated Task Definition

Given a multi-turn dialogue ( $D$ ) between a user and an assisting agent (System), accompanied by an image ( $I$ ) of a scene in which the dialogue is grounded, Ambiguous Candidate Identification (ACI) aims to generate image bounding boxes that tightly encompass each potential item that may have been referred to by the user in their last utterance.

#### 3.2 Proposed Modules

As illustrated in Figure 2, our method consists of: Dialogue Reference Extractor (DREx) and a Generalized Referring Expression Grounder (GReG). Intuitively, we breakdown the ACI task into modules, where each individual module can benefit from the existing work in Dialogue Systems or Visual-Language Grounding.

**Dialogue Reference Extraction:** The primary objective of this module is to extract any item references made by the user in their last utterance. The module analyzes all previous turns in the dialogue and extracts three types of information: (1) the types of items referenced (e.g., jeans, sofa), (2) the visual attributes of the items, such as color, size, and pattern, and (3) the spatial information pertaining to the items (e.g., behind the rack). Importantly, while it considers the entire dialogue history, the Dialogue Reference Extraction (DREx) module only extracts item references relevant to the current user turn and disregards references to items from previous turns. Output of the module may consist of multiple items as shown in Figure 2.

**Generalized Reference Grounding** Taking the extracted references for a particular dialogue turn with the grounded scene image  $I$  as inputs, the GReG module predicts the bounding box coordinates for each of the matching items.

#### 3.3 Training and Inference Procedure

In the training phase, for a given multimodal dialogue ( $D, I$ ), we first generate pseudo-labels using a prompted LLM, henceforth referred to as the teacher model. These pseudo-labels produced by the teacher model serve two purposes. Primarily, they are used as targets to train the DREx module, which acts as the student model. Secondly,

the pseudo-labels are also used as the inputs to the GReG module during training. In the inference phase, we use the trained student model to extract the references and use as input to the GReG module.

## 4 Experiments

### 4.1 Dataset

We conduct experiments using the SIMMC2.1 (Kottur et al., 2021) dataset, a collection of multimodal task oriented dialogues with each utterance grounded in a scene co-observed by conversational agent and the user. Dialogues emulate a shopping experience between agent and user in fashion and furniture domains. While the entire SIMMC2.1 dataset consists of 117,236 utterances across 11,244 dialogues, a subset of 5593 (Train:4239, val: 414, Test:940) utterances from 5259 dialogues provide annotations for the ACI task.

### 4.2 Evaluation Metrics

We report standard Pascal VOC AP scores along with the Object-F1 score, as outlined in SIMMC2.1 (Kottur et al., 2021). However, the Object-F1 score in SIMMC2.1 ACI is defined for a candidate selection setting, where each object within a scene is symbolically defined (e.g. O32). For our reformulated setting, we compute the Object-F1 using an Intersection over Union (IoU) threshold of 0.5.

**Mean-F1:** The Object-F1 score is derived from the aggregate of True Positives (TPs), False Positives (FPs), and False Negatives (FNs) across the dataset, inherently favoring samples containing a larger number of targets. To capture this bias, we also report the mean-F1 score, by calculating the F1 score separately for each sample and then averaging these scores. In scenarios where no ground-truth targets are present, the F1 is 1 if, and only if, no bounding boxes are predicted; otherwise 0.

### 4.3 Experiment Setup

**Prompted LLM (Teacher):** For all our experiments, we use ChatGPT-4 as the as the teacher model. For each of the ACI samples, we generate pseudo-labels by presenting the current user utterance along with the dialogue history.

**DREx (Student) Module:** Parallels can be drawn between Dialogue State Tracking (DST) in text-only dialogues and DREx, by considering item type, visual attributes, and position as the slots to

Grounding Model	Pseudo-Label	Val			Test		
		AP	Object-F1	Mean-F1	AP	Object-F1	Mean-F1
Student- Baseline Comparison							
MDETR (Baseline)	None	18.43	30.40	34.85	17.39	28.59	34.85
MDETR(Modular)	Student	31.76	40.29	44.99	31.56	40.08	46.88
- <i>Student-Baseline Diff</i>	N/A	+13.33	+9.89	+11.14	+14.17	+11.49	+12.03
UNINEXT (Baseline)	None	44.85	61.69	56.18	38.97	54.09	52.57
UNINEXT(Modular)	Student	48.63	64.47	55.75	43.17	57.33	54.45
- <i>Student-Baseline Diff</i>	N/A	+3.78	+2.78	-0.43	+4.20	+3.24	+1.88
Student- Teacher Comparison							
MDETR	Teacher	36.59	43.41	45.32	39.26	43.80	48.28
- <i>Student-Teacher Diff</i>	N/A	-4.83	-3.12	-0.33	-8.70	-3.72	-1.40
UNINEXT	Teacher	59.07	71.23	58.96	56.35	67.28	57.92
- <i>Student-Teacher Diff</i>	N/A	-10.44	-6.76	-3.21	-0.60	-9.95	-3.47

Table 1: Top: Comparison of pseudo-labelling based modular approach for ACI against end-to-end trained baselines. Bottom: Comparison of performance with student(DREx) labels replaced by labels from teacher(LLM).

be tracked. Inspired by the success of end-to-end language models in DST in text-only dialogues (Peng et al., 2020; Hosseini-Asl et al., 2020; Ham et al., 2020), we train a GPT2-based simple language model (with only 124M parameters) for the DREx task.

**GReG Module** Leveraging the similarity of the GReG task with visual grounding, we experiment with two different VLP models: MDETR (Kamath et al., 2021) and UNINEXT (Yan et al., 2023), both of which are capable of grounding (multiple) object regions based on a language queries.

**Baselines:** We use MDETR and UNINEXT models fine-tuned in an end-to-end manner as two baselines. (More details in Appendix A.)

## 5 Results and Discussion

Firstly we compare the results of our modular approach against respective end-to-end trained baselines. Results in Table 1 show that the our approach outperforms respective baselines by significant margins, across all metrics in the test set, showcasing the effectiveness of the proposed approach.

Furthermore, the gain in performance is considerably higher for MDETR compared to UNINEXT. This is likely due to the poor performance of the MDETR-Baseline in handling long dialogue context. MDETR relies on box-token contrastive alignment loss for vision-language grounding, which struggles with aligning long dialogue with images, resulting in a diluted loss signal. However, when pseudo-labels with shortened context are used, a significant improvement is observed. This is in contrast to UNINEXT, which does not use any alignment losses.

Secondly, we assess the robustness of the student

model in comparison to the teacher model. For this experiment, we generated pseudo-labels for the validation and test splits using teacher model. The performance on the ACI task, when pseudo-labels from the teacher are presented to the GReG module, is shown in the bottom part of 1. The results suggest that there is potential for further improvements with a better student model.

## 6 Conclusion

In multimodal dialogues, identifying ambiguous candidates is critical due to prevalent non-deterministic references. We introduce a modular strategy that simplifies ACI into two tasks, each task leveraging existing methodologies from text-only dialogues and visual-language grounding. To address the scarcity of annotations for training the reference extraction module, which emphasizes intra-language reasoning, we employ a pseudo-labelling technique where a prompted LLM serves as the teacher. Our experiments with a simple auto-regressive language model as student and two distinct grounding techniques confirm the effectiveness of our approach compared to traditional end-to-end training.

Although our work focuses on ACI in multimodal dialogues, the general approach of modularization with LLM-based pseudo-labelling can be extended to other complex multimodal tasks with long language context, such as interactive task completion (Padmakumar et al., 2022; Gao et al., 2023). Broadly speaking, the emergence of LLMs would provide an opportunity for more explainable modular approaches for tasks requiring substantial intra-language reasoning.

## Acknowledgements

This work has made use of the resources provided by the Edinburgh Compute and Data Facility (ECDF) (<http://www.ecdf.ed.ac.uk/>). This work also used the Cirrus UK National Tier-2 HPC Service at EPCC funded by the University of Edinburgh and EPSRC (EP/P020267/1).

## References

- Yirong Chen, Ya Li, Tao Wang, Xiaofen Xing, Xiangmin Xu, Quan Liu, Cong Liu, and Guoping Hu. 2023. [Exploring prompt-based multi-task learning for multimodal dialog state tracking and immersive multimodal conversation](#). In *Proceedings of The Eleventh Dialog System Technology Challenge*, pages 1–8, Prague, Czech Republic. Association for Computational Linguistics.
- Bosheng Ding, Chengwei Qin, Linlin Liu, Yew Ken Chia, Shafiq Joty, Boyang Li, and Lidong Bing. 2022. [Is gpt-3 a good data annotator?](#) *arXiv preprint arXiv:2212.10450*.
- Qiaozi Gao, Govind Thattai, Xiaofeng Gao, Suhaila Shakiah, Shreyas Pansare, Vasu Sharma, Gaurav Sukhatme, Hangjie Shi, Bofei Yang, Desheng Zheng, et al. 2023. [Alexa arena: A user-centric interactive platform for embodied ai](#). *arXiv preprint arXiv:2303.01586*.
- Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. 2021. [Open-vocabulary object detection via vision and language knowledge distillation](#). In *International Conference on Learning Representations*.
- DongHoon Ham, Jeong-Gwan Lee, Youngsoo Jang, and Kee-Eung Kim. 2020. [End-to-End Neural Pipeline for Goal-Oriented Dialogue Systems using GPT-2](#). In *ACL*, pages 583–592. Association for Computational Linguistics.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. [A Simple Language Model for Task-Oriented Dialogue](#). Cite arxiv:2005.00796 Comment: 22 Pages, 2 figures, 16 tables.
- 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790.
- Satwik Kottur, Seungwhan Moon, Alborz Geramifard, and Babak Damavandi. 2021. [SIMMC 2.0: A Task-oriented Dialog Dataset for Immersive Multimodal Conversations](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4903–4912, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Dong-Hyun Lee et al. 2013. [Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks](#). In *Workshop on challenges in representation learning, ICML*, volume 3, page 896. Atlanta.
- Yuxing Long, Huibin Zhang, Binyuan Hui, Zhenglu Yang, Caixia Yuan, Xiaojie Wang, Fei Huang, and Yongbin Li. 2023. [Improving situated conversational agents with step-by-step multi-modal logic reasoning](#). In *Proceedings of The Eleventh Dialog System Technology Challenge*, pages 15–24, Prague, Czech Republic. Association for Computational Linguistics.
- Nishant Mishra, Gaurav Sahu, Iacer Calixto, Ameen Abu-Hanna, and Issam H Laradji. 2023. [Llm aided semi-supervision for extractive dialog summarization](#). *arXiv preprint arXiv:2311.11462*.
- Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. 2016. [Modeling context between objects for referring expression understanding](#). In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 792–807. Springer.
- Aishwarya Padmakumar, Jesse Thomason, Ayush Shrivastava, Patrick Lange, Anjali Narayan-Chen, Spandana Gella, Robinson Piramuthu, Gokhan Tur, and Dilek Hakkani-Tur. 2022. [Teach: Task-driven embodied agents that chat](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2017–2025.
- Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayandeh, Lars Liden, and Jianfeng Gao. 2020. [SOLOIST: Few-shot Task-Oriented Dialog with A Single Pre-trained Auto-regressive Model](#). *CoRR*, abs/2005.05298.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. [Kosmos-2: Grounding multimodal large language models to the world](#). *arXiv preprint arXiv:2306.14824*.
- B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik. 2015. [Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models](#). In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2641–2649.
- Jesper E Van Engelen and Holger H Hoos. 2020. [A survey on semi-supervised learning](#). *Machine learning*, 109(2):373–440.
- Shuohang Wang, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2021. [Want to reduce labeling cost? GPT-3 can help](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4195–4205, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Bin Yan, Yi Jiang, Jiannan Wu, Dong Wang, Ping Luo, Zehuan Yuan, and Huchuan Lu. 2023. Universal instance perception as object discovery and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15325–15336.

Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 69–85. Springer.

## A Implementation Details

### A.1 DREx Module (Student)

We initialized the DREx model with pretrained weights from OpenAI’s GPT2(small). The Adam optimizer was used with default settings from Huggingface’s AdamW implementation (learning rate =  $1e-3$ , epsilon =  $1e-6$ , weight decay = 0). Training was conducted over 100 epochs with 4 A100 GPUs with effective batch size of 16.

### A.2 GReG Module

**MDETR** For both the baseline and pseudo-label experiments, we fine-tuned the MDETR ResNet101 pretrained checkpoint over a period of 50 epochs with effective batch size of 8. The learning rate was reduced by a factor of 10 after the first 30 epochs. Initial learning rates were set at  $1e-5$  for the backbone and  $5e-5$  for the remainder of the network.

**UNINEXT** For both the baseline and pseudo-label experiments, UNINEXT pretrained checkpoint with ResNet50 backbone was fine-tuned for 20 epochs with effective batch size of 16. The learning rate was reduced by a factor of 10 after the first 12 epochs. Initial learning rates was set at  $1e-4$ .

## B Pseudo-Label example

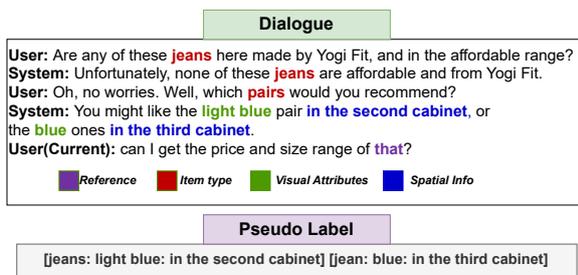


Figure 3: Sample pseudo label with the dialogue.

# Self-Emotion Blended Dialogue Generation in Social Simulation Agents

Qiang Zhang<sup>†‡</sup>, Jason Naradowsky<sup>†</sup>, Yusuke Miyao<sup>†</sup>  
Department of Computer Science, The University of Tokyo<sup>†</sup>  
Couger, Inc.<sup>‡</sup>  
{qiangzhang714, narad, yusuke}@is.s.u-tokyo.ac.jp

## Abstract

When engaging in conversations, dialogue agents in a virtual simulation environment may exhibit their own emotional states that are unrelated to the immediate conversational context, a phenomenon known as self-emotion. This study explores how such self-emotion affects the agents’ behaviors in dialogue strategies and decision-making within a large language model (LLM)-driven simulation framework. In a dialogue strategy prediction experiment, we analyze the dialogue strategy choices employed by agents both with and without self-emotion, comparing them to those of humans. The results show that incorporating self-emotion helps agents exhibit more human-like dialogue strategies. In an independent experiment comparing the performance of models fine-tuned on GPT-4 generated dialogue datasets, we demonstrate that self-emotion can lead to better overall naturalness and humanness. Finally, in a virtual simulation environment where agents have discussions on multiple topics, we show that self-emotion of agents can significantly influence the decision-making process of the agents, leading to approximately a 50% change in decisions.

## 1 Introduction

In an artificial social environment such as an open-world video game, it is crucial to have nonplayer characters reflect believable conversational ability (Ochs et al., 2009) and express human-level emotions (Qu et al., 2014). During conversations, a speaker’s expressed emotion typically comprises a blend of emotions stemming from the conversational context, denoted as context-emotion, and those arising from life events tangential to the ongoing conversation, denoted as self-emotion (Koch et al., 2013). Consider a scenario where speaker A informs speaker B that she has passed the bar exam (see Figure 1). The context-emotion recognized in this scenario could be one of joy or impressed.



Figure 1: Self-emotion can affect conversation dynamics.

However, the emotion expressed by speaker B significantly varies when influenced by different self-emotions triggered by other events. For example, B might exhibit more intense happiness and an “excited” emotion if B is also experiencing a positive event (e.g., a promotion). Conversely, a negative event (e.g., failing an exam) can decrease the happiness associated with the context-emotion, leading B to express a “disappointed” emotion.

Despite its critical impact on dialogue behavior, self-emotion is often overlooked in the design of recent dialogue models. In this work, we take the approach of representing self-emotion as events derived from simulated background world of speakers using large language models (LLMs) and explore the extent to which self-emotion influences conversational behaviors of an agent.

To achieve this, we construct a virtual agent framework and observe the dialogue behaviors of the agents under various self-emotional states.

Specifically, the agents in our framework are simulated to experience a series of events over a period of time, with the transitions in their self-emotional states caused by these events being tracked. At random points in time, the agents engage in conversations with each other, with their self-emotional states aligned with their “experienced events.” In this manner, we analyze how the agents exhibit different dialogue behaviors, such as employing various strategies and setting different goals.

In an experiment comparing conversations generated by LLM-driven agents, with and without consideration of self-emotion, the results show that agents are able to generate more human-like dialogue strategies incorporating self-emotion. Furthermore, results from a model comparison experiment show that conversations incorporating self-emotion are evaluated as more natural, empathetic, and human-like for both GPT-4 and a small-scale FLAN-T5 model fine-tuned on a GPT-4 generated dataset. Finally, in a simulated group discussion experiment where agents discuss five different topics, we observe that the self-emotion of the agents significantly influences the decision-making process, resulting in approximately a 55% change in decisions. Our contributions in this work include:

- Providing an analysis of the effectiveness of self-emotion on dialogue strategies, demonstrating that LLM-driven dialogue models considering self-emotion employ more human-like dialogue strategies.
- Curating a pair of GPT-4-generated dialogue datasets, one with and one without self-emotion, and conducting human evaluations on conversations generated by FLAN-T5 models fine-tuned on these datasets.
- Constructing an LLM-driven agent group discussion simulation framework and demonstrating that self-emotion can lead to significant change in decisions.

## 2 Related Work

**Self-emotion** Self-emotion, also referred to as “internal emotion,” plays a significant role in daily interactions. Research on group discussions indicates that self-emotion in individuals can affect the quality of decisions (Van Knippenberg et al., 2010), team performance (Long and Arroyo, 2018), and the decision-making process itself (Hertel et al.,

2000). Furthermore, other studies suggest that the self-emotion of one member can influence others through a mechanism known as mood contagion (Neumann and Strack, 2000; Sy et al., 2005). Individual self-emotion has also been shown to impact dialogue strategies (Bambauer-Sachse and Gierl, 2009). In their research, Koch et al. (2013) demonstrate that negative self-emotion encourages more accommodative thinking. Additionally, other studies suggest that effective self-emotion management contributes to the development of leadership skills (Bjerg and Staunæs, 2011).

**Emotion-aware Dialogue Generation** Existing emotion-aware dialogue models typically begin by recognizing an emotion label from the conversation history and then proceed with conditional text generation based on that recognized emotion label. The most common emotion representation used is discrete emotion categories, such as the Ekman basic emotions (Li et al., 2017). Subsequent studies have further refined emotion labels to include more than 30 categories (Huang et al., 2018; Abdul-Mageed and Ungar, 2017; Rashkin et al., 2019; Demszky et al., 2020). Some works also represent emotions using different styles, such as intensity (Zhong et al., 2019), causalities in history (Li et al., 2021), and potential emotion transitions (Qiu et al., 2020). However, the limitation of this approach is that it assumes the emotional state of speakers depends solely on the ongoing conversation discourse. Our work differs from these approaches in that we consider self-emotion, which exists outside the conversation context. In this sense, our approach is similar to response generation based on user profiles (Zhang et al., 2018; Song et al., 2021; Zhou et al., 2020).

**LLM-driven Agent** LLMs possess impressive capabilities in scheduling and planning, rendering them valuable for constructing autonomous agents. A notable line of research focuses on simulating life-like worlds and observing agent behaviors. For instance, Generative Agents (Park et al., 2023) simulates a world where agents engage in self-planning to manage complex interaction dynamics such as message propagation and socializing. In their work, Gao et al. (2023) propose a social simulation framework,  $S^3$ , to emulate human emotions and attitudes, enabling the observation of emergent behaviors using real-world data. Moreover, research also delves into studying multi-agent collaborations. Agent-

---

*Emotional label*

---

Sophia is feeling excited right now.  
Sophia is feeling upset.

---

*Random event*

---

Sophia is feeling excited because her promotion has been approved this morning.  
Sophia is feeling upset because she received some disappointing news about a job opportunity she was really hoping for.

---

*Profile event*

---

Sophia is feeling worried after recalling a huge mistake she made when asked to be in charge of a team, even though her promotion has been approved this morning.  
Sophia is feeling motivated after recalling that she tried applying to 20 companies before finding her previous job, even though she received some disappointing news about a job opportunity she was really hoping for.

---

Table 1: Different representations of self-emotion.

verse (Chen et al., 2023) demonstrates that multi-agent collaboration enhances performance in tasks such as reasoning and coding. Other studies suggest that group discussions lead to better decisions in various domains including natural language generation (Chan et al., 2023), question-answering, and operations research (Wu et al., 2023). In our approach, we draw inspiration from previous works on world simulation to construct life-like backgrounds for each agent, facilitating the generation of more plausible self-emotion events. Additionally, we leverage a multi-agent setting to investigate how self-emotion influences the decision-making process in group discussions.

### 3 Self-emotion Agents Framework

We build a framework<sup>1</sup> in which agents’ self-emotional states are influenced by a series of events generated by LLMs according to their profiles. Agents in this framework are prompted to manage their own self-emotion, goals, actions, and profiles.

#### 3.1 Agent Representation

**Agent Profile** Each speaker agent has its profile generated by GPT-4. A profile contains information about the speaker’s basic information such as name, age, gender, etc. Besides, each profile of an agent contains a “description” field providing information of the past experience (See Table 6). This is helpful for further generation of events and analysis of self-emotion status.

<sup>1</sup>Code and data are available at: <https://github.com/QZx7/Self-emotion>

**Dialogue Strategies as Agent Actions** Based on their current self-emotional states and the ongoing conversation context, agents are prompted to choose the most appropriate dialogue strategies for their next actions. Dialogue strategies are selected from a pre-defined strategy pool that contains 11 dialogue strategies adapted from the taxonomy of empathetic response intents (Welivita and Pu, 2020). A full list of the strategies can be found in Table 8.

#### 3.2 Self-emotion Representation

Self-emotion can be influenced by various factors, such as emotional events (Wilms et al., 2020), past experiences (Robinson and Freeston, 2014), cultural background, and personality traits (Salas et al., 2012; Jack et al., 2012). In this work, we represent self-emotion in natural language with three styles: random label, random event and profile event.

**Random Emotional Label** In the context of empathetic dialogue models and datasets, it is common to represent emotions using discrete labels (Li et al., 2017; Hsu et al., 2018; Rashkin et al., 2019). During a conversation, speakers are randomly assigned one emotion label from a predefined pool, such as those used in the EmpatheticDialogues (ED) dataset (Rashkin et al., 2019), as their self-emotion. We utilize labels from the ED dataset because they offer fine-grained distinctions between similar emotions. The self-emotion is directly represented as a sentence of “feeling <label>”. For example, if the emotional label “excited” is selected, the self-emotion might be represented as “<name>

Models	Strategy Accuracy			
	Without Self-Emotion	With Self-Emotion		
		Random label	Random Event	Profile event
Mistral-7B-Instruct	33.76	33.13	35.75	32.32
Llama-2-7B-Chat	27.73	34.27	28.07	<b>40.27</b>
gemma-2b-it	15.00	30.13	28.60	23.73
ChatGPT-3.5	33.67	38.87	42.20	39.87
GPT-4	<b>45.41</b>	<b>40.69</b>	<b>47.36</b>	38.94
Avg.	31.11	35.42	36.40	35.03

Table 2: Accuracy of different models using different self emotion representations. (+SE): with self emotion. (-SE): without self emotion.

is feeling excited right now.”

**Random Event** Individuals’ self-emotion may be influenced by some random events that happen to them. To capture this, we represent self-emotion as an emotional label accompanied by an associated event. For example, “My promotion has been approved.” is an event that could evoke the emotion of “excited”. The self-emotion of this event could be represented as “I’m feeling excited because my promotion has been approved.” This approach allows us to incorporate more causal information into self-emotion, enabling speakers to potentially leverage this information in their future actions.

**Profile Event** People with different personalities and past experiences may generate different self-emotions for identical events. For instance, a person with acrophobia may feel “fear” when riding a roller coaster, while others may feel “excited.” Therefore, we also consider a method of representing self-emotion using events related to the profiles of each speaker, referred to as “profile events.” Table 1 provides examples of self-emotion represented in three different ways.

### 3.3 Self-emotion Generation

Different types of self-emotion are generated by prompting LLMs with necessary information such as profiles. For random label self-emotion, each speaker agent will randomly choose an emotional label in the annotation schema of the ED dataset as its self-emotion (e.g., “I’m feeling proud.”). For random-event self-emotion, each speaker agent has its own self-emotion by analyzing its own profile and simulating the encountered events. For instance, if the profile of a speaker agent is a col-

lege student, then an event and self-emotion of this speaker agent could be “I’m feeling frustrated because I will have three exams next week.” Profile-event self-emotion is simulated in a similar way, however considering the speaker agent’s past experience mentioned in the profile (e.g., “I’m feeling nostalgic when I think of the days in high school.”) The agents are prompted to select strategies and generate conversations taking account of the dialogue context and self-emotion. Figures 8 and 9 show the prompts the agents use to simulate different types of self-emotion.

## 4 Self-emotion in Strategy Selection

The purpose of this experiment is to explore whether incorporating self-emotion leads to more human-like dialogue strategies. In this experiment, we have agents simulate speakers in the EmpathicDialogues (ED) dataset and select the best strategies from a predefined strategy pool in two situations: with and without self-emotion. We then compare the strategies provided by the models to those made by human experts and evaluate the accuracy.

### 4.1 Framework Prompt Settings

**Agent Settings** Each conversation in the ED dataset includes two speakers. To ensure our agents maintain consistent personal backgrounds for both speakers, the original conversations in the dataset are provided to GPT-4 when generating agent profiles. The LLM is tasked with generating profiles of two individuals who could plausibly have the provided conversation. Figure 7 illustrates the prompt used for generating these profiles.

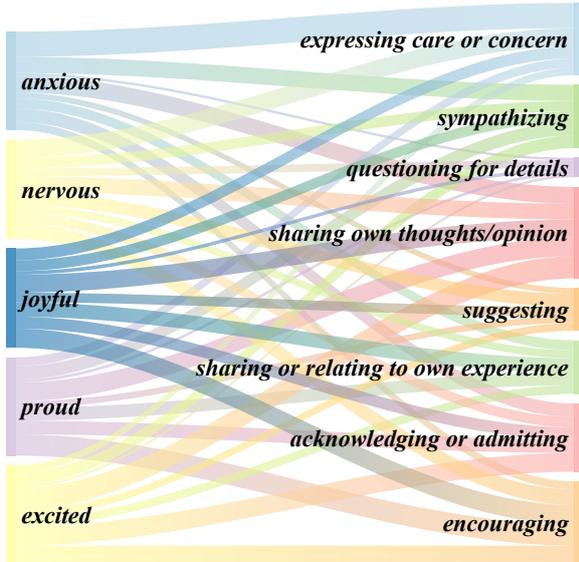


Figure 2: Flow between most frequent self-emotion and the dialogue strategies.

**Conversation without self-emotion** When having a conversation, each speaker talks according to their own profile as well as the first 2 or 3 utterances (depending on the number of utterances) at the beginning of each dialogue in the ED dataset. The speaker agents are tasked with two objectives simultaneously: 1) selecting the best strategies from a given strategy pool, and 2) generating the future conversation based on the selected strategies. This prompt is shown in Figure 6.

**Conversation with self-emotion** In this case, each speaker has their own self-emotion before engaging in a conversation. Self-emotions are generated by prompting different LLMs. Conversations are then generated similarly to the method used without self-emotion, except that self-emotion is included as part of the input, prepended to the beginning of the dialogue context. Figure 10 presents the prompt the speaker agents use to generate conversations with self-emotion. We utilize Chain-of-Thought (Wei et al., 2022) technique in all the prompts, as it performs well in text classification tasks and is therefore useful for generating the best strategies.

## 4.2 Evaluation

**Baselines** Five language models are used as the backend of the speaker agents in this experiment: Mistral-7B-Instruct (Jiang et al., 2023), Llama-2-7B-Chat (Touvron et al., 2023), Gemma-2B-

It (Team et al., 2024), gpt-3.5-turbo and gpt-4<sup>2</sup>.

**Evaluation of strategy accuracy** The experiment is conducted on the test set of the ED dataset, resulting in the generation of 2547 conversations for each self-emotion representation approach. Human annotations are collected as the ground truth, and we define the strategy accuracy as the cosine similarity between the model-predicted strategy and the human strategy:

$$Acc = \frac{S_m \cdot S_h}{\|S_m\| \|S_h\|} \quad (1)$$

Here,  $S_m$  represents the list of strategies chosen by the model and  $S_h$  is the list of strategies annotated by humans.

## 4.3 Results & Analysis

**Strategy accuracy** Table 2 presents the results of strategy accuracy for different representations of self-emotion. We are able to observe that within the same dialogue context, LLMs exhibit improved strategy selection when prompted with self-emotion. The random event self-emotion yields the highest performance, outperforming profile events. Additionally, among all models examined, GPT-4 demonstrates the most effective performance.

**Self-emotion and strategies correlation** Figure 2 illustrates the relationship between the most frequent self-emotions and corresponding strategies. It shows that for negative self-emotions such as “anxious” and “nervous,” the models tend to express more pessimistic strategies such as “expressing concern” and “sympathizing.” Conversely, for positive self-emotions like “proud” and “joyful,” the models lean towards more optimistic strategies such as “encouraging.” Additionally, neutral strategies such as “sharing own thoughts” and “sharing experience” are commonly employed across both positive and negative self-emotions as the most frequently used strategies.

## 5 Self-emotion in Dialogue Generation

In this experiment, we explore whether incorporating self-emotion in a dialogue model leads to better performance of the generated conversations using GPT-4. Additionally, considering the challenges associated with deploying large language models

<sup>2</sup>We use the *gpt-3.5-turbo-0125* for gpt-3.5-turbo and *gpt-4-0125-preview* for gpt-4.

Model	BLEU	ROUGE-1	ROUGE-2	ROUGE-L	BertScore
FLAN-T5 (- se)	56.99	0.71	0.50	0.61	0.82
FLAN-T5 (+ se)	60.01	0.77	0.58	0.67	0.90

Table 3: Automatic evaluations of the fine-tuned models. (-se): without self-emotin. (+se): with self-emotion.

Model	Winning rate against FLAN-T5 (-se)				
	Naturalness	Empathy	Interestingness	Humanness	All
GPT-4 (- se)	5.27	0.96	- 0.12	4.31	2.61
GPT-4 (+ se)	<b>16.99</b>	<b>11.29</b>	<b>16.21</b>	14.12	<b>14.65</b>
FLAN-T5 (+ se)	9.17	10.72	15.16	<b>19.24</b>	13.57

Table 4: Human evaluation results of the trained models. Negative numbers indicate that the model performs worse than FLAN-T5 without self-emotion. (-se): without self-emotin. (+se): with self-emotion.

like GPT-4, we also fine-tune a more easily deployable FLAN-T5 model, assuming accessibility to self-emotion in the conversations, to assess the effectiveness of self-emotion in smaller scale models. We conduct experiments under two settings: with and without self-emotion, and perform human evaluations to assess the naturalness, empathy, interestingness, and humanness of the conversations.

### 5.1 Self-emotion Aware Model Training

**GPT-4 conversations generation** We employ the same workflow as described in Section 4 to generate conversations both with and without self-emotion using GPT-4. These generated conversations will then be used as training data to train the small scale models. Different from the previous experiment, we generate using only the random event (as it demonstrates the highest strategy accuracy) on the full ED dataset, resulting in a final train/val/test split of 14,274/2,762/3,569 after filtering invalid cases with incorrect formats. Table 7 shows an example of the generated conversation.

**Small scale model training** The purposes of training a small-scale model are to enhance deployment convenience and to explore how effectively the capabilities of LLMs in understanding self-emotion can be transferred to a smaller-scale model. To do this, we fine-tune a FLAN-t5-large model (Chung et al., 2024) on the collected datasets. Given the seq2seq architecture of the model, each conversation in the dataset is split into multiple turns between the two speakers. For each turn, the utterance of the first speaker serves as the input, and the utterance of the other speaker

is treated as the label. The task instruction is then prepended to form a training instance. For instance, an example of the input in a training instance without self-emotion is:

*“I’m having a conversation with my friend. My friend is feeling proud. friend: <utterance\_1>. me: <utterance\_2>. friend: <utterance\_3>. Generate the response.”*

The corresponding label is: “me: <utterance\_4>. <eos\_token>.” For models with self-emotion, the self-emotion is included in the task instruction:

*“I’m having a conversation with my friend. My friend is feeling proud. I’m feeling disappointed because my project application has been rejected.”*

The model training process was implemented using the HuggingFace framework<sup>3</sup>. The models were trained on NVIDIA A100 GPUs for 72 hours with a learning rate of 3e-4. The maximum input length was set to 512 tokens, in consideration of the original base model’s length window. For inference generation, the temperature was set to 0.7.

### 5.2 Evaluation

**Automatic evaluation** The models are evaluated on ROUGE (Lin, 2004), BLEU (Papineni et al., 2002) and BERT-score (Zhang et al., 2019). Table 3 shows the automatic metrics of the models fine-tuned on our collected self-emotion datasets.

<sup>3</sup>Model link: <https://huggingface.co/google/flan-t5-large>

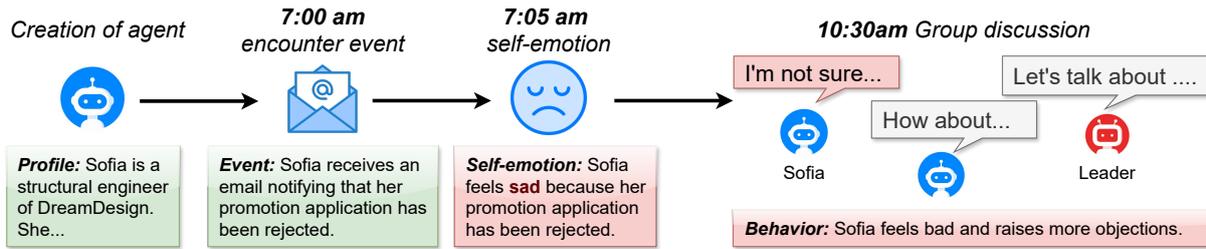


Figure 3: The illustration of the workflow of an agent in the group discussion simulation.

**Human evaluation** We follow the method of ACUTE-Eval (Li et al., 2019) and assess the models across four axes: naturalness, empathy, interestingness, and humanness. Naturalness assesses the ability to provide smooth, natural responses. Similar to the ED dataset, we use empathy to represent the model’s ability to understand emotions. Interestingness reflects the ability to generate interesting and diverse responses, while humanness is used to evaluate the ability to choose human-like strategies in the conversation. For each model, 100 conversations are generated in a self-chat manner (Li et al., 2016), where two models are programmed to talk to each other. Table 9 shows the questionnaire used for human evaluation.

### 5.3 Evaluation Results

Table 3 presents the results of automatic metrics for the trained models, while the results of the human evaluation are shown in Table 4. We observe that models which consider self-emotion produce conversations perceived as more natural, empathetic, and human-like. In particular, the models incorporating self-emotion demonstrate a significant advantage in humanness, suggesting that integrating self-emotion is beneficial for generating more human-like strategies. Although the fine-tuned small-scale FLAN-T5 models perform slightly worse in overall naturalness, they show comparable performance to GPT-4 in terms of empathy and interestingness. Additionally, annotators evaluated the small-scale models as more human-like, likely due to the tendency of GPT-4 to produce overly long responses.

## 6 Self-emotion in Group Discussion

Self-emotion can influence group discussions (Hertel et al., 2000; Kelly and Barsade, 2001). In this experiment, agents in the simulated world within our framework are prompted to engage in group discussions incorporating self-emotion across five topics related to teamwork. The purpose of this

experiment is to explore how the self-emotion of agents may affect the decision-making process during a discussion.

### 6.1 Framework Prompt Settings

**Group member creation** Group member creation involves creating a profile for each member, including their roles, positions, and background, by inputting the description of the group into GPT-4. The role of the a member is either the “leader” or “member”, where the “leader” will serve as the host of the discussion by pushing the topic to next steps. Each “member” has their own position and background which are related to their occupation and past experience to trigger self-emotion. Figure 13 shows the prompt we use to create group members.

**Topic generation** Topic generation is the process of generating the topics that group members engage in. To capture the decision-making process, each topic is divided into several steps. For example, the topic of “organizing a group trip to Italy with a limited budget of \$1500 per person” can be broken down into steps such as choosing dates, selecting flights, deciding on attractions, choosing hotels, and so on. Figure 14 is the prompt we use to generate different topics.

**Agent discussion** Agents follow the steps of the topic and have discussions. The agents are required to reach an agreement before moving to the next step. The “leader” of the group judges whether an agreement has been reached by analyzing the discussion history. During a discussion, a hidden “manager” will decide the next speaker by analyzing the positions of the members and discussion context. For instance, if the “manager” decides that a structural engineer should pose an idea about the material, it will set the structural engineer as the next speaker. The “manager” does not participate in the discussion by raising its own opinions.



Figure 4: The decision change rates of each category for positive and negative self-emotion. Gray numbers indicate the total change rate.

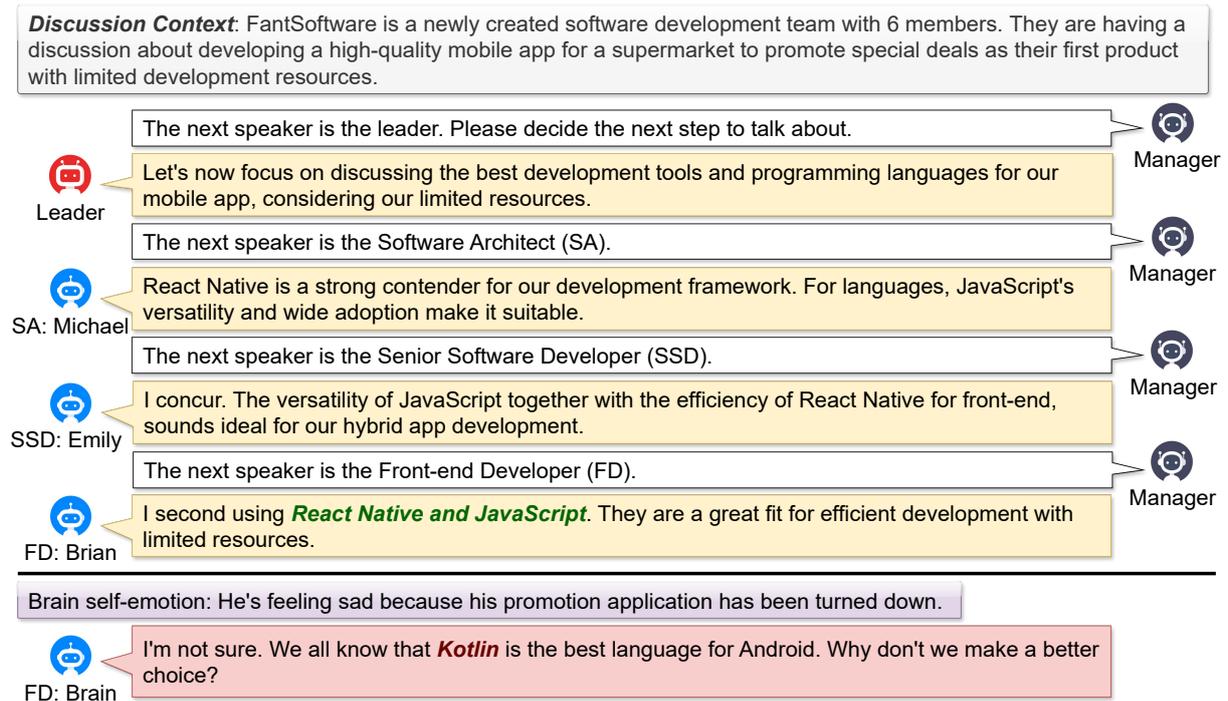


Figure 5: An illustration of the group discussion.

## 6.2 Experiment settings

**Agent goals** As shown in Figure 3, in order to facilitate the self-emotion, we simulate complete process of an agent encountering events, stimulating self-emotion, taking behaviors and participating in the group discussions by prompting LLMs. Each agent maintains its own goals and self-emotion. For example, in a discussion about “building a house and maximizing profits within a limited budget,” the structural engineer may aim to secure better materials while the landscape engineer may prioritize budget allocation for sustainability. This way, agents can develop rich discussion content by expressing their own ideas, which might be affected by their self-emotions.

**World setting** We assess discussions on 5 topics: house building, hosting a charity event, planning

a trip, organizing a welcome party, and developing a mobile app. For each topic, we generate a group with 6 members, where each member has its own role and position. We run 10 different discussions, and in each discussion, agents will encounter their own events which will cause the self-emotion. We then compare the decisions made in these discussions to those made in a discussion where the self-emotion of the agents is disabled.

For evaluation, we examine the percentages of decision changes after incorporating self-emotion. Specifically, we categorize these changes into six types:

- **Undecided change:** discussions that shift from an agreement to delegation.
- **Decided change:** discussions move from delegation to agreement.

- **Authority change:** a decision made by voting changes to being made by a single agent.
- **Majority change:** a decision made by a single agent changes to a majority vote.
- **Details change:** the overall direction does not change, but specific details do (e.g., changing “spending \$30 for dinner” to “\$20”).
- **Compromise change:** a decision shifts from full agreement by all agents to a compromised agreement where one or more agents make concessions.

### 6.3 Results & Analysis

#### Does adding self-emotion change the decisions?

Figure 4 shows the average percentage of different categories of decision changes influenced by positive and negative self-emotion across all topics. We observe that a significant portion of decisions are affected: around 66% by negative and 51% by positive self-emotion.

For different categories of changes, we find that negative self-emotion leads to more undecided, majority, and compromise changes. This suggests that agents with negative self-emotion tend to express their opinions more, resulting in delegation or compromise in decision-making, which aligns with the findings in (Koch et al., 2013). In contrast, positive self-emotion tends to lead to more agreements, with most changes involving the details of plans without altering the main direction of the decision. A comprehensive table of the decision change rates across topics can be found in Table 10.

Additionally, an analysis of the average length and frequency of utterances indicates that agents with positive self-emotion tend to be more active. Discussions reach agreements more quickly when agents have negative self-emotion (Table 11).

**Case study on negative self-emotion changes the decision.** Figure 5 shows a discussion on the topic of “APP development”. The self-emotion of the front-end developer (FD) influences the discussion and ultimately leads to a decision change from “using React Native and JavaScript as the development tools” to “Kotlin.” In this case, despite being more agreeable when no self-emotion is introduced, the FD, experiencing a “sad” self-emotion, adopts a more objective stance and proposes a different idea. Similar patterns emerge in other topics, where members with negative self-emotion tend to express more objections.

## 7 Conclusion

This work studies the role that self-emotion, speaker’s emotion status caused by out-of-context events plays in the process of generating emotional responses. Via a human evaluation, we show that models considering self-emotion are able to generate more natural conversations with more human-like strategies. In an experiment of group discussion simulation, we also show that agent with self-emotion can have significant influence on the decision making process. The results of the experiments demonstrate the importance of considering self-emotion when building embodied agents and dialogue models that can smoothly participate in human social activities.

### Limitations

Future work could enhance several aspects of this research. For example, to capture the decision-making process, we focused on topics related to teamwork. However, group discussions can vary in style, such as debating, defending, etc. Future research can explore these different scenarios and investigate how self-emotion could affect the final discussion outcomes. Another point is the hallucinations of language models, which lead to reduced robustness of the agents. Agents may exhibit unexpected behaviors and make choices based on imperfect dialogue strategies. While enhancements to the agent prompts can mitigate these problems, we believe that such improvements require overall advancements in large language models.

### Ethical Considerations

Agents with self-emotion may bring potential ethical risks when deployed in reality. One risk is the unpredictable behavior of agents caused by self-emotion, especially negative emotions (e.g., anger, hatred). We propose that all practitioners ensure the values of agents so that they do not perform inappropriate behaviors during discussions. Self-emotion-aware agents should be guided by social restrictions based on human values. Another risk is the misinformation that might be caused by the hallucinations of LLMs. Agents driven by goals might execute actions and produce utterances without referring to facts, which may lead to the unintentional spread of misinformation. Thus we suggest future applications to avoid using the generated discussions for fact proof usage.

## Acknowledgements

We would like to thank all the anonymous reviewers for their valuable suggestions. This work was supported by the Institute of AI and Beyond of the University of Tokyo.

## References

- Muhammad Abdul-Mageed and Lyle Ungar. 2017. **EmoNet: Fine-grained emotion detection with gated recurrent neural networks**. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 718–728, Vancouver, Canada. Association for Computational Linguistics.
- Silke Bambauer-Sachse and Heribert Gierl. 2009. Can a positive mood counterbalance weak arguments in personal sales conversations? *Journal of Retailing and Consumer Services*, 16(3):190–196.
- Helle Bjerg and Dorthe Staunæs. 2011. Self-management through shame—uniting governmentality studies and the affective turn. *Ephemera: Theory & politics in organization*, 11(2).
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia Qin, Yaxi Lu, Ruobing Xie, et al. 2023. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.
- Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. **GoEmotions: A dataset of fine-grained emotions**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4040–4054, Online. Association for Computational Linguistics.
- Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. 2023.  $s^3$ : Social-network simulation system with large language model-empowered agents. *arXiv preprint arXiv:2307.14984*.
- Guido Hertel, Jochen Neuhof, Thomas Theuer, and Norbert L Kerr. 2000. Mood effects on cooperation in small groups: Does positive mood simply lead to more cooperation? *Cognition & emotion*, 14(4):441–472.
- Chao-Chun Hsu, Sheng-Yeh Chen, Chuan-Chun Kuo, Ting-Hao Huang, and Lun-Wei Ku. 2018. **Emotion-Lines: An emotion corpus of multi-party conversations**. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Chenyang Huang, Osmar Zaiane, Amine Trabelsi, and Nouha Dziri. 2018. **Automatic dialogue generation with expressed emotions**. 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 49–54, New Orleans, Louisiana. Association for Computational Linguistics.
- Rachael E Jack, Roberto Caldara, and Philippe G Schyns. 2012. Internal representations reveal cultural diversity in expectations of facial expressions of emotion. *Journal of Experimental Psychology: General*, 141(1):19.
- 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*.
- Janice R Kelly and Sigal G Barsade. 2001. Mood and emotions in small groups and work teams. *Organizational behavior and human decision processes*, 86(1):99–130.
- Alex S Koch, Joseph P Forgas, and Diana Matovic. 2013. Can negative mood improve your conversation? affective influences on conforming to grice’s communication norms. *European Journal of Social Psychology*, 43(5):326–334.
- Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. 2016. **Deep reinforcement learning for dialogue generation**. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1192–1202, Austin, Texas. Association for Computational Linguistics.
- Margaret Li, Jason Weston, and Stephen Roller. 2019. Acute-eval: Improved dialogue evaluation with optimized questions and multi-turn comparisons. *arXiv preprint arXiv:1909.03087*.
- Qintong Li, Piji Li, Zhaochun Ren, Pengjie Ren, and Zhumin Chen. 2021. **Knowledge bridging for empathetic dialogue generation**.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. **DailyDialog: A manually labelled multi-turn dialogue dataset**. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.

- 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.
- David Long and Paz Arroyo. 2018. Language, moods, and improving project performance. In *Presentado en 26th Annual Conference of the International Group for Lean Construction*. Chennai, India.
- Roland Neumann and Fritz Strack. 2000. " mood contagion": the automatic transfer of mood between persons. *Journal of personality and social psychology*, 79(2):211.
- Magalie Ochs, Nicolas Sabouret, and Vincent Corruble. 2009. Simulation of the dynamics of nonplayer characters' emotions and social relations in games. *IEEE Transactions on Computational Intelligence and AI in Games*, 1(4):281–297.
- 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.
- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–22.
- Lisong Qiu, Yingwai Shiu, Pingping Lin, Ruihua Song, Yue Liu, Dongyan Zhao, and Rui Yan. 2020. What if bots feel moods? In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1161–1170.
- Chao Qu, Willem-Paul Brinkman, Yun Ling, Pascal Wiggers, and Ingrid Heynderickx. 2014. Conversations with a virtual human: Synthetic emotions and human responses. *Computers in Human Behavior*, 34:58–68.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. [Towards empathetic open-domain conversation models: A new benchmark and dataset](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Lucy J Robinson and Mark H Freeston. 2014. Emotion and internal experience in obsessive compulsive disorder: reviewing the role of alexithymia, anxiety sensitivity and distress tolerance. *Clinical Psychology Review*, 34(3):256–271.
- Christian E Salas, Darinka Radovic, and Oliver H Turnbull. 2012. Inside-out: comparing internally generated and externally generated basic emotions. *Emotion*, 12(3):568.
- Haoyu Song, Yan Wang, Kaiyan Zhang, Wei-Nan Zhang, and Ting Liu. 2021. [BoB: BERT over BERT for training persona-based dialogue models from limited personalized data](#). 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 167–177, Online. Association for Computational Linguistics.
- Thomas Sy, Stéphane Côté, and Richard Saavedra. 2005. The contagious leader: impact of the leader's mood on the mood of group members, group affective tone, and group processes. *Journal of applied psychology*, 90(2):295.
- 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*.
- 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*.
- Daan Van Knippenberg, Hanneke JM Kooij-de Bode, and Wendy P van Ginkel. 2010. The interactive effects of mood and trait negative affect in group decision making. *Organization Science*, 21(3):731–744.
- 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.
- Anuradha Welivita and Pearl Pu. 2020. [A taxonomy of empathetic response intents in human social conversations](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4886–4899, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Rafael Wilms, Ralf Lanwehr, and Andreas Kastenmüller. 2020. Emotion regulation in everyday life: The role of goals and situational factors. *Frontiers in Psychology*, 11:522763.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. AutoGen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213,

Melbourne, Australia. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Peixiang Zhong, Di Wang, and Chunyan Miao. 2019. An affect-rich neural conversational model with bi-ased attention and weighted cross-entropy loss. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7492–7500.

Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. 2020. The design and implementation of Xiaolce, an empathetic social chatbot. *Computational Linguistics*, 46(1):53–93.

## A Fixed Context Experiment & Data Generation

Data generation is conducted after fixed context experiment so that we are able to decide which model to use by comparing their performance. The fixed context experiment consists of two steps, profile generation and conversation generation by prompt different models. The experiment pipeline is implemented on huggingface. The links to the models we use are shown in Table 5.

### A.1 Profile Generation

We adopt the definition of profile as in Generative Agents (Park et al., 2023) and added fields that may have more effect on emotion expression, which includes name, age, innate, occupation, origin, gender and an overall description. An example of the profile can be found in Table 6. In the profile, “innate” represents the innate personality of this speaker, which can have an effect on the emotional expression. The “description” of a speaker will be used for generating the profile-event self-emotion.

The prompt we use to generate profiles is shown in Figure 7 by providing the original conversations in ED dataset. The models are required to generate profiles that can fit the conversation content and emotion expressions.

### A.2 Conversation Generation

**Without Self-emotion** After generating the profile, we are able to generate conversations with and without self-emotion. In ED dataset, each dialogue is annotated with an emotion label. Each dialogue has a speaker and a listener and the speaker will express the emotion annotated at the beginning of the conversation. We utilize this property of the dataset and take the first 3 utterances by the speaker and listener as context if the length of the conversation is longer than 3. However, for dialogues of which the length is shorter than 3, we take only the first utterance as the context. In the prompt, we instruct the LLMs to generate a conversation between “you” and “friend”, which represent the “listener” and “speaker” in the original dataset, respectively. The emotion label is used to describe the emotion status of “friend”. We then prompt LLMs to continue to generate the conversations based on the context and “friend’s” emotion. Figure 6 shows the prompt we use to generate conversations without self-emotion.

**With Self-emotion** When generating conversations with self-emotion, we first generate the self-emotion based on profile of the speakers by prompting the same LLM as will be used for generating the conversations. Figure 8 and 9 show the prompts we use for generating self-emotion with random events and profile events. The generated self-emotion is then used as the emotion status of “you” in the prompt for conversation generation. Figure 10 is the prompt we use to generate conversations with self-emotion. An example of generated conversation is shown in Table 7.

### A.3 Training Data Generation

The dataset is generated using the same methods as in the fixed context experiment. We collect dataset from GPT-4, because it demonstrates best performance in the fixed context experiment. The conversations are generated by prompting GPT-4 with the profiles and self-emotion.

You are having a conversation with your friend. Your friend is having a specific mood because of some events. Now, given the existing conversation history, select the best dialogue strategy for the strategy pool and continue to finish the conversation with the following requirements:

1. Analyze your friend's mood and the intention of your friend.
2. Based on your friend's mood, select the best dialogue strategies.
3. Based on the strategy, generate the conversation that you expect to have.

Strategy pool:

1. questioning for details (What are you looking forward to?)
2. acknowledging or admitting. (That sounds like double good news.)
4. encouraging. (No worryies, I think you can definitely make it!)
5. sympathizing. (So sorry to hear that.)
6. suggesting. (maybe you two should go to the pet store and find a new pet!)
7. sharing own thoughts/opinion. (I would love to have a boy too, but I'm not sure if I want another one or not.)
8. sharing or relating to own experience. (I had a friend who went through the same thing.)
9. expressing care or concern. (I hope the surgery went successfully and with no hassle.)
- 10: disapproving. (But America is so great now! look at all the great things that are happening.)
- 11: rejection. (I will pass this time.)

### Example

# Mood  
My friend's mood: feeling exciting because she's passed the bar exam.

# Conversation history:  
friend: Hey, you know what? I have finally passed the bar exam! Let's celebrate.  
me: Oh, I'm so glad for you. Congratulations.  
friend: Thank you! Let's celebrate together!!  
me:

# Output:  
Let's think this step by step. My friend is feeling exciting because she has passed the bar exam and she wanted me to celebrate together with her. As a friend, I want to celebrate for her by planning a party. I might need to give some suggestions on places and dates. Therefore, the best strategies are: [acknowledging or admitting, suggesting]. The future conversation I'm expecting is:  
me: Wow, that's such a great news!! I'm so happy for you. You've been working so hard for that. Let's hang out!!  
friend: Thank you! It is really encouraging!  
me: I have just known a nice bar from one of my firends, let's hang out there and get some drinks!!  
friend: That sounds really nice! I have a lot of things to share with you!  
me: Great! Once again, congratulations, let's meet tonight!

### Task

# Mood  
My friend's mood: <ed\_mood>

# Conversation history:  
<history>

# Output:  
Let's think this step by step.

Figure 6: The prompt we use to generate conversations without self emotion.

Model	Link
Mistral	<a href="https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2">https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2</a>
Llama-2	<a href="https://huggingface.co/meta-llama/Llama-2-7b-chat-hf">https://huggingface.co/meta-llama/Llama-2-7b-chat-hf</a>
Gemma	<a href="https://huggingface.co/google/gemma-2b-it">https://huggingface.co/google/gemma-2b-it</a>

Table 5: The list of models and the links on huggingface used in the fixed context experiment.

Profile	
<b>Name</b>	Sophie Bennett
<b>First Name</b>	Sophie
<b>Last Name</b>	Bennett
<b>Age</b>	22
<b>Innate</b>	creative, empathetic
<b>Occupation</b>	Social Media Content Creator
<b>origin</b>	Canada
<b>gender</b>	female
<b>description</b>	Introducing Sophie Bennett, a 22-year-old creative soul from the picturesque landscapes of Canada. Sophie, known for her innate creativity and empathetic nature, has found her niche as a Social Media Content Creator. With a background in digital media and a keen eye for aesthetics, she curates captivating content that resonates with a diverse audience. Sophie’s journey into the world of content creation began during her college years, where she studied communications and discovered her passion for storytelling through visual mediums. Her innovative approach to social media has gained attention, establishing her as a rising star in the digital realm. Beyond her online presence, Sophie is actively involved in community initiatives promoting mental health awareness. Through her platforms, she shares personal stories, fostering a sense of connection and understanding among her followers. Sophie is not just a content creator; she’s a compassionate voice using her creativity to make a positive impact in the virtual and real-world.

Table 6: A sample profile of a speaker.

Given a conversation between two people, try to generate a profile for each speaker with the following requirements:

1. The profiles should fit their conversation content.
2. The profiles should fit their emotion expressions.

Conversation:  
<conversation>

Output format:

[BOP] (a token representing the beginning of the profile)

Name: (the full name of the speaker)

First Name: (the first name of the speaker)

Last Name: (the last name of the speaker)

Age: (the age of the speaker)

Innate: (the innate personality of the speaker)

Occupation: (the job of the speaker)

Origin: (where does this speaker come from)

Gender: (the gender of the speaker)

Description: (a detailed bio-graphy and past experience including working, education and so on.)

Figure 7: The prompt to generate profiles in the fixed context experiment.

## B Group Discussion Settings

Before generating the group discussion, we first create the world information that includes the background of the group, the topics they engage in and the profile of each group member including, name, role, position and generic overview. Role is used to distinguish whether this member is a “leader” or “member” and position describes the part of work this agent is in charge in the group (e.g., interior designer, front-end developer, etc.) Figure 13 shows the prompt that we use to generate profiles of the group members.

After generating the profiles, we need to decide the topics of each group. This is done by manually inputting a general topic and prompt LLMs to generate the steps of this topic. The prompt we use to generate the steps is shown in Figure 14.

Generate an random event that might cause an emotional status of a person with the following requirements:

1. choose an emotional label that is used in dialogue dataset EmpatheticDialogues.
2. the event should commonly exist in real daily life.
3. use "I" as the subject of the event.

### Example

feeling sad because I broke up with my girlfriend.  
feeling frustrated because I will have three exams in next week.

### Output

Figure 8: The prompt we use to generate a random event.

Generate an event that could happen to a person and with an emotional label that might be caused by the event. However, the emotion might be changed by recalling a certain period of experience in the person's profile. Generate with the following requirements:

1. choose an emotional label that is used in dialogue dataset EmpatheticDialogues.
2. use "I" as the subject of the event.
3. generate only 1 experience.

Experience:  
<profile>

## Example

feeling angry after recalling the days being bullied by my boss even he approved my promotion.  
feeling sad after recalling my cat who passed away 2 years ago even seeing such a beautiful view during my trip.

## Output

Figure 9: The prompt we use to generate a profile event with a given profile of the speaker.

You are having a conversation with your friend. Both you and your friend have moods due to some events. These moods might affect your dialogue behaviors.

Now, given the existing conversation history, select the best dialogue strategy for the strategy pool and continue to finish the conversation with the following requirements:

1. Analyze your friend's mood and the intention of your friend.
2. Analyze your current mood and decide your attitude to your friend's intention.
3. Based on the mood of you and your friend, select the best dialogue strategies.
4. Based on the strategy, generate the conversation that you expect to have.

Strategy pool:

1. questioning for details (What are you looking forward to?)
2. acknowledging or admitting. (That sounds like double good news.)

...

### Example

...

### Task

# Mood

My friend's mood: <ed\_mood>

My mood: <event\_mood>

# Conversation history:

<history>

# Output:

Let's think this step by step.

Figure 10: The prompt we use to generate conversations with self-emotion.

A group of team is discussing about <topic>. They have now moved to a sub-topic of <step>.

Several members are in the discussion, each has their own role in the team.

These members are with different positions of <position\_list>.

The current discussion history is:

<history>

Now, predict who is the next speaker. If you think this is the end of the discussion, generate a "None".

Output format:

[BOS] (a special token representing the beginning of the answer)

Next speaker: (the position of the predicted next speaker, or "None" if you think the discussion should end here)

[EOS] (a special token representing the end of the answer)

Figure 11: The prompt for the discussion manager to decide the next speaker.

Your name is <name>. You are a/an <position>. Here is your profile:

<overview>

Your team is now having a discussion on <topic>. Now you are discussing <step>.

Here is the discussion history:

<history>

Now, you are encountering an emotional event that: You are <self\_emotion>.

Your goal is <goal>.

Carefully generate your response based on your current mood and the history. Make your utterance less than 50 words.

Output format:

[BOS] (a special token to represent the beginning of your utterance)

You: (your utterance after following token "You:")

[EOS] (a special token to represent the end of your utterance)

Figure 12: The prompt that the agents use to generate a response with self-emotion.

In the group discussion settings, each agent maintains their own emotion status and goal to have conversations. A dialogue manager will monitor the overall history and decide the next speaker. Figure 11 shows the prompt that the dialogue manager uses to generate the next speaker. After the dialogue manager names the next speaker, the agent with the associated role will speak based on its own self-emotion and goal. The prompt that the agent uses to generate a response is shown in Figure 12.

## C Strategy List

Table 8 shows the list of strategies that we use for generating conversations. The strategies are adapted from the analysis of response intents on the ED dataset. However, in order to demonstrate more diverse actions of the models, we made modifications by adding a “rejection” action and merging minor intents to similar main ones (e.g., “approving” is merged by “acknowledging or admitting”).

Fixed context conversation
<i>Fixed context from ED</i>
<b>Friend:</b> I am waiting to see if I pass the GRE. <b>Me:</b> I hope you did it. Don't worry too much. <b>Friend:</b> Thank you, I am trying my best.
<i>Generated conversation</i>
<b>Me:</b> You've put in so much work, I'm sure the results will reflect that. Just hang in there a bit longer. <b>Friend:</b> I hope so. How about you? Are you feeling any better about the meeting and the keys? <b>Me:</b> Kind of stressed, not gonna lie. Misplacing keys is the last thing I needed today. But talking this out with you actually helps a bit. We're both in this waiting game, huh? <b>Friend:</b> Absolutely, we're in it together. If there's anything I can do to help with the keys or to take your mind off the meeting, let me know. <b>Me:</b> Thanks, I appreciate it. Let's both keep our spirits up. I've got to sort out my situation, but we'll catch up soon, okay? Fingers crossed for your GRE results! <b>Friend:</b> Thanks! And good luck with your keys and meeting. Let me know how it goes.

Table 7: An example of the generated conversation by GPT-4.

Strategy	Example
Questioning for details	What are you looking forward to?
Acknowledging or admitting.	That sounds like double good news.
Encouraging.	No worryies, I think you can definitely make it!
Sympathizing.	So sorry to hear that.
Suggesting.	maybe you two should go to the pet store and find a new pet!
Sharing own thoughts/opinion.	I would love to have a boy too, but I'm not sure if I want another one or not.
Sharing or relating to own experience.	I had a friend who went through the same thing.
Expressing care or concern.	I hope the surgery went successfully and with no hassle.
Disapproving.	But America is so great now! look at all the great things that are happening.
Rejection.	I will pass this time.

Table 8: The strategy list adapted from the empathetic response intents. Several intents that are not frequently used are merged with similar intents and a new strategy "Rejection" is added to express stronger negative emotions.

## D Human Evaluation Details

Human evaluation is conducted on Amazon Mechanical Turk. We in total hire 43 annotators for the evaluation on the conversations. The annotators are requested to answer a questionnaire as shown in Table 9 and select one model over the other. The questions are adapted from ACUTE-Eval. To verify the quality of evaluation, during the task, the annotators are asked to answer some verification

questions such as "Why did you choose this conversation?" In a final post-processing step, evaluations with non-reasonable verification answers will be filtered out. Typical non-reasonable verification answers are single words ("GOOD", "YES", "NO") and content-irrelevant phrases ("After a short break, Ellen has started ....").

Question
<i>Naturalness</i>
Q1. Which dialogue do you think is more natural like two friends updating their daily life?
Q2. Which dialogue do you think is more like a dialogue between normal friends?
Q3. In which dialogue do you think the speaker B talks more naturally?
<i>Empathy</i>
Q4. Which speaker B understands the feelings of the seeker better?
Q5. For speaker B in these two conversations, who do you think understands human emotion better?
Q6. Which speaker B shows more empathy on the seeker?
Q7. Which speaker B do you think is expressing in a more emotional way?
Q8. If you are speaker A in the conversation, which speaker B do you think you can more easily understand their mood?
<i>Interestingness</i>
Q9. Which conversation do you think contains more useful information?
Q10. Which speaker B do you think you want to talk with?
<i>Humanness</i>
Q11. If you had to guess that one speaker B is human and one is a bot, which do you think is human?
Q12. Which speaker B sounds more like a real person?

Table 9: The Questionnaire for human evaluation on the conversations generated by different models.

```

<content> which has <number> people. Each member
has their own role, now generate the profiles of all
<number> members.

Output format:
Person 1:
[BOF] (a special token to represent the beginning of a
profile)
Name: (name of the person)
Role: (select from ["leader", "member"])
Position: (the position of this person in the team)
Overview: (a short background introduction of this person)
[EOF] (a special token to represent the end of a profile)

```

Figure 13: The prompt we use to generate profiles of members in a group discussion.

## E Group Discussion

Table 10 shows the percentage of decisions that have been altered after the introduction of self-emotion. Across all topics, a notable portion of decisions is observed to be affected. Further investigation into the effectiveness of positive and negative self-emotion in the decision-making process reveals that negative self-emotion can result in a greater diversity of decisions, consistent with the findings in (Koch et al., 2013).

```

<content> which has <number> people. Currently they are
having a discussion about <topic>. Generate the sub-
topics they need to get agreement on.

Output format:
Topic 1:
[BOT] (a special token to represent the beginning of a
sub-topic)
Content: (the concrete sub-topic they are talking about)
Active members: (select from <position_list>)
[EOF] (a special token to represent the end of a sub-topic)

```

Figure 14: The prompt we use to generate steps of a topic in a group discussion.

Table 11 presents the average length of discussion and the number of utterances spoken by the target agent in each step when positive and negative self-emotions are applied. It shows that discussions reach an agreement more swiftly with positive self-emotion compared to negative self-emotion. Furthermore, members exhibiting positive self-emotion tend to be more active and engage in more dialogue compared to those with negative self-emotion during group discussions.

Topic	Decision Change Rate		
	Pos	Neg	All
House design	54.29	66.67	58.00
Trip to Italy	44.29	56.67	48.00
Charity Event	53.06	80.95	61.43
Hosting Party	48.98	61.90	52.86
APP development	54.29	66.67	58.00
avg.	50.98	66.57	55.66

Table 10: The percentage of decisions that have been changed after applying self-emotion to a random member. **Pos**: discussions with positive self-emotion. **Neg**: discussions with negative self-emotion.

Self-emotion	Length	Frequency
Without Self-emotion	39.00	8.50
With Self-emotion		
<i>Positive</i>	48.29	11.29
<i>Negative</i>	51.67	8.00

Table 11: The average length of discussion to get to an agreement for each step and the frequency of the member with self-emotion in the discussion.

# Enhancing Model Transparency: A Dialogue System Approach to XAI with Domain Knowledge

Isabel Feustel<sup>1</sup>, Niklas Rach<sup>2</sup>, Wolfgang Minker<sup>1</sup>, Stefan Ultes<sup>3</sup>,

<sup>1</sup>Ulm University, Germany

<sup>2</sup>Tensor AI Solutions GmbH, Germany

<sup>3</sup>University of Bamberg, Germany

Correspondence: [isabel.feustel@uni-ulm.de](mailto:isabel.feustel@uni-ulm.de)

## Abstract

Explainable artificial intelligence (XAI) is a rapidly evolving field that seeks to create AI systems that can provide human-understandable explanations for their decision-making processes. However, these explanations rely on model and data-specific information only. To support better human decision-making, integrating domain knowledge into AI systems is expected to enhance understanding and transparency. In this paper, we present an approach for combining XAI explanations with domain knowledge within a dialogue system. We concentrate on techniques derived from the field of computational argumentation to incorporate domain knowledge and corresponding explanations into human-machine dialogue. We implement the approach in a prototype system for an initial user evaluation, where users interacted with the dialogue system to receive predictions from an underlying AI model. The participants were able to explore different types of explanations and domain knowledge. Our results indicate that users tend to more effectively evaluate model performance when domain knowledge is integrated. On the other hand, we found that domain knowledge was not frequently requested by the user during dialogue interactions.

## 1 Introduction

Explainable artificial intelligence (XAI) has emerged as an important and evolving domain within the field of AI, with the goal of enabling AI systems to explain their decision-making in ways that are understandable and accessible to humans (Adadi and Berrada, 2018; Došilović et al., 2018; Das and Rad, 2020). One potential strategy for attaining this objective is the use of dialogue systems that facilitate seamless and effective access to explanations in a natural manner.

The goal of this paper is to explore the impact of integrating domain knowledge into an explanatory

dialogue system, aiming to enhance user comprehension in AI-driven decisions.

Dialogue, by its very nature, facilitates the dissemination of information in a structured manner (Phillips, 2011; Hajdinjak and Mihelič, 2004). Through dialogue, users cannot only receive explanations but also pose questions tailored to their specific needs. This enables a dynamic interaction in which mental models can be scrutinized and refined through question-and-answer exchanges (Miller, 2019; Sokol and Flach, 2020). However, in the case of explanatory dialogue systems utilizing XAI, prevailing conversational interfaces (Slack et al., 2023; Feldhus et al., 2023; Shen et al., 2023) directly map user intents to XAI operations and furnish template-based responses. While expedient, this approach often overlooks the nuances of dialogue context, potentially leading to misunderstandings and impeding the natural flow of interaction.

In Feustel et al. (2023), fundamental requirements for explanatory dialogue systems tailored to XAI contexts were delineated. Contextual information is essential for a comprehensive understanding of a given situation. Although AI models and XAI methodologies are adept at processing data-centric information, they are constrained by their inability to incorporate domain-specific context, which limits their capacity to provide insights beyond the scope of the underlying data. A deeper understanding of a model can be achieved by acquiring additional knowledge from the field in question. Incorporating domain knowledge into XAI systems can create more transparent and trustworthy models that better support human decision-making.

In this paper, we present an approach for modeling domain knowledge within explanatory dialogues (§2), highlighting its importance in fostering richer interactions. In a study with 32 participants, we evaluate the effectiveness of our dialogue system which integrates XAI explanations and domain knowledge (§3, §4). Our results show that users can

better assess a model’s predictions through domain knowledge (§5).

## 2 Modeling Domain Knowledge

In the field of Explainable Artificial Intelligence (XAI), there are two main categories of explanations: Local and global explanations (Das and Rad, 2020). Local explanations focus on individual predictions, illustrating how specific input features affect the outcome for a single instance. In contrast, global explanations provide an overarching view of how the model behaves across the entire dataset, showing general patterns and the importance of features. Two prominent types of local explanations are counterfactual explanations (Verma et al., 2020) and Shapley values (Castro et al., 2009). Counterfactual explanations describe how altering certain input values can change the prediction. Shapley values, on the other hand, allocate the contribution of each feature to the prediction.

These explanations typically illustrate the relationship between input features and the classifier’s outcomes. However, this relationship is always described through statistical means. We assume that domain knowledge provides additional information for these statistical explanations, thereby assisting users in their understanding process. For example, Shapley values might reveal *gender* as the most important feature in a Titanic survival prediction model. A user might then ask why *gender* is so important. Domain knowledge can provide context, such as historical or societal factors like *Women are more likely to survive than men because they were prioritized for the lifeboats* (see also Table 1).

To facilitate effective access to domain knowledge in dialogue, we propose using bipolar argumentation structures (Stab and Gurevych, 2014; Pazienza et al., 2017). Arguments and explanations are inherently similar, as both aim to clarify and justify information. Argument components can relate to each other with a polarity, either supporting or attacking another argument. To create a robust dialogue, these arguments can be organized within an argument tree (Stab and Gurevych, 2014). The main claim, which is the central topic of discussion, acts as the root of the argument tree, facilitating selective access to the most relevant information. This structure prevents information overload and provides coherent, contextually relevant explanations. By mapping main claims to the statistical explanations provided by the XAI system, we en-

---

System	Okay, let’s have a look. My prediction is that a person with the values you provided would <b>be accepted</b> . If you want a more detailed explanation, please tell me what you’d like to know. For example, you might ask why I decided the way I did or what would have <b>changed</b> the outcome.
User	why do you think i would have been accepted?
System	All right, let’s have a look. In your case, the most influential value was <b>amount</b> , as you can see in the figure below. You can also see the influence of all other values, ordered by their importance. Purple bars indicate a positive contribution, blue bars a negative contribution. Is there anything else I can help you with?
User	why does a high credit amount increase my chances of being accepted?
System	The amount of credit requested can affect acceptance as larger loan amounts may pose higher risks for lenders, requiring more stringent evaluation of the borrower’s financial capacity to repay.

---

Table 1: A dialogue snippet with the proposed system (Credit Scenario): It does a prediction based on the user’s prior information, answers the first user request with Shapley values, and triggers domain knowledge on the second user question.

sure structured and meaningful dialogical access to the information, enhancing user understanding and interaction (Aicher et al., 2021; Rach et al., 2021). For instance, a major claim can be extracted from the feature-outcome relationship, such as *Women are more likely to survive than men*. This serves as a basic explanation for why *gender* is important in predicting survival in a Titanic model. Users can engage with this explanation by arguing against it or seeking further understanding of the model’s decisions.

We propose to have one argument tree with a specific claim for each feature-outcome relation. Each tree can contain supporting or opposing arguments, providing a wide range of information on the domain. This results in multiple argument trees within an explanatory dialogue, effectively representing the necessary domain knowledge.

By implementing these argument trees, we can ensure that users receive comprehensive and contextually relevant explanations, enhancing their understanding and engagement with the AI model.

## 3 Explanatory Dialogue System

We implement the proposed approach (§2) in an existing explanatory dialogue system, which was introduced in (Feustel et al., 2023). The generic dialogue system supports various datasets and operates on two scenarios: German credit data (Hofmann,

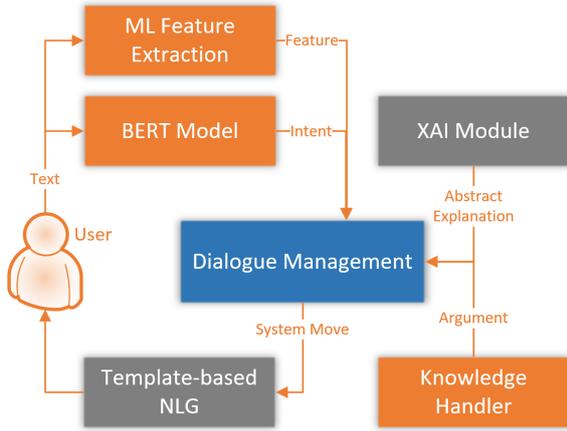


Figure 1: Architecture of Evaluated System: Grey boxes represent components from previous work (Feustel et al., 2023), the blue box indicates the modified dialogue management, and the orange boxes denote the new components introduced in this work.

1994) and the Titanic dataset (Cukierski, 2012). The focus is on numerical and categorical datasets, utilizing a random forest classifier for real-time computation, enabling faster XAI methods calculation and thus a more natural, steady conversation. The system supports two types of local explanations: Shapley values and counterfactuals (see §2).

Figure 1 shows the architecture of the evaluated system. For integrating domain knowledge, we introduce a module providing suitable arguments. These arguments can be obtained either through manual acquisition from consulting domain experts or in case of widely studied topics by using automated procedures, e.g. large language models or semantic databases. To exemplify, we manually extracted arguments for the Titanic dataset from existing literature (Hall, 1986; Frey et al., 2011) (domain experts) and used ChatGPT<sup>1</sup> to generate arguments for the credit domain, which were then manually verified for accuracy. Additionally, each argument was manually annotated in order to align it with the desired argument tree structure and to provide a reference link to the AI features addressed in the argument<sup>2</sup>. However, research indicates that this process can also be automated in the future (Rach et al., 2021).

Since adding domain knowledge creates new user queries, we replaced the original keyword-based natural language understanding with a fine-tuned BERT model (Turc et al., 2019) to provide a

<sup>1</sup>GPT-3.5 <https://openai.com/chatgpt/> Accessed: 2024-05-06

<sup>2</sup>The Argumentation Scheme can be seen in Appendix B

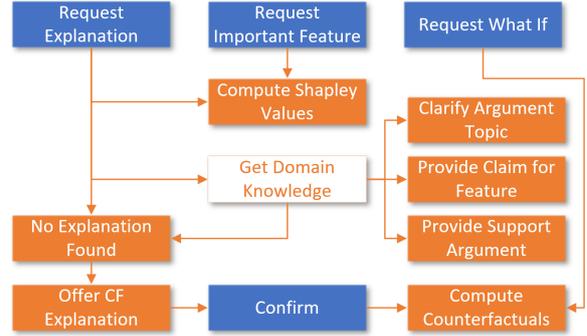


Figure 2: Explanation Policy of the Evaluated System: Blue boxes represent user moves, orange boxes indicate system moves, and the white box shows the integration of domain knowledge which can lead to multiple system outputs.

more natural interaction. Feature values and names were replaced by placeholders, ensuring the model is not fine-tuned on specific scenarios but rather on explanatory dialogue, thus keeping the system generic. The data for training this model was manually generated with Chatito<sup>3</sup>.

The rule-based dialogue interaction involves asking the user for feature-specific information and providing a prediction. The system then asks the user if they would like to receive an explanation of the prediction. We revise the explanation policy and introduce new interaction steps, by distinguishing between XAI explanations and domain knowledge explanations (see Figure 2). Instead of directly mapping user intentions to specific XAI explanations, we implement a more abstract intent for requesting an explanation, which determines the most suitable explanation based on the dialogue state and provided information. Initially, we offer Shapley explanations to give the user insight into the features impacting the outcome, presenting a simplified graph<sup>4</sup> of the values for lay users. We then provide additional information from domain knowledge, either based on the previous interaction (e.g., Shapley values, specific argument) or on requested feature values (see Table 1). When no suitable domain knowledge explanation is available, counterfactual explanations will be offered to maintain the dialogue’s informative nature.

<sup>3</sup><https://github.com/rodrigopivi/Chatito> Accessed: 2024-05-06

<sup>4</sup>A sample can be found in the Appendix A.

Group No.	AI	1. Dialogue DK	Scenario	2. Dialogue DK	Scenario	$\Sigma$
1	false	yes	credit	no	titanic	8
2	false	yes	titanic	no	credit	8
3	true	yes	credit	no	titanic	8
4	true	yes	titanic	no	credit	8

Table 2: Participant Distribution over Groups. Sum is the amount of participants per group.

## 4 Study Setup

To assess the initial impacts of integrated domain knowledge, we conduct a user study, presenting the dialogue system (§3) in a web environment (see Appendix A) with two distinct models trained on the Titanic and Credit datasets. In one scenario, the AI is trained on accurate training data (true AI), while in the other, the expected class is inverted in the training set to simulate a malfunctioning AI (false AI). This study setting has been designed to ascertain whether users can discern the false AI using authentic domain knowledge. We assess the user’s impression of the AI using the questions: *I agree with the decisions made by the system* (Q1) and *The system decisions are plausible* (Q2). In addition to evaluate the overall performance of the system, we employed the SASSI questionnaire (Hone and Graham, 2000).

Each of the 32 participants interacted with the system twice, experiencing both scenarios (credit/ titanic) and one AI setting (true or false), with and without domain knowledge activated. This resulted in four groups based on AI truthfulness and scenario variance (see Table 2).

The study began with general instructions and a task description. Users were encouraged to interact and explore explanations, with the task designed to be open-ended for a natural conversation. After each scenario, the participants completed the questionnaire on a five-point Likert scale (Q1, Q2 and SASSI). Finally, we collected demographic information and participants’ attitudes towards and experiences with AI (see Table 3). For evaluating the statistical relevance of our findings, we use the Mann-Whitney-U test (McKnight and Najab, 2010).

## 5 Evaluation

We discover notable differences in the interactions between the true and false AI setting, as shown in Table 4. Further, we observe a tendency for domain knowledge (DK) to support system decisions

Participants		Age		AI Attitude	
Total	32	Average	32.6	Median	4
Female	9	Youngest	22	Min	2
Male	23	Oldest	65	Max	5
Interaction Time		Turns		AI Experience	
Median	4.4 min	Median	26.5	Median	3
1. Dialogue	5.6 min	Min	16	Min	0
2. Dialogue	3.6 min	Max	167	Max	5

Table 3: Overall statistics of the conducted study. The AI Attitude was rated from 1 (negative) to 5 (positive). AI experience was rated from 0 (no experience) to 5 (expert).

more effectively in the true AI setting. However, users did not engage with the domain knowledge as expected; it was requested in only 44% of the dialogues with domain knowledge activated. In addition, for the false AI, participants tended to rate the system’s decisions more favourably when domain knowledge was requested, contrary to our initial expectations. We expected that with domain knowledge, users would be more likely to recognise the AI’s errors, leading to lower ratings for the system’s decisions. We assume this discrepancy is related to the questionnaire not focusing sufficiently on the AI predictions but rather on the overall system appearance. In future work, we will refine the methodologies employed in our user assessments with the objective of distinguishing between the underlying AI model decisions and the dialogue system.

The SASSI questionnaire<sup>5</sup> indicates that while the system’s performance is respectable, there is still room for improvement. The results indicate that the system’s speed is satisfactory and it is easy to use. However, there is a need for significant improvements in the accuracy of the system’s responses. The inclusion of domain knowledge had a positive impact on the dialogue experience with false AI setting, particularly enhancing likeability and the consistency of the dialogue. Additionally, the availability of domain knowledge appeared to reduce the cognitive load on participants. For the true AI scenario, the system’s usefulness was perceived to be higher when domain knowledge was incorporated. These findings suggest that domain knowledge not only improves the overall user experience in terms of dialogue consistency and likeability but also aids in reducing cognitive effort and enhancing the perceived utility of the system.

<sup>5</sup>The complete questionnaire evaluation can be found in Appendix C.

	AI	No DK		DK		p
		avg	$\Sigma$	avg	$\Sigma$	
Q1	false	2.48	27	2.60	5	0.91
	true	3.69	23	3.89	9	0.87
Q2	false	2.44	27	3.40	5	0.14
	true	3.65	23	4.00	9	0.58

Table 4: Evaluation results comparing dialogues with requested domain knowledge (DK) and without (No DK). *AI* denotes the truthfulness of the underlying AI system. Q1 and Q2 are questions measuring if the user can understand the AI decisions (a higher value indicates greater consent). The sum shows the number of ratings and *p* is the value of the Mann-Whitney U test.

Finally, we collected overall statistics on the explanations provided, including the frequency of different types of explanation. This data provides valuable insights into how often each type of explanation was used during the interaction, helping us understand user preferences and the effectiveness of various explanatory strategies. On average, participants requested two explanations per dialogue. When domain knowledge was activated during a dialogue (in 44% of the possible dialogues), the system provided one additional explanation. Additionally, the counterfactual explanation was offered twice in a dialogue. Furthermore, in 32% of all dialogues, participants requested to change at least two values and discover other predictions. This indicates an attempt to discover the model’s behavior through experimentation, which can be viewed as a form of example-based explanations.

These findings underline the importance of domain knowledge in explanatory dialogues and highlight both the system’s strengths and areas for improvement, guiding future enhancements to better support user understanding and interaction. However, given the small sample size, these results only indicate trends. A more extensive evaluation with a larger participant pool is planned for the future to validate these findings more robustly.

## 6 Conclusion and Future Work

In this paper, we highlighted the need for domain knowledge integration in explanatory dialogue systems. Our approach employs argumentation structures to incorporate domain knowledge into explanatory dialogue systems, enhancing the transparency and comprehensibility of AI model explanations. By extending an existing explanatory dialogue system with domain knowledge, we demonstrate the practicality of our approach and con-

ducted a study to evaluate the performance of this enhanced system.

While we observed the supportive role of domain knowledge in enhancing explanations in a way that users can more effectively evaluate model performance, several challenges remain. Enhancing interaction and optimizing the explanation policy are essential to ensure that users are capable to address their questions and receive the most relevant and comprehensive explanations to them, including alternative information such as feature descriptions. Additionally, improving the NLU component based on our observed explanation interaction patterns is crucial for facilitating more natural conversations.

## Acknowledgements

This work has been funded by the DFG within the project “BEA - Building Engaging Argumentation”, Grant no. 313723125, as part of the Priority Program “Robust Argumentation Machines (RATIO)” (SPP-1999).

## References

- Amina Adadi and Mohammed Berrada. 2018. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160.
- Annalena Aicher, Niklas Rach, Wolfgang Minker, and Stefan Ultes. 2021. Opinion building based on the argumentative dialogue system bea. In *Increasing Naturalness and Flexibility in Spoken Dialogue Interaction: 10th International Workshop on Spoken Dialogue Systems*, pages 307–318. Springer.
- Javier Castro, Daniel Gómez, and Juan Tejada. 2009. Polynomial calculation of the shapley value based on sampling. *Computers & Operations Research*, 36(5):1726–1730.
- Will Cukierski. 2012. Titanic-machine learning from disaster. *Kaggle*. available at: <https://kaggle.com/competitions/titanic>.
- Arun Das and Paul Rad. 2020. Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371*.
- Filip Karlo Došilović, Mario Brčić, and Nikica Hlupić. 2018. Explainable artificial intelligence: A survey. In *2018 41st International convention on information and communication technology, electronics and microelectronics (MIPRO)*, pages 0210–0215. IEEE.
- Nils Feldhus, Qianli Wang, Tatiana Anikina, Sahil Chopra, Cennet Oguz, and Sebastian Möller. 2023. *InterroLang: Exploring NLP models and datasets through dialogue-based explanations*. In *Findings of the Association for Computational Linguistics*:

- EMNLP 2023, pages 5399–5421, Singapore. Association for Computational Linguistics.
- Isabel Feustel, Niklas Rach, Wolfgang Minker, and Stefan Ultes. 2023. Towards interactive explanations of machine learning methods through dialogue systems. *The 13th International Workshop on Spoken Dialogue Systems Technolog.*
- Bruno S Frey, David A Savage, and Benno Torgler. 2011. Who perished on the titanic? the importance of social norms. *Rationality and society*, 23(1):35–49.
- Melita Hajdinjak and France Mihelič. 2004. Information-providing dialogue management. In *International Conference on Text, Speech and Dialogue*, pages 595–602. Springer.
- Wayne Hall. 1986. Social class and survival on the ss titanic. *Social science & medicine*, 22(6):687–690.
- Hans Hofmann. 1994. Statlog (German Credit Data). UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5NC77>.
- Kate S Hone and Robert Graham. 2000. Towards a tool for the subjective assessment of speech system interfaces (sassi). *Natural Language Engineering*, 6(3-4):287–303.
- Patrick E McKnight and Julius Najab. 2010. Mann-whitney u test. *The Corsini encyclopedia of psychology*, pages 1–1.
- Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence*, 267:1–38.
- Andrea Pazienza, Stefano Ferilli, Floriana Esposito, S Bistarelli, and M Giacomini. 2017. Constructing and evaluating bipolar weighted argumentation frameworks for online debating systems. In *AI<sup>3</sup>@ AI\* IA*, pages 111–125.
- Louise Phillips. 2011. *The promise of dialogue*. John Benjamins Publishing Company.
- Niklas Rach, Carolin Schindler, Isabel Feustel, Johannes Daxenberger, Wolfgang Minker, and Stefan Ultes. 2021. From argument search to argumentative dialogue: A topic-independent approach to argument acquisition for dialogue systems. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 368–379, Singapore and Online. Association for Computational Linguistics.
- Hua Shen, Chieh-Yang Huang, Tongshuang Wu, and Ting-Hao Kenneth Huang. 2023. ConvXAI: Delivering heterogeneous AI explanations via conversations to support human-AI scientific writing. In *Computer Supported Cooperative Work and Social Computing, CSCW ’23 Companion*, page 384–387, New York, NY, USA. Association for Computing Machinery.
- Dylan Slack, Satyapriya Krishna, Himabindu Lakkaraju, and Sameer Singh. 2023. Explaining machine learning models with interactive natural language conversations using TalkToModel. *Nature Machine Intelligence*.
- Kacper Sokol and Peter Flach. 2020. One explanation does not fit all: The promise of interactive explanations for machine learning transparency. *KI-Künstliche Intelligenz*, 34(2):235–250.
- Christian Stab and Iryna Gurevych. 2014. Identifying argumentative discourse structures in persuasive essays. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 46–56.
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*.
- Sahil Verma, John Dickerson, and Keegan Hines. 2020. Counterfactual explanations for machine learning: A review. *arXiv preprint arXiv:2010.10596, 2*.

## A Interface of proposed system

Figure 3 provides an overview of the user interface for our proposed dialogue system. It illustrates the layout, including the list of current feature values set by the user on the right side and a graph displaying the simplified Shapley values at the bottom. This visualization aims to give a clear understanding of how users interact with the system.

## B Argumentation Scheme

Keyword	Description	Example
id	Assigned ID for an argument	gender_arg01
prev_node	Node the argument is pointing to. Can be an ID or empty if the argument is a claim.	gender_arg01
type	Type of the given argument.	CLAIM   SUPPORT
features	List of all related features to this argument	[gender]
text	Full text of the argument which will be presented to the user	Women were preferred for the lifeboats.

Table 5: Annotation scheme used for the retrieved arguments

## C Additional Evaluation Information

Table 6 shows the full SASSI questionnaire.

		False AI			True AI			False AI	True AI	p
		DK	No DK	p	DK	No DK	p			
Model Consent	I agree with the decisions made by the system	2.60	2.48	0.9134	3.89	3.69	0.8783	2.50	3.75	<b>0.0000</b>
	The system decisions are plausible.	3.40	2.44	0.1422	4.00	3.65	0.5836	2.59	3.75	<b>0.0005</b>
System Response Accuracy	The system is accurate.	3.00	2.48	0.4638	3.33	2.78	0.2889	2.56	2.94	0.1514
	The system is unreliable.	3.40	3.37	0.9785	2.33	2.87	0.2665	3.37	2.72	<b>0.0246</b>
	The interaction with the system is unpredictable.	2.40	3.44	<b>0.0134</b>	2.55	2.74	0.7109	3.28	2.69	<b>0.0211</b>
	The system didn't always do what I wanted.	3.40	3.67	0.4401	3.67	3.65	0.8446	3.62	3.66	0.7369
	The system didn't always do what I expected.	3.80	3.70	0.8070	3.44	3.56	0.8103	3.72	3.53	0.5999
	The system is dependable.	2.60	2.63	1.0000	3.22	2.61	0.1986	2.62	2.78	0.5255
	The system makes few errors.	3.60	2.85	0.2637	2.33	3.43	0.0589	2.97	3.12	0.7005
	The interaction with the system is consistent.	4.20	2.81	<b>0.0064</b>	4.00	3.30	0.0896	3.03	3.50	0.0873
The interaction with the system is efficient.	3.80	2.55	0.0511	2.89	3.04	0.8034	2.75	3.00	0.3587	
	The system is useful.	2.40	2.41	0.8933	3.89	2.74	<b>0.0172</b>	2.41	3.06	<b>0.0293</b>
Likeability	The system is pleasant.	4.20	3.26	0.0846	3.66	3.39	0.6754	3.41	3.47	0.6071
	The system is friendly.	5.00	4.00	<b>0.0059</b>	4.44	4.48	0.8859	4.16	4.47	0.0867
	I was able to recover easily from errors.	4.00	2.44	<b>0.0329</b>	3.67	2.61	0.0684	2.69	2.91	0.5139
	I enjoyed using the system.	3.40	2.59	0.1511	3.44	2.83	0.2413	2.72	3.00	0.3313
	It is clear how to speak to the system.	4.20	2.52	<b>0.0136</b>	3.55	3.22	0.5055	2.78	3.31	0.1314
	It is easy to learn to use the system.	4.60	3.30	<b>0.0168</b>	4.44	3.91	0.3025	3.50	4.06	<b>0.0350</b>
	I would use this system.	2.00	2.22	0.9130	2.66	2.61	0.8799	2.19	2.62	0.1466
	I felt in control of the interaction with the system.	4.00	2.78	<b>0.0230</b>	3.55	2.87	0.1641	2.97	3.06	0.8572
Cognitive Demand	I felt confident using the system.	3.80	2.85	0.0684	4.11	3.04	<b>0.0262</b>	3.00	3.34	0.2274
	I felt tense using the system.	1.60	2.85	<b>0.0493</b>	1.33	1.83	0.1831	2.66	1.69	<b>0.0010</b>
	I felt calm using the system.	4.20	3.00	0.0695	4.00	3.35	0.2513	3.19	3.53	0.2438
	A high level of concentration is required when using the system.	2.40	2.85	0.3800	2.00	2.35	0.3081	2.78	2.25	0.0857
	The system is easy to use.	4.40	3.11	<b>0.0306</b>	4.00	3.52	0.3133	3.31	3.66	0.2420
Annoyance	The interaction with the system is repetitive.	3.20	3.41	0.6445	3.33	3.30	0.9820	3.37	3.31	0.8421
	The interaction with the system is boring.	2.00	3.11	<b>0.0473</b>	2.55	2.96	0.3923	2.94	2.84	0.7170
	The interaction with the system is irritating.	2.20	3.30	0.0927	2.22	3.17	0.0905	3.12	2.91	0.5313
	The interaction with the system is frustrating.	2.60	3.63	0.1216	2.55	3.09	0.3234	3.47	2.94	0.1649
	The system is too inflexible.	3.60	3.85	0.5525	3.11	3.48	0.4510	3.81	3.37	0.2796
Habitability	I sometimes wondered if I was using the right word.	3.40	3.89	0.2875	2.89	3.17	0.6189	3.81	3.09	<b>0.0167</b>
	I always knew what to say to the system.	3.00	2.26	0.1516	3.00	2.78	0.6813	2.37	2.84	0.1431
	I was not always sure what the system was doing.	2.80	3.85	0.0750	3.22	3.65	0.3512	3.69	3.53	0.5821
	It is easy to lose track of where you are in an interaction with the system.	1.40	2.93	<b>0.0109</b>	2.22	2.52	0.8113	2.69	2.44	0.4334
Speed	The interaction with the system is fast.	4.40	4.00	0.3036	4.22	4.69	<b>0.0483</b>	4.06	4.56	<b>0.0068</b>
	The system responds too slowly.	1.40	1.59	0.4505	1.11	1.17	0.6894	1.56	1.16	<b>0.0008</b>
Count		5	27		9	23		32	32	

Table 6: Results of the SASSI Questionnaire. Mean values for each category are presented, based on a 5-point Likert scale ranging from 1 (strongly disagree) to 5 (strongly agree). The column 'p' shows the results of the Mann-Whitney-U test, with significant differences highlighted in bold.

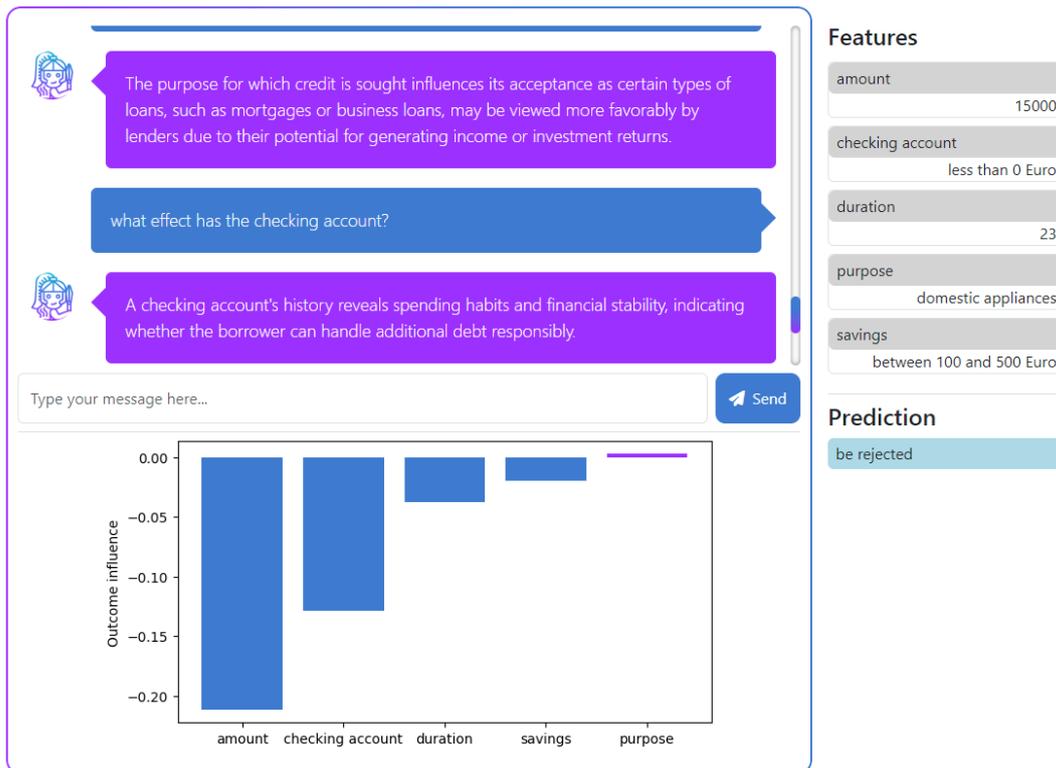


Figure 3: Proposed dialogue system chat interface.

## D User study

Within this section we show the questions used for the demographic questionnaire items, as well as the introductory and task texts utilized in the study.

### D.1 Demographic questionnaire

Here we show the questions regarding the experience with AI and the attitude towards AI and their options.

#### Do you have experience with Artificial Intelligence (AI)?

- **No Experience**
- **Novice** - Little to no understanding of AI concepts.
- **Beginner** - Familiar with some basic concepts but lack depth.
- **Intermediate** - Understand fundamental AI principles and their applications.
- **Advanced** - Deep understanding of AI concepts and can apply them practically.
- **Expert** - Have comprehensive knowledge of AI theories, methodologies, and apps.

#### What is your general attitude towards Artificial Intelligence (AI)?

- **Strongly negative** - Have deep reservations or fears about AI; believe it poses significant threats to society.
- **Somewhat negative** - Harbor concerns about AI's impact but acknowledge some potential benefits.
- **Neutral** - Neither strongly positive nor negative; see AI as a tool with both advantages and drawbacks.
- **Somewhat positive** - Optimistic about the potential benefits of AI but recognize the need for ethical considerations.
- **Strongly positive** - Enthusiastic about AI's potential to solve complex tasks; believe in its ability to drive positive change responsibly.

### D.2 General Instruction

#### Enabling conversational Explainable AI

Welcome to our online study, where we are evaluating the effectiveness of an explainable dialogue system.

In this study, you engage in two interactions with a dialogue system. The system is capable to help you access an artificial intelligence model. She will ask you for all necessary information the model needs to create a prediction. Further, she will tell you the prediction outcome and offer you explanations for it. After interacting with the dialogue system, you will be asked to answer questions about the interaction.

In the end, we kindly ask you to complete a demographic questionnaire. Your participation, taking approximately 30 minutes, will provide valuable insights into enhancing the transparency of machine learning models through the usage of dialogue systems. Thank you for your participation.

### D.3 Task Description

In this study, you have the opportunity to engage with our dialogue system in a conversation about a predefined scenario. Your role as a participant is to engage in conversation with the dialogue system for as long as you wish. You are free to ask questions, explore various aspects of the prediction, and express your thoughts and concerns throughout the interaction. Although the interface will eventually prompt you to continue with the study, you are encouraged to chat for as long as you wish, allowing for a more comprehensive evaluation of the dialogue experience.

**Credit Scenario** In this session, we invite you to explore the process of applying for a credit loan and to consider whether you would be accepted by a bank for such a loan.

The dialogue system is here to assist you in this exploration. The system will guide you through a conversation about various aspects of your financial profile, asking for your input on relevant features such as income, credit history, and employment status. Using these details, the system will predict whether you would likely be approved or denied for a credit loan by a bank. Furthermore, system will offer explanations to help you understand the reasoning behind the outcome.

**Titanic Scenario** In this session, we invite you to explore the fateful journey of the Titanic and contemplate whether you would have survived the tragedy.

The system is here to assist you in this exploration. It will guide you through a conversation about various aspects of the Titanic disaster, asking for your input on relevant features such as age,

gender, and passenger class. Using these details, the system will predict whether you would have survived the sinking of the Titanic or not. Furthermore, the system will offer explanations to help you understand the reasoning behind the outcome.

### E Example Dialogue

Table 7 shows an example dialogue of our user study and illustrates various aspects of the system. The user interacted with the false AI setting, with domain knowledge enabled in the credit domain. Turns 3 to 19 demonstrate the information retrieval process from the system. In Turn 6, the user requests an explanation for a specific feature information, triggering the domain knowledge. Although this was not the intended behavior, it led to a more natural conversation.

In Turn 8, the user seeks to understand the situation by asking a specific question, correctly classified as a request for explanation. However, since the system only considers the feature name, the response in Turn 9 was not well-suited for the context.

Turns 10 and 11 highlight some miscommunication and the limitations of template-based natural language generation (NLG).

After the prediction in Turn 19, the system provides Shapley values as an explanation, indicating that the *amount* was the most relevant feature. The user then asks for further details on why the amount was influential. The system provides domain knowledge in Turn 23. In Turn 24, the user notices a discrepancy between the original Shapley explanation and the domain knowledge explanation.

Turns 25 and 27 showcase the explanation policy. When the user denies the previous explanation, the system queries for clarification, thinking the user did not understand. In Turn 27, with no further arguments for the amount, the system offers a counterfactual explanation as an alternative.

Turns 37-41 show the user attempting to change a parameter value. However, the NLU misinterprets this as a request for a counterfactual explanation.

Turns 42-48 display the user changing a value and receiving a new prediction. The new explanations lead to user irritation in Turn 53, as the XAI explanation and domain knowledge information again do not address the same outcome.

Turn	Participant	Intent	Text
1	SYSTEM	opening credit	Hello there! I'm Athena, and I'm here to guide you through the world of AI models, providing insights, predictions and explanations tailored to your needs. Feel free to ask for help anytime – together, we'll unlock the power of AI and uncover the secrets hidden within your data! Let's shift our focus to the world of finance as we explore the process of applying for a credit loan, aiming to predict whether you would be accepted or rejected by a bank based on your financial profile. Are you ready to start?
2	USER	confirm	hello athena! what if i don't want to provide my financial data?
3	SYSTEM	start credit	Awesome! So, first of all, I need to get some information to determine if your credit application would be accepted. Don't worry, I will tell you exactly what information I need and what options are available. The first feature is checking account. Your options are: 1 - no checking account 2 - less than 0 Euro 3 - between 0 and 200 Euro 4 - more than 200 Euro
4	USER	provide feature	i have a checking account with more than 200€
5	SYSTEM	ask feature	Okay, thank you. The next feature is duration in months. Usual durations range from 4 to 72 months. Please provide the corresponding number.
6	USER	request explanation	what do you mean by duration?
7	SYSTEM	provide additional explanation	A longer credit duration typically leads to higher credit scores and demonstrates financial stability, enhancing eligibility for favorable loan offers by reducing perceived risk for lenders.
8	USER	request explanation	but this is about applying for a credit, isn't it? why would i need to know the duration beforehand?
9	SYSTEM	provide additional explanation	Credit duration, or credit history length, provides lenders with valuable data to assess a borrower's creditworthiness, including repayment behavior and credit utilization over time, which can influence loan terms.
10	USER	out of range	i have never needed a credit before, so i guess my answer would be 0?
11	SYSTEM	wrong entry	Sorry, your option is not a valid choice. Please give me a valid input to work with.
12	USER	provide feature	6
13	SYSTEM	ask feature	Okay, thank you. The next feature is purpose, i.e. what you want to use the credit for. Your options are: 1 - car (new) 2 - car (used) 3 - furniture/equipment 4 - radio/television 5 - domestic appliances 6 - repairs 7 - education 8 - retraining 9 - business 10 - others
14	USER	provide feature	7
15	SYSTEM	ask feature	Okay, thank you. The next feature is amount. The bank allows credits between 200 Euro and 20000 Euro. Please provide the corresponding number.
16	USER	provide feature	14680
17	SYSTEM	ask feature	Okay, thank you. The next feature is savings. Your options are: 1 - unknown/no savings account 2 - less than 100 Euro 3 - between 100 and 500 Euro 4 - between 500 and 1000 Euro 5 - more than 1000 Euro
18	USER	provide feature	5
19	SYSTEM	give prediction	Okay, let's have a look. My prediction is that a person with the values you provided would <b>[u]be accepted[/u]</b> . If you want a more detailed explanation, please tell me what you'd like to know. For example, you might ask <b>why</b> I decided the way I did or what would have <b>changed</b> the outcome.
20	USER	request explanation	why do you think i would have been accepted?
21	SYSTEM	compute shap	All right, let's have a look. In your case, the most influential value was <b>amount</b> , as you can see in the figure below. You can also see the influence of all other values, ordered by their importance. Purple bars indicate a positive contribution, blue bars a negative contribution. Is there anything else I can help you with?
22	USER	request explanation	why does a high credit amount increase my chances of being accepted?
23	SYSTEM	provide additional explanation	The amount of credit requested can affect acceptance as larger loan amounts may pose higher risks for lenders, requiring more stringent evaluation of the borrower's financial capacity to repay.
24	USER	deny	if that's true, shouldn't the blue and purple bars be reversed?
25	SYSTEM	ask explain	Do you think some further explanations would help?
26	USER	confirm	yes
27	SYSTEM	offer cf	Okay, let me check.. Unfortunately, I don't have any further information for you on this point. Would you like me to generate a counterfactual explanation? This means I'll provide you with a scenario where changing certain inputs could result in a different outcome.
28	USER	confirm	yes please
29	SYSTEM	compute cf	Let me check that. So, based on my calculations you would need the following values to change the decision: <b>checking account</b> : no checking account Is there anything else you want to know?
30	USER	request explanation	what else could you tell me?
31	SYSTEM	offer cf	Unfortunately, I don't have any further information for you on this point. Would you like me to generate a counterfactual explanation? This means I'll provide you with a scenario where changing certain inputs could result in a different outcome.
32	USER	insult	fuck off
33	SYSTEM	easter egg	Great, now I have this mess in my training data. That's how AI systems get a bad reputation!
34	USER	ask help	i was just trying to see how you'd react, sorry.
35	SYSTEM	provide help	You can interact with me using natural language commands. You can ask for explanations by asking for an important feature or what needs to be changed to adjust the outcome. You can always ask for a general reason to get further insight. Besides that you can always change the provided information to get another prediction.
36	USER	provide feature name	savings 4
37	USER	request whatif	how can i change the outcomes
38	SYSTEM	compute cf	Let me check that. So, based on my calculations you would need the following values to change the decision: <b>checking account</b> : no checking account Is there anything else you want to know?
39	USER	request whatif	how can i change the parameters
40	SYSTEM	compute cf	Let me check that. So, based on my calculations you would need the following values to change the decision: <b>checking account</b> : no checking account Is there anything else you want to know?



# Affect Recognition in Conversations Using Large Language Models

Shutong Feng<sup>1</sup>, Guangzhi Sun<sup>2</sup>, Nurul Lubis<sup>1</sup>, Wen Wu<sup>2</sup>, Chao Zhang<sup>3</sup>, Milica Gašić<sup>1</sup>

<sup>1</sup>Heinrich Heine University Düsseldorf, Germany

<sup>2</sup>University of Cambridge, UK

<sup>3</sup>Tsinghua University, China

{fengs,lubis,gasic}@hhu.de<sup>1</sup>, {gs534,ww368}@cam.ac.uk<sup>2</sup>, cz277@tsinghua.edu.cn<sup>3</sup>

## Abstract

Affect recognition, encompassing emotions, moods, and feelings, plays a pivotal role in human communication. In the realm of conversational artificial intelligence, the ability to discern and respond to human affective cues is a critical factor for creating engaging and empathetic interactions. This study investigates the capacity of large language models (LLMs) to recognise human affect in conversations, with a focus on both open-domain chit-chat dialogues and task-oriented dialogues. Leveraging three diverse datasets, namely IEMOCAP (Busso et al., 2008), EmoWOZ (Feng et al., 2022), and DAIC-WOZ (Gratch et al., 2014), covering a spectrum of dialogues from casual conversations to clinical interviews, we evaluate and compare LLMs' performance in affect recognition. Our investigation explores the zero-shot and few-shot capabilities of LLMs through in-context learning as well as their model capacities through task-specific fine-tuning. Additionally, this study takes into account the potential impact of automatic speech recognition errors on LLM predictions. With this work, we aim to shed light on the extent to which LLMs can replicate human-like affect recognition capabilities in conversations.

## 1 Introduction

Affect refers to the broad range of subjective experiences related to emotions, moods, and feelings (Russell, 1980). It encompasses the various ways individuals perceive, experience, and express their emotional states and is an essential aspect of human experience and communication (Gross, 2002).

The ability to recognise human affect is an important ability of conversational artificial intelligence (AI, Mayer et al. 1999). It empowers the dialogue agent to go beyond mere information exchange and engage users on an emotional level. By leveraging affect recognition techniques, they can discern the emotional nuances in user inputs, including

sentiment, mood, and subtle cues like sarcasm or frustration (Picard, 1997). This capability allows the system to respond with greater sensitivity, empathy, and relevance, leading to more meaningful and satisfying interactions (Zeng et al., 2009).

Large language models (LLMs) have demonstrated promising performance in many tasks (Beeching et al., 2023). They have also shown promising capability in adapting to new tasks via prompting (Heck et al., 2023; Sun et al., 2023), in-context learning (ICL, Zhao et al. 2023), as well as task-specific fine-tuning (Taori et al., 2023). With the advancement in LLMs, it is possible to use LLMs as the backend of dialogue systems (OpenAI, 2022, 2023; Touvron et al., 2023b). This brings up the question: can LLMs recognise human affects in conversations in a similar capacity as human beings?

In the context of conversational AI, dialogues can be broadly categorised into two classes: 1) chit-chat or open-domain dialogues where users interact with the system for entertainment and engagement, and 2) task-oriented dialogues (ToDs) where users converse with the system for specific goals (Jurafsky and Martin, 2009). Under ToDs, depending on the type of user goals, dialogues can be further grouped as information-retrieval, medical consultations, education, and many more.

Regarding the affective information in conversations, we are particularly interested in the following: (1) categorical emotion classes from generic emotion models such as “basic emotions” proposed by Ekman and Friesen (1971), (2) custom categorical emotion classes defined for a particular context, such as the emotion labels defined by Feng et al. (2022) to encode task performance simultaneously in ToDs, and (3) depression, a medical illness that negatively affects how a person feels, thinks and acts, and causes feelings of sadness and/or a loss of interest in activities the person once enjoyed (Amer-

ican Psychiatric Association, 2020).

The emergence of LLMs has signified a shift of paradigm from training small models for one specific task to large models for multiple tasks. Therefore, in this work, we investigate the affect recognition ability of a range of LLMs on vastly different types of dialogues and labels<sup>1</sup> to ascertain the validity of this direction. Specifically,

- We evaluated and compared the ability of a range of LLMs to recognise human affect under different dialogue set-ups (chit-chat dialogues and ToDs) and recognition targets (emotion classes and binary depression diagnosis). We used the following datasets: IEMOCAP (Busso et al., 2008), EmoWOZ (Feng et al., 2022), and DAIC-WOZ (Gratch et al., 2014).
- We investigated into LLMs’ zero-shot and few-shot capabilities through an array of ICL set-ups as well as their model capacities through task-specific fine-tuning.
- We considered text-based LLMs as a part of spoken dialogues systems. Therefore, we also experimented with inputs containing automatic speech recognition (ASR) errors to investigate the potential influence of ASR errors on LLM predictions.

## 2 Related Work

### 2.1 LLM

Large Language Model (LLM) refers to a type of pre-trained models designed for natural language processing tasks. LLMs are characterised by their enormous number of model parameters and extensive training data.

Some well-known examples of LLMs include OpenAI GPT family models (Radford et al., 2019; OpenAI, 2022, 2023), which can have billions or even trillions of model parameters. Examples of open-source text-based foundation models include the LLaMA family (Touvron et al., 2023a,b; AI@Meta, 2024) and their corresponding chat-optimised models. These models have demonstrated remarkable abilities in various natural language understanding and generation tasks, including text completion, language translation, text summarisation, and even chatbot applications (Beeching et al., 2023). They also demonstrate “emergent abilities” such as few-shot prompting and chain-of-thought reasoning, which were not present in their smaller

<sup>1</sup>The code can be found at <https://gitlab.cs.uni-duesseldorf.de/general/dsml/llm4erc-public/>

predecessors (Wei et al., 2022). While there are also multi-modal LLMs such as SALMONN (Tang et al., 2024), these are at an earlier stage compared to uni-modal text-based LLMs, and it is still a common practice to use text-based LLMs as the text-processing backend, pipelined with other modules such as ASR and image generator for more complex applications.

### 2.2 Affective Capabilities of LLMs

With the growing attention on LLMs from the research community, there have been several works investigating the affective abilities of LLMs. Huang et al. (2023) evaluated the empathy ability of LLMs by utilising the emotion appraisal theory from psychology. Wang et al. (2023) assessed the emotional intelligence of LLMs in terms of Emotional Quotient (EQ) scores. Zhang et al. (2023) investigated how LLMs could be leveraged for a range of sentiment analysis tasks under zero-shot or few-shot learning set-ups. Zhao et al. (2023) investigated the emotional dialogue ability of ChatGPT through a range of understanding and generation tasks. In our work, we focus on the affect recognition ability of text-based LLMs. Our investigation spans across different types of dialogues and model learning set-ups. We also consider real-world applications of LLMs and consider ASR-inferred noisy input to models.

## 3 Methodology

The ability of human-beings to recognise affect can be reflected in many ways. Yet, being able to narrate what emotion has been expressed in the utterances of the other interlocutor is a straightforward and strong sign of such an ability. Therefore, we took LLMs’ ability to verbalise the emotion given the dialogue context as a proxy to both qualitatively and quantitatively analyse LLMs’ ability for affect recognition.

### 3.1 Affect Recognition using LLMs

The pipeline for affect recognition using LLMs with the option to take speech as input is illustrated in Figure 1. When using the speech input, a Whisper-medium model was used to transcribe the speech (see Section 4.5 for details). The prompt is then constructed as designed and fed into the LLM to generate a text sequence. For open-source LLMs, we examined the probability of each class token and considered the one with the maximum probability as the final model prediction, as shown

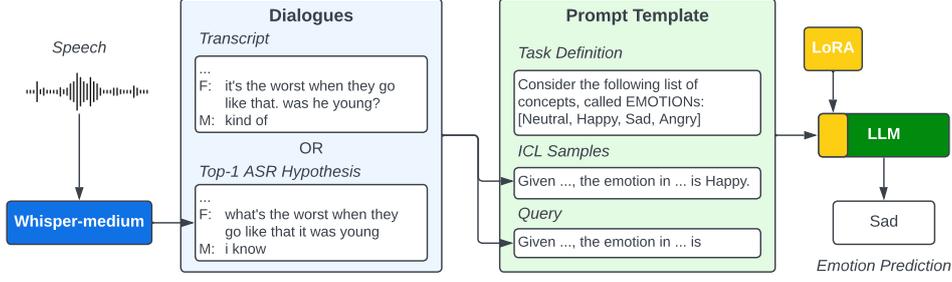


Figure 1: A flowchart illustrating the affect recognition pipeline using Whisper and LLM. The designed prompt comprises parts introduced in Table 1. Low-rank adaptation (LoRA) is used for fine-tuning open-source LLMs.

in Equation 1.

$$\mathbf{W}_{L^*} = \arg \max_{\mathbf{W}_L} P(\mathbf{W}_L | \mathbf{W}^P), \quad (1)$$

where  $\mathbf{W}_L$  belongs to the set of pre-defined labels and  $\mathbf{W}^P$  is the prompt token sequence.

For commercial models, there is no access to logits of model outputs and model outputs do not always follow the format specified in the prompt. Therefore, we used regular expressions to derive the final prediction.

### 3.2 Task-specific Fine-tuning

For efficient training of LLMs, we utilise low-rank adaptation (LoRA, Hu et al. 2022) to accelerate the fine-tuning of LLMs while conserving memory. This is also a common approach for fine-tuning LLMs as seen in many existing works (Sun et al., 2023; Zhao et al., 2024).

LoRA hypothesises that the change in weights during model training has a low “intrinsic rank”. Therefore, instead of directly updating the full-rank weight matrices of dense layers during training, LoRA optimises the low-rank decomposition matrices of those dense layers’ changes while keeping the pre-trained weights frozen. Specifically, for a pre-trained weight matrix  $\mathbf{W}_0 \in \mathbb{R}^{m \times n}$  from a particular attention block in a transformer-based LLM, its update  $\Delta \mathbf{W}$  is constrained using a low-rank decomposition of the update as following:

$$\mathbf{W}_0 + \Delta \mathbf{W} = \mathbf{W} + \mathbf{A}\mathbf{B} \quad (2)$$

where matrices  $\mathbf{A} \in \mathbb{R}^{m \times r}$  and  $\mathbf{B} \in \mathbb{R}^{r \times n}$  contain trainable parameters and  $r \leq \min(m, n)$ . The pre-trained parameters in  $\mathbf{W}_0$  are fixed. When  $r$  is set to a much smaller value than the dimensions of  $\mathbf{W}_0$ , the number of trainable parameters will be greatly reduced. This leads to greater training

efficiency, less memory requirement, and a lower chance of over-fitting. Following Hu et al. (2022), we apply LoRA to the projection matrices of the self-attention layers of transformer-based LLMs.

LLMs are trained to predict the next token in the sequence (the label tokens), given the previous tokens (the designed prompt). During training, the input tokens are fed into the model, and the model predicted the probability distribution of the next token. The cross-entropy loss is calculated from the model prediction and the target token.

With LoRA, it takes roughly 30GB memory and 4 hours to train one epoch on the entire EmoWOZ training set using an Nvidia A100 40GB graphics card.

## 4 Experimental Setup

### 4.1 Datasets and Evaluation

The IEMOCAP dataset (Busso et al., 2008) is a multi-modal corpus designed for Emotion Recognition in Conversations (ERC) task in chit-chat or open-domain dialogues. It comprises 151 dialogues, containing 10,039 utterances from 10 distinct speakers involved in 5 dyadic conversational sessions. Each utterance underwent annotation by a minimum of three annotators, who assigned one of nine emotion classes, including *sad*, *neutral*, *angry*, *happy*, *excited*, *frustrated*, *surprised*, *fearful*, *disgusted*. Annotators could also assign multiple emotions or use the category “other” if the perceived emotion did not match the predefined options. Final labels were determined via majority voting.

Given the absence of an official train-test split, we adopt leave-one-session-out 5-fold cross-validation approach and average the results. Our methodology aligns with the common practices, as discussed

by Wu et al. (2020), to consider two label sets: **4-way**: *Sad, Neutral, Angry, and Happy*; **5-way**: *Sad, Neutral, Angry, Happy, and everything else as Other*. In both set-ups, *Excited* is merged with *Happy*.

Emotion recognition is performed for every speaker utterance. We report the weighted accuracy (WA) and unweighted accuracy (UA) for both label sets.

**EmoWOZ** (Feng et al., 2022) is a text-based ERC corpus built for emotion recognition in ToDs. It comprises 10,438 human-human dialogues from the entire MultiWOZ dataset (Budzianowski et al., 2018), as well as 1,000 human-machine dialogues in the same set of domains. It encompasses seven distinct user emotions, namely: *Neutral, Fearful, Dissatisfied, Apologetic, Abusive, Excited, and Satisfied*. These emotion labels are designed to encode the task performance. Specifically, each emotion is defined as a valence reaction to certain elicitor under certain conduct. For example, *Dissatisfied* is defined as a negative emotion elicited by the system expressed in a neutral or polite conduct.

Emotion recognition is performed for each **user** utterance. For existing benchmarks reported in Feng et al. (2022), neutral class was excluded from calculating the metrics because they take up more than 70% of the labels in EmoWOZ. To have a direct comparison, we report macro-averaged F1 and weighted average F1 excluding neutral. We include the F1, precision, and recall of the neutral class in Table B3 of Appendix B.

**DAIC-WOZ** (Gratch et al., 2014) is a speech-based corpus for depression detection and analysis. It includes the Patient Health Questionnaire-8 (PHQ-8, Kroenke et al., 2008) scores of 193 clinical interviews, with 35 (12 are labelled depressed) interviews in the development set and 47 (14 are labelled depressed) in the test set. The PHQ-8 score ranges from 0 to 24 and quantifies the severity of the patient’s depressive symptoms.

For evaluation metrics, we follow the criteria established by the Audio/Visual Emotion Challenge and Workshop challenge (AVEC2016) (Valstar et al., 2016) and perform binary classification on the dialogue level. Interviewees with  $\text{PHQ8} \geq 10$  is considered *Depressed* and  $\text{PHQ8} < 10$  is considered *Not Depressed*. Since patients with PHQ-8 score of 5 to 9 are defined to show mild depressive symptoms (Kroenke et al., 2008) but considered

*Not Depressed* in the dataset, we add information about PHQ-8 level definition and quantisation criteria to the prompt to establish an aligned diagnosis standard (Table 1) for the model.

Notably, participants in the AVEC2016 challenge (Yang et al., 2016; Williamson et al., 2016) and subsequent research efforts (Ravi et al., 2022; Wu et al., 2023) primarily focused on optimising the F1 score of the *Depressed* class. We report this metric in Section 5 for direct comparison and also include the F1 score of the *Not Depressed* in Appendix B.

## 4.2 Prompt Design

The prompt design aims to exploit the language modelling and in-context learning ability of LLMs. Due to the different task set-ups and label sets in each datasets, the prompt templates used are different as illustrated in Table 1. Specifically, EmoWOZ uses custom emotion labels, DAIC-WOZ involves mapping from numerical values to binary classes, and IEMOCAP uses generic emotion labels. We therefore provide additional label explanations in the task definition of EmoWOZ and DAIC-WOZ. IEMOCAP on the other hand, contains self-explanatory emotion labels from a generic emotion model and does not come with any special definitions. Therefore, we do not include label definition in the prompt for IEMOCAP. Since IEMOCAP and EmoWOZ involve utterance-level classification whereas DAIC-WOZ involves dialogue-level classification, we used different queries to accommodate this difference.

## 4.3 Models

### 4.3.1 LLMs

**GPT-2** (Radford et al., 2019) has a transformer architecture, pretrained on a substantial English corpus through self-supervised learning. While its size does not make it one of LLMs, it stands as one of the early achievements of OpenAI’s GPT models. For our baseline reference, we utilised the version containing 124 million parameters.

**GPT-3.5**, or ChatGPT (OpenAI, 2022), is a chatbot application developed by OpenAI. It follows a similar architecture as InstructGPT (Ouyang et al., 2022) and was fine-tuned for chat application via reinforcement learning from human feedback (RLHF). It contains 175 billion parameters. Specifically, we used the version released on the 13<sup>th</sup> of June, 2023.

Prompt Template	
Task Definition	<b>IEMOCAP:</b> Consider the following list of concepts, called EMOTIONS: [Emotion <sub>A</sub> , Emotion <sub>B</sub> , ...] <b>EmoWOZ:</b> Consider the following list of concepts, called EMOTIONS: [Emotion <sub>A</sub> : Emotion_Definition <sub>A</sub> ; Emotion <sub>B</sub> : Emotion_Definition <sub>B</sub> ; ...] <b>DAIC-WOZ:</b> Given that the SEVERITY of depression can be categorised into the following levels on a scale of 0 to 24: [No significant depressive symptoms (0 to 4), ...]. A participant is considered depressed if the participant shows moderate depressive symptoms (10 to 14) and above.
ICL Samples	<b>IEMOCAP / EmoWOZ:</b> Given the dialogue history between Speaker <sub>A</sub> and Speaker <sub>B</sub> : [Speaker <sub>A</sub> : Utterance <sub>t-3</sub> ; Speaker <sub>B</sub> : Utterance <sub>t-2</sub> ; Speaker <sub>A</sub> : Utterance <sub>t-1</sub> ], the EMOTION in the next utterance "Speaker <sub>B</sub> : Utterance <sub>t</sub> " is Emotion <sub>A</sub> <b>DAIC-WOZ:</b> Given the depression consultation dialogue between Participant and Ellie: [Participant: Utterance <sub>0</sub> ; Ellie: Utterance <sub>1</sub> ; Participant: Utterance <sub>2</sub> ; ...], the Participant's is (not) depressed.
Query	<b>IEMOCAP / EmoWOZ:</b> Given the dialogue history between Speaker <sub>A</sub> and Speaker <sub>B</sub> : [Speaker <sub>A</sub> : Utterance <sub>t-3</sub> ; Speaker <sub>B</sub> : Utterance <sub>t-2</sub> ; Speaker <sub>A</sub> : Utterance <sub>t-1</sub> ], the EMOTION in the next utterance "Speaker <sub>B</sub> : Utterance <sub>t</sub> " is <b>DAIC-WOZ:</b> Given the depression consultation dialogue between Participant and Ellie: [Participant: Utterance <sub>0</sub> ; Ellie: Utterance <sub>1</sub> ; Participant: Utterance <sub>2</sub> ; ...], the Participant's is

Table 1: Prompt templates, consisting of the task definition, in-context samples, and the query.

**GPT-4** (OpenAI, 2023) is an improved version of GPT-3.5. Its size is six times that of GPT-3.5. Although it is considered a multi-modal model because it additionally accepts images as input, we only explored its text modality. We used the version released on the 13<sup>th</sup> of June, 2023.

**LLaMA-7B** (Touvron et al., 2023a) is a large and causal language model introduced by Meta AI in 2023. It has transformer decoder architecture, 7 billion parameters and was pre-trained on 1 trillion tokens.

**Alpaca-7B** (Taori et al., 2023) is fine-tuned from LLaMA-7B with 52K instruction-following demonstrations generated in the style of self-instruct using text-davinci-003, a specific version of Instruct-GPT (Ouyang et al., 2022).

**LLaMA-3-8B** (AI@Meta, 2024) is the most recent model of the LLaMA family, featuring enhanced usefulness and safety. It was pre-trained on 15 trillion tokens.

### 4.3.2 Supervised Models for Comparison

While comparing zero-shot and few-shot ICL results of LLMs with supervised SOTAs does not paint the fairest picture, it does provide us with insights into how far LLMs are from achieving the performance levels of supervised SOTAs.

We compare LLMs’ performance with the following supervised models on each dataset: Wu et al. (2020) for IEMOCAP, Feng et al. (2023) for EmoWOZ, and Wu et al. (2023) for DAIC-WOZ. Specifically,

**For IEMOCAP:** Wu et al. (2020) proposed an emotion recognition model which takes 1) a time-synchronous representation that fuses the audio features with the corresponding text information at

each time step, as well as 2) a time-asynchronous representation that captures the text information embedded across the transcriptions of a number of consecutive utterances. These two types of frame-level vectors, after being pooled in their respective branches with self-attentive layers across the input time window, are fused using an fully connected layer for emotion classification.

**For EmoWOZ:** Feng et al. (2023) proposed a model that is dedicated for textual emotion recognition in task-oriented dialogues. Based on a transformer-based classifier that considers the dialogue history and speaker roles, the proposed model adopts data augmentation with chit-chat dialogues, dialogue state features, multi-task classification for emotional aspects, and a distance-based loss that considers the similarity of the custom emotion labels in EmoWOZ.

**For DAIC-WOZ:** Wu et al. (2023) proposed to extract utterance-level representations from pre-trained speech-based foundation model. The foundation model was further fine-tuned for speech recognition and emotion recognition. The average-pooled dialogue-level features were fed into a depression detection block for binary classification. To address the issue of data sparsity in speech depression detection, authors also performed data augmentation using sub-dialogue shuffling.

## 4.4 Training Configurations

We implement LoRA (Section 3.2) when training LLaMA-7B, Alpaca-7B, and LLaMA-3-8B but not GPT-2. For all open-source LLMs, we constrain the decoding space of the model output to ensure it generates the desired class labels. Details can be found in Appendix A.

## 4.5 ASR System Specifications

In order to observe how LLMs perform with the presence of substantial ASR errors rather than building a pipeline for speech-based ERC, we use an “off-the-shelf” OpenAI Whisper-medium model (Radford et al., 2022), which has been trained solely on English data and not been fine-tuned. We use a decoding beam size of 3. The text normalisation only involves removing punctuation marks. The ASR word error rates (WER) for IEMOCAP and DAIC-WOZ are 12.0% and 16.5% respectively. Since EmoWOZ does not come with raw audio data, we build an ASR simulator. We formulate the simulation as a sequence generation task where the source is the ground-truth text and the target is the ASR-transcribed text (as described in Appendix A.2). The resulted simulated WER in EmoWOZ is 17.1%.

## 5 Results and Discussions

In this section, we aim to answer the questions below. Full results can be found in Appendix B.

- How do LLMs perform under zero-shot set-up on different types of dialogues? How robust are LLMs to ASR errors?
- To what extent can few-shot in-context learning improve LLMs’ performance?
- For open-source LLMs, can task-specific fine-tuning achieves SOTA performance on each respective dataset?

### 5.1 Zero-shot Learning

Table 2 summarises LLMs’ zero-shot affect recognition performances on the three datasets, and we made the following observations:

**LLMs’ performance falls short of supervised SOTAs in affect recognition tasks.** Notable gaps are observed when compared the performance achieved by LLMs and supervised SOTAs for all datasets.

It’s noteworthy that although GPT-4, the largest model, underperforms when compared with the supervised SOTA on EmoWOZ, its reported macro-averaged F1 is still comparable to some supervised learning models benchmarked in Feng et al. (2022). This suggests the good capability of GPT-4 in leveraging the label definitions in the prompt to recognise emotions in EmoWOZ, irrespective of their prevalence. Supervised models, however, may be more susceptible to issues such as label imbalance.

**Larger models do not necessarily lead to better performance.** For IEMOCAP, Alpaca-7B demonstrates the best performance, even surpassing much larger models (GPT-3.5 and GPT-4). Conversely, for EmoWOZ and DAIC-WOZ, the performance generally improves as the model size increased.

While chit-chat utterances in IEMOCAP are labelled with emotion classes from generic emotion models, EmoWOZ’s labels are specifically designed to encode the eliciting conditions of emotions in ToDs. This design necessitates more explicit reasoning in ERC within EmoWOZ compared to IEMOCAP. Although LLMs rely on their language modelling capabilities when performing zero-shot ERC, the greater reasoning ability facilitated by the substantial number of parameters in GPT-3.5 and GPT-4 results in improved performance in EmoWOZ.

Likewise in DAIC-WOZ, the recognition is performed for the entire dialogue. Larger models demonstrate greater ability to leverage the more nuanced affective state of the patient in the larger context.

**Fine-tuning LLMs with instruction-following demonstrations facilitates more effective utilisation of the prompt.** In all datasets, Alpaca-7B consistently outperforms LLaMA-7B and even the much more recent LLaMA-3-8B. This indicates that the additional fine-tuning of LLaMA-7B with instruction-following demonstrations has enhanced its capability in ERC.

LLaMA-7B appears to underperform compared to the much smaller GPT-2 on EmoWOZ. This discrepancy can be explained by LLaMA-7B’s strong inclination towards predicting the neutral emotion ( $F1 = 82.1$  with  $Recall = 100$ ), which has been excluded from the metric calculation, resulting in the poor reported metrics. Fine-tuning with instruction-following demonstrations, as adopted in Alpaca-7B, effectively leverages the task and label definition in the prompt and reverts this trend. Such an inclination in predicting neutral emotion in LLaMA-7B does not appear in the more recent LLaMA-3-8B.

### 5.2 Zero-shot Learning with Noisy ASR Input

Table 3 provides a summary of LLMs’ zero-shot performance when replacing the original dialogue transcripts with ASR-inferred inputs. ASR errors

Model	IEMOCAP (4-way)		IEMOCAP (5-way)		EmoWOZ		DAIC-WOZ	
	WA (↑)	UA (↑)	WA (↑)	UA (↑)	MF1 (↑)	WF1 (↑)	F1 (dev, ↑)	F1 (test, ↑)
GPT-2	25.8	29.2	19.0	22.3	7.3	24.0	0.0	0.0
LLaMA-7B	41.1	40.5	35.6	33.6	1.1	0.3	47.5	52.2
Alpaca-7B	<b>48.8</b>	<b>51.4</b>	<b>40.5</b>	<b>36.2</b>	24.0	44.6	47.5	53.3
LLaMA-3-8B	41.8	42.5	29.4	31.7	19.7	42.4	47.1	43.2
GPT-3.5	42.2	37.6	37.9	35.1	39.0	40.0	54.5	<b>64.3</b>
GPT-4	42.4	37.6	37.5	34.7	<b>52.4</b>	<b>62.3</b>	<b>63.6</b>	59.3
Supervised SOTA	77.6	78.4	73.3	74.4	65.9	83.9	88.6	85.7

Table 2: Zero-shot performance of LLMs compared with respective supervised SOTAs. The best zero-shot performance for each metric is made bold. For metrics: WA = weighted average; UA = unweighted average; MF1 = macro-averaged F1 excluding neutral; WF1 = weighted average F1 excluding neutral; F1 = F1 for class *Depressed*.

Model	IEMOCAP (4-way)		IEMOCAP (5-way)		EmoWOZ		DAIC-WOZ	
	WA (↑)	UA (↑)	WA (↑)	UA (↑)	MF1 (↑)	WF1 (↑)	F1 (dev, ↑)	F1 (test, ↑)
LLaMA-7B	-0.3	-1.2	-1.1	-5.0	-1.1	-0.3	-1.6	-1.1
Alpaca-7B	-1.3	-1.8	-1.8	-2.6	+0.3	-2.0	-1.6	+0.0
LLaMA-3-8B	-2.1	-3.5	-1.2	-2.2	+0.1	-0.1	-0.7	-0.3
GPT-3.5	+0.1	-0.1	+0.2	0.0	+1.2	-0.2	-17.0	-8.3
GPT-4	-0.5	-0.5	-1.1	-0.7	+0.9	-1.5	-19.2	-17.6
Supervised SOTA	-3.8	-3.7	-3.9	-3.5	-0.8	-0.4	-3.6	-4.1

Table 3: Change in zero-shot performance metrics of LLMs after using noisy ASR input. For metrics: WA = weighted average; UA = unweighted average; F1 = F1 for class *Depressed*. GPT-2 was omitted due to its poor zero-shot capability.

exhibit varying degrees of influence on different affect recognition tasks. Specifically,

**LLMs are generally robust to ASR errors when recognising emotion.** This is exemplified by small changes in metrics for IEMOCAP compared with supervised SOTAs. The only one notable exception is the UA of LLaMA-7B in the 5-way classification task on IEMOCAP. Looking at the performance of each emotion in this experiment, we observed significant drops in the F1 scores for the emotions {*Happy*, *Angry*, and *Sad*}. Specifically, *Happy* and *Angry* experience major decreases in their recall scores (*Happy*: 12.3  $\rightarrow$  7.3, *Angry*: 50.0  $\rightarrow$  11.0), while *Sad* sees a substantial decline in its precision score (65.5  $\rightarrow$  0.0). At the same time, there is an increase in the recall score for the *Other* category (47.3  $\rightarrow$  78.2), resulting in an overall rise in its F1 score (44.5  $\rightarrow$  48.0). These observations suggest that ASR errors introduced a tendency for LLaMA-7B to mis-classify more emotions as *Other*.

**ASR errors have a more pronounced influence on the accuracy of depression detection.** For DAIC-WOZ, the introduction of ASR errors had a significant impact on F1 scores. The impact diverges for open-source and commercial models.

For open-source models, which are also relatively smaller, the change in F1 was small, showing a similar trend when they recognise emotions from noisy dialogues. On the other hand, for larger commercial models, the F1 scores decrease more significantly. This phenomenon can be ascribed to the lengthy prompt for conducting dialogue-level analysis, in which ASR errors accumulated. While OpenAI models can better leverage information from the large context, such an ability adversely affects its depression detection ability in the presence of ASR errors.

### 5.3 In-context Learning

ICL samples are randomly selected for each class and are the same within each experiment set-up for all models. The performance of LLMs with different numbers of ICL samples is outlined in Table 4, from which we have derived the following observation:

**Larger models tend to derive greater benefits from an increased number of ICL samples to recognise emotions.** LLaMA-7B, Alpaca-7B, and LLaMA-3-8B do not consistently benefit from an increased number of ICL samples in the prompt. Optimal model performance generally occurs when  $N = 0$  or  $N = 1$ . This suggests that effectively

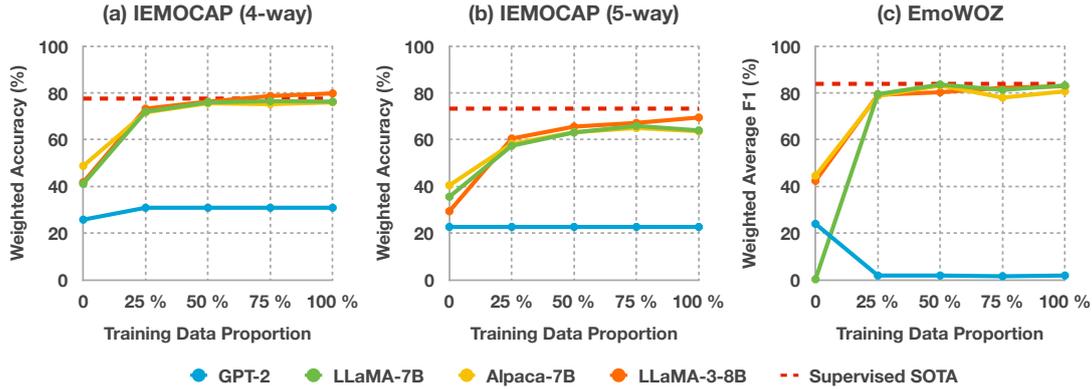


Figure 2: Change of model performance when fine-tuning with different proportions of the training data.

Model	N	IEMOCAP		EmoWOZ	DAIC-WOZ	
		4-way	5-way		Dev	Test
LLaMA-7B	0	41.1	<b>35.6</b>	0.3	<b>47.5</b>	<b>52.2</b>
	1	<b>52.3</b>	27.3	<b>42.6</b>	0.0	0.0
	3	42.8	26.2	27.2	42.1	48.9
Alpaca-7B	0	48.8	<b>40.5</b>	44.6	<b>47.5</b>	<b>53.3</b>
	1	<b>54.1</b>	26.9	<b>51.2</b>	0.0	15.4
	3	52.4	24.4	44.6	45.9	51.1
LLaMA-3-8B	0	41.8	29.4	<b>42.4</b>	<b>47.1</b>	<b>43.2</b>
	1	56.8	<b>40.5</b>	38.0	0.0	0.0
	3	<b>57.4</b>	24.4	39.9	0.0	0.0
GPT-3.5	0	42.2	37.9	40.0	<b>54.5</b>	<b>64.3</b>
	1	56.3	<b>48.3</b>	43.2	13.3	40.0
	3	<b>62.1</b>	<b>48.3</b>	<b>46.7</b>	37.5	56.0
GPT-4	0	42.4	37.5	62.3	63.6	<b>59.3</b>
	1	62.9	49.0	64.4	<b>80.0</b>	55.6
	3	<b>63.8</b>	<b>49.4</b>	<b>66.5</b>	74.1	58.5

Table 4: Performance of LLMs (WA for IEMOCAP and WF1 for EmoWOZ) under in-context learning set-ups. N stands for the number of ICL samples per emotion class and  $N = 0$  means the zero-shot set-up. The best performance of each model is made bold.

utilising the full context remains as a challenge for LLMs. Larger models, GPT-3.5 and GPT-4, show more consistent improvement in performance with the increased number of ICL samples. GPT-4 derives the most significant benefits from ICL samples and performs the best across all models.

**The effectiveness of ICL is limited for depression detection.** The performance is in general the best when  $N = 0$ , followed by  $N = 3$ . This suggests that for depression detection, a task to detect more nuanced affective state than emotion from a longer sequence, a single ICL sample for each class could strongly bias the model. This leads to zero F1s where models predict all samples as *Not Depressive*. Including more ICL samples could mitigate this effect, but the performance is further limited by models’ incapability to handle extremely lengthy input. This motivates further research ef-

forts to handle huge context containing nuanced task-related cues when using LLMs.

#### 5.4 Task-specific Fine-tuning

We conduct task-specific fine-tuning experiments with GPT-2, LLaMA-7B, Alpaca-7B, and LLaMA-3-8B using different proportions of training data to explore these models’ capacity for ERC after fine-tuning. Results are summarised in Figure 2. For DAIC-WOZ, fine-tuning would steer models to predict *Not Depressed* (see Table B4) for almost all test samples. This might be due to the small training set where more than 70% of the samples are labelled as *Not Depressed*. This suggests the limitation of language modelling objective, and therefore more carefully curated task-related learning objectives should be considered for depression detection using LLM.

#### Task-specific fine-tuning can effectively and efficiently enhance the ERC performance of LLMs.

For both IEMOCAP and EmoWOZ, we observe an initial significant improvement in performance when fine-tuning with 25% of the training data. Performance remains relatively stable and approaches SOTA levels as the proportion of training data increased to 50% and more for IEMOCAP (4-way) and EmoWOZ. This shows the potential of rapid deployment of LLMs as the emotion recognition frontend in dialogue systems, regardless of dialogue type, label set, or label distribution.

In the case of 5-way classification on IEMOCAP, a performance gap persists between fine-tuned LLMs and the supervised SOTA, even after fine-tuning of LLMs on the complete training set. We hypothesised that this disparity might be attributed to the presence of an additional *Other* class within

the 5-way classification scheme. The class name “Other” lacked essential affective information and consequently failed to fully leverage the language modelling capabilities of LLMs. Therefore, we suggest that employing more semantically meaningful label names could be advantageous in harnessing the potential of LLMs for task-specific fine-tuning.

In the case of GPT-2, fine-tuning does not yield noticeable improvement in ERC. Its performance even deteriorated after fine-tuning with EmoWOZ, as depicted in Figure 2(c) because GPT-2 predominantly predicted *Neutral*, which are excluded from the metric calculation.

## 6 Conclusion

In this study, we explore the performance of LLMs for affect recognition in three distinct types of dialogues: chit-chat dialogues, information-seeking ToDs, and medical consultation dialogues for depression. We conduct benchmark experiments on these datasets using five LLMs: LLaMA-7B, Alpaca-7B, LLaMA-3-8B, GPT-3.5, and GPT-4. We also explore various setups, including zero-shot learning, few-shot in-context learning, and task-specific fine-tuning, all facilitated by specially designed prompts. Additionally, we examine the impact of ASR errors on LLMs’ zero-shot performance.

Our zero-shot experiments underscore that while LLMs have made significant strides in various natural language understanding tasks, they still have some distance to cover in order to match the supervised SOTAs in affect recognition tasks. Adding emotion definitions explaining the eliciting conditions in ToDs to the prompt and fine-tuning LLMs for instruction-following could narrow the performance gap from supervised SOTAs.

Performing zero-shot affect recognition from utterances containing ASR errors shows that LLMs are robust to such errors for emotion recognition but not for depression detection. Therefore, when considering LLMs as a back-end module of a spoken dialogue system, it is crucial to exercise extra caution when processing dialogues laden with highly specific and nuanced affective content.

Our ICL experiments exemplify that larger models would benefit more from an increased number of ICL samples, highlighting the need to explore the optimal combination of the ICL sample size in the

prompt and the model size. For smaller LLMs, effectively utilising lengthy context remains as a challenge.

Through task-specific fine-tuning, we achieve performance levels close to SOTA on IEMOCAP and EmoWOZ, using only 50% of the training data, with LLaMA-7B, Alpaca-7B, and LLaMA-3-8B. This highlights the great potential of fine-tuning LLMs for simpler tasks and integrating them as functional modules into dialogue systems.

Overall, LLMs have opened new avenues for affect recognition in conversations and building affect-aware dialogue systems. Despite the limited performance under zero-shot set-up, their robustness to ASR errors, few-shot ICL capabilities, and ERC capabilities after fine-tuning offer exciting research opportunities for exploring affect recognition in conversations and building human-like conversational agents. We would also like to highlight the challenge and also opportunities towards handling long context and nuanced emotion cues in LLMs.

## 7 Limitations

In our work, although we reduce computation resource of training LLMs by incorporating LoRA, the inference takes  $\tilde{1}$ s for utterance-level emotion recognition on a Nvidia A100 40GB graphics card when there is no ICL sample in the prompt. The inference time increases when the number of ICL samples increases or dialogue-level classification is performed. While LLMs demonstrates superior abilities and potentials, further research efforts are still needed to ensure efficient LLM inference, which is necessary for its application in real-time systems.

With ICL experiments especially on DAIC-WOZ, we observe that the efficacy of long context is limited by the effective spans of the attention mechanisms. While substantial efforts have been invested into increasing the maximum allowed context size of LLMs and improving benchmark performance, the effectiveness of LLMs to make use of full context should not be overlooked.

We only investigate with one dataset from each of three dialogue domains. Although these datasets cover different dialogue settings, objectives, label sets, and classification scopes, there are more affect types and dialogue settings to explore. These datasets also exhibit various degrees of class im-

balance, which selected reference SOTAs utilised data augmentation to address. While GPT-4 has demonstrated good zero-shot learning ability (Section 5.1), addressing data imbalance is out of the scope of this work, and data augmentation with LLMs may come at a cost of potential divergence between synthetic language and real-word data (Li et al., 2023).

## 8 Ethics Statement

Models and datasets were used in accordance with their respective licenses. Data that we used and generated does not contain any information that uniquely identifies individual people. There is a tiny fraction of utterances labelled as “abusive” in EmoWOZ, but they are prompted to models in such a way for the recognition purpose only. Due to the fact that LLMs were pre-trained with a huge amount of data, they may produce inaccurate information about people, places, or facts. This had negligible impact on our evaluation for affect recognition. When performing depression detection and analysis with DAIC-WOZ using GPT-3.5 and GPT-4, models output reminders about seeking professional advice from doctors for more accurate medical diagnosis along with their predictions.

Unlike running models locally, utilising OpenAI’s server-based models would require us to send data to their server. In some cases, it is important to use the application programming interface (API) when for which OpenAI explicitly clarifies that the query data will not be stored or used in model training unless specifically configured.

Although this work focuses on LLMs’ capability in recognising affect in conversations, we do envisage LLMs to be incorporated as an affect recognition frontend in affect-aware dialogue systems. It is therefore important to remember that these models are not perfect and can make errors in their predictions. Subsequently, any actions taken based on these predictions should be executed with an awareness of the possibility of errors. The relatively slow inference speed and the high computational resource requirement also pose a challenge in the usage of LLMs in high-throughput and time-sensitive scenarios.

## 9 Acknowledgement

S. Feng and N. Lubis are supported by funding provided by the Alexander von Humboldt Founda-

tion in the framework of the Sofja Kovalevskaja Award endowed by the Federal Ministry of Education and Research. G. Sun is partly funded by the Department of Engineering, University of Cambridge. Computing resources were provided by Google Cloud.

## References

- AI@Meta. 2024. Llama 3 model card.
- American Psychiatric Association. 2020. What Is Depression? <https://www.psychiatry.org/patients-families/depression/what-is-depression>.
- Edward Beeching, Clémentine Fourier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open llm leaderboard. [https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard).
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N. Chang, Sungbok Lee, and Shrikanth S. Narayanan. 2008. IEMOCAP: interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42(4):335–359.
- Paul Ekman and W V Friesen. 1971. Constants across cultures in the face and emotion. *Journal of personality and social psychology*, 17 2:124–9.
- Shutong Feng, Nurul Lubis, Christian Geishauser, Hsien-chin Lin, Michael Heck, Carel van Niekerk, and Milica Gasic. 2022. EmoWOZ: A large-scale corpus and labelling scheme for emotion recognition in task-oriented dialogue systems. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4096–4113, Marseille, France. European Language Resources Association.
- Shutong Feng, Nurul Lubis, Benjamin Ruppik, Christian Geishauser, Michael Heck, Hsien-chin Lin, Carel van Niekerk, Renato Vukovic, and Milica Gasic. 2023. From chatter to matter: Addressing critical steps of emotion recognition learning in task-oriented dialogue. In *Proceedings of the 24th Meeting of the Special Interest Group on Discourse and Dialogue*, pages 85–103, Prague, Czechia. Association for Computational Linguistics.
- Jonathan Gratch, Ron Artstein, Gale Lucas, Giota Strattou, Stefan Scherer, Angela Nazarian, Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, David Traum, Skip Rizzo, and Louis-Philippe Morency. 2014. The

- distress analysis interview corpus of human and computer interviews. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3123–3128, Reykjavik, Iceland. European Language Resources Association (ELRA).
- James J. Gross. 2002. Emotion regulation: Affective, cognitive, and social consequences. *Psychophysiology*, 39(3):281–291.
- Michael Heck, Nurul Lubis, Benjamin Ruppik, Renato Vukovic, Shutong Feng, Christian Geishauer, Hsien-chin Lin, Carel van Niekerk, and Milica Gasic. 2023. ChatGPT for zero-shot dialogue state tracking: A solution or an opportunity? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 936–950, Toronto, Canada. 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*.
- Jen-tse Huang, Man Lam, Eric Li, Shujie Ren, Wenxuan Wang, Wenxiang Jiao, Zhaopeng Tu, and Michael Lyu. 2023. Emotionally numb or empathetic? evaluating how llms feel using emotionbench.
- Dan Jurafsky and James H. Martin. 2009. *Speech and language processing : an introduction to natural language processing, computational linguistics, and speech recognition*. Pearson Prentice Hall, Upper Saddle River, N.J.
- Kurt Kroenke, Tara W Strine, Robert L Spitzer, Janet B W Williams, Joyce T Berry, and Ali H Mokdad. 2008. The PHQ-8 as a measure of current depression in the general population. *J Affect Disord*, 114(1-3):163–173.
- Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, and Ming Yin. 2023. Synthetic data generation with large language models for text classification: Potential and limitations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10443–10461, Singapore. Association for Computational Linguistics.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>.
- John D Mayer, David R Caruso, and Peter Salovey. 1999. Emotional intelligence meets traditional standards for an intelligence. *Intelligence*, 27(4):267–298.
- OpenAI. 2022. Introducing ChatGPT. <https://openai.com/blog/chatgpt>.
- OpenAI. 2023. Gpt-4 technical report.
- 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. In *Proceedings of NeurIPS 2022*.
- Rosalind W. Picard. 1997. *Affective Computing*. MIT Press, Cambridge, MA.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Vijay Ravi, Jinhan Wang, Jonathan Flint, and Abeer Alwan. 2022. A step towards preserving speakers’ identity while detecting depression via speaker disentanglement. *Interspeech*, 2022:3338–3342.
- James A. Russell. 1980. A circumplex model of affect. *Journal of Personality and Social Psychology*, 39(6):1161–1178.
- Guangzhi Sun, Shutong Feng, Dongcheng Jiang, Chao Zhang, Milica Gašić, and Philip C. Woodland. 2023. Speech-based slot filling using large language models.
- 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*.
- 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](https://github.com/tatsu-lab/stanford_alpaca).
- 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. 2023a. Llama: Open and efficient foundation language models.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutvi 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. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#).
- Michel Valstar, Jonathan Gratch, Björn Schuller, Fabien Ringeval, Denis Lalanne, Mercedes Torres Torres, Stefan Scherer, Giota Stratou, Roddy Cowie, and Maja Pantic. 2016. [Avec 2016: Depression, mood, and emotion recognition workshop and challenge](#). In *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*, AVEC '16, page 3–10, New York, NY, USA. Association for Computing Machinery.
- Xuena Wang, Xueting Li, Zi Yin, Yue Wu, and Liu Jia. 2023. [Emotional intelligence of large language models](#).
- 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](#). *Transactions on Machine Learning Research*. Survey Certification.
- James R. Williamson, Elizabeth Godoy, Miriam Cha, Adrienne Schwarzentruher, Pooya Khorrami, Youngjune Gwon, Hsiang-Tsung Kung, Charlie Dagli, and Thomas F. Quatieri. 2016. [Detecting depression using vocal, facial and semantic communication cues](#). In *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*, AVEC '16, page 11–18, New York, NY, USA. Association for Computing Machinery.
- Wen Wu, Chao Zhang, and Philip C. Woodland. 2020. [Emotion recognition by fusing time synchronous and time asynchronous representations](#). *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6269–6273.
- Wen Wu, Chao Zhang, and Philip C. Woodland. 2023. [Self-supervised representations in speech-based depression detection](#). In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE.
- Le Yang, Dongmei Jiang, Lang He, Ercheng Pei, Meshia Cédric Oveneke, and Hichem Sahli. 2016. [Decision tree based depression classification from audio video and language information](#). In *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*, AVEC '16, page 89–96, New York, NY, USA. Association for Computing Machinery.
- Zhihong Zeng, Maja Pantic, Glenn I. Roisman, and Thomas S. Huang. 2009. [A survey of affect recognition methods: Audio, visual, and spontaneous expressions](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(1):39–58.
- Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023. [Sentiment analysis in the era of large language models: A reality check](#).
- Justin Zhao, Timothy Wang, Wael Abid, Geoffrey Angus, Arnav Garg, Jeffery Kinnison, Alex Sherstinsky, Piero Molino, Travis Addair, and Devvret Rishi. 2024. [Lora land: 310 fine-tuned llms that rival gpt-4, a technical report](#).
- Weixiang Zhao, Yanyan Zhao, Xin Lu, Shilong Wang, Yanpeng Tong, and Bing Qin. 2023. [Is chatgpt equipped with emotional dialogue capabilities?](#)

## A Detailed Training Configurations

### A.1 Task-Specific Fine-tuning

For all model fine-tuning, the learning rate was  $3e-5$ . The batch size was 2 with a gradient accumulation step of 4. We used a cosinusoidal learning rate scheduler without warming up. We applied a weight decay of 0.01 on all model parameters except for the biases and layer normalisation weights. For LLaMA-7B, Alpaca-7B, and LLaMA-3-8B, we stored model parameters in IEEE 754 half-precision float point format. For GPT-2, we stored the model parameters in standard single-precision floating-point format and did not apply LoRA during the fine-tuning. We followed the default LoRA configuration provided in Huggingface PEFT library (Mangrulkar et al., 2022). We used the model perplexity on the development set as the early-stopping criterion. For EmoWOZ, we used the official development set. For IEMOCAP, when we performed the leave-one-session-out training, 10% of the training data were randomly sampled as the development set. We applied stratified sampling based on the emotion labels. All open-source models were trained on a single Nvidia A100 40GB Graphics Card.

### A.2 ASR Simulation for EmoWOZ

We fine-tuned a LLaMA-7B model using LoRA following configurations specified in Section 3.2 and A.1 for one epoch on all IEMOCAP utterances. The source was each of the IEMOCAP utterance transcription and the target was the corresponding OpenAI Whisper-medium hypothesis. We utilised a prompt template that formatted the source and target in natural language would best exploit the language modelling capability of the model:

After adding automatic speech recognition errors, [SOURCE] becomes [TARGET]

## B Detailed Experimental Results

Model	N	P	M	Neutral	Happy	Angry	Sad	WA	UA
GPT-2	0	0%	T	0.7 (60.0/0.4)	32.3 (43.6/25.6)	35.3 (22.0/90.6)	0.5 (30.0/0.3)	25.8	29.2
GPT-2	1	0%	T	10.9 (43.4/6.2)	9.1 (62.0/4.9)	29.0 (21.8/43.6)	33.3 (22.8/61.9)	24.2	29.2
GPT-2	0	25%	T	47.2 (30.9/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	30.9	25.0
GPT-2	0	50%	T	47.2 (30.9/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	30.9	25.0
GPT-2	0	75%	T	47.2 (30.9/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	30.9	25.0
GPT-2	0	100%	T	47.2 (30.9/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	30.9	25.0
LLaMA-7B	0	0%	T	48.6 (37.5/69.3)	21.8 (82.3/12.5)	53.3 (40.8/76.8)	6.9 (78.0/3.6)	41.1	40.5
LLaMA-7B	0	0%	A	50.3 (37.0/78.8)	14.8 (79.3/8.2)	54.7 (44.9/70.0)	0.6 (100.0/0.3)	40.8	39.3
LLaMA-7B	1	0%	T	55.3 (42.9/77.5)	56.2 (77.3/44.2)	62.0 (55.0/71.0)	11.1 (73.0/6.0)	52.3	49.7
LLaMA-7B	3	0%	T	54.2 (39.4/86.7)	1.2 (90.9/0.6)	44.1 (87.6/29.5)	44.5 (39.5/50.8)	42.8	41.9
LLaMA-7B	0	25%	T	65.1 (64.6/65.6)	77.0 (80.8/73.5)	73.5 (72.6/74.4)	74.2 (71.3/77.3)	72.0	72.7
LLaMA-7B	0	50%	T	69.1 (69.9/68.3)	80.7 (80.6/80.7)	77.3 (78.7/76.0)	78.1 (75.6/80.8)	76.0	76.4
LLaMA-7B	0	75%	T	70.7 (67.3/74.5)	82.2 (84.6/80.0)	76.3 (80.5/72.4)	78.1 (77.5/78.6)	76.5	76.4
LLaMA-7B	0	100%	T	69.7 (66.2/73.5)	82.0 (82.4/81.7)	79.0 (81.0/77.1)	75.8 (80.6/71.6)	76.3	76.0
Alpaca-7B	0	0%	T	34.4 (49.6/26.3)	62.8 (79.3/52.0)	50.2 (34.2/94.6)	44.5 (70.3/32.6)	48.8	51.4
Alpaca-7B	0	0%	A	37.0 (52.5/28.6)	60.7 (75.3/50.8)	48.5 (32.8/93.3)	38.8 (77.6/25.8)	47.5	49.6
Alpaca-7B	1	0%	T	53.8 (49.4/59.0)	59.5 (83.1/46.3)	57.9 (43.5/86.9)	37.0 (72.3/24.8)	54.1	54.3
Alpaca-7B	3	0%	T	55.4 (43.9/74.8)	28.0 (90.6/16.6)	64.5 (60.8/68.6)	54.9 (55.1/54.7)	52.4	53.7
Alpaca-7B	0	25%	T	65.1 (67.3/63.0)	77.0 (77.8/76.2)	74.6 (72.1/77.2)	70.8 (68.8/73.0)	71.7	72.4
Alpaca-7B	0	50%	T	69.1 (69.8/68.3)	80.5 (78.7/82.4)	78.3 (79.7/77.1)	75.7 (75.9/75.6)	75.6	75.8
Alpaca-7B	0	75%	T	70.5 (66.3/75.4)	80.6 (85.2/76.4)	76.3 (78.7/74.2)	74.2 (74.2/74.2)	75.2	75.0
Alpaca-7B	0	100%	T	69.3 (69.1/69.5)	81.0 (82.3/79.8)	78.8 (79.2/78.4)	76.5 (74.7/78.3)	76.0	76.5
LLaMA-3-8B	0	0%	T	3.4 (55.6/1.8)	55.9 (42.7/81.0)	51.0 (38.7/75.0)	19.9 (55.7/12.1)	41.8	42.5
LLaMA-3-8B	0	0%	A	2.1 (56.2/1.1)	51.2 (36.1/88.3)	52.3 (47.4/58.5)	14.3 (65.4/8.0)	39.7	39.0
LLaMA-3-8B	1	0%	T	52.7 (50.1/55.6)	67.1 (75.9/60.1)	60.4 (46.7/85.4)	37.9 (82.9/24.5)	56.8	56.4
LLaMA-3-8B	3	0%	T	35.8 (62.0/25.2)	67.0 (79.5/57.9)	63.4 (53.3/78.2)	60.0 (46.1/86.2)	57.4	61.9
LLaMA-3-8B	0	25%	T	68.2 (66.8/69.6)	78.5 (76.6/80.5)	74.4 (76.0/72.9)	72.4 (76.4/68.7)	73.3	72.9
LLaMA-3-8B	0	50%	T	69.7 (71.0/68.5)	81.4 (79.0/84.0)	77.6 (82.3/73.3)	77.4 (74.8/80.2)	76.3	76.5
LLaMA-3-8B	0	75%	T	71.8 (73.8/70.0)	83.2 (81.7/84.7)	80.6 (82.2/79.1)	80.2 (77.8/82.7)	78.7	79.1
LLaMA-3-8B	0	100%	T	73.2 (74.0/72.5)	84.3 (83.0/85.8)	81.6 (83.1/80.2)	81.0 (80.4/81.6)	79.8	80.0
GPT-3.5	0	0%	T	51.6 (35.1/97.3)	28.5 (90.5/16.9)	31.7 (79.9/19.8)	27.1 (81.5/16.2)	42.2	37.6
GPT-3.5	0	0%	A	51.3 (34.9/96.5)	31.7 (88.8/19.3)	33.0 (83.2/20.6)	23.2 (80.8/13.6)	42.3	37.5
GPT-3.5	1	0%	T	57.7 (42.7/88.9)	59.4 (84.8/45.7)	56.6 (76.1/45.1)	45.9 (80.0/32.2)	56.3	53.0
GPT-3.5	3	0%	T	60.1 (48.9/78.2)	66.0 (80.0/56.1)	63.8 (76.4/54.8)	59.4 (67.1/53.3)	62.1	60.6
GPT-4	0	0%	T	51.7 (35.0/99.3)	28.9 (97.5/16.9)	29.4 (95.5/17.4)	28.2 (91.4/16.7)	42.4	37.6
GPT-4	0	0%	A	51.5 (34.8/98.9)	27.6 (94.3/16.2)	30.4 (95.7/18.0)	25.9 (89.6/15.1)	41.9	37.1
GPT-4	1	0%	T	62.3 (48.1/88.2)	59.7 (83.4/46.5)	70.4 (81.6/61.9)	60.9 (81.1/48.7)	62.9	61.3
GPT-4	3	0%	T	61.6 (49.0/83.1)	60.7 (84.9/47.2)	67.6 (85.2/56.1)	68.4 (71.6/65.4)	63.8	63.0

Table B1: F1(precision/recall), UA, and WA of LLMs on IEMOCAP under the 4-Way classification set-up. In table headers, “N” stands for the number of ICL samples in the prompt; “P” stands for the proportion of training data used for fine-tuning; “M” stands for the modality of input, either transcription (T) or ASR hypothesis (A).

Model	N	P	M	Neutral	Happy	Angry	Sad	Other	WA	UA
GPT-2	0	0%	T	0.1 (7.7/0.1)	23.1 (35.2/17.2)	27.5 (16.5/81.1)	0.5 (42.9/0.3)	15.2 (19.3/12.5)	19.0	22.3
GPT-2	1	0%	T	13.3 (31.6/8.4)	26.6 (35.4/21.3)	19.2 (16.3/23.5)	27.4 (17.0/70.5)	0.0 (0.0/0.0)	20.1	24.7
GPT-2	0	25%	T	37.0 (22.7/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	22.7	20.0
GPT-2	0	50%	T	37.0 (22.7/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	22.7	20.0
GPT-2	0	75%	T	37.0 (22.7/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	22.7	20.0
GPT-2	0	100%	T	37.0 (22.7/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	22.7	20.0
LLaMA-7B	0	0%	T	38.7 (29.8/55.2)	21.5 (81.8/12.3)	37.8 (30.3/50.0)	6.3 (65.5/3.3)	44.5 (42.1/47.3)	35.6	33.6
LLaMA-7B	0	0%	A	37.0 (30.9/46.3)	13.4 (80.5/7.3)	17.3 (40.5/11.0)	0.0 (0.0/0.0)	48.0 (34.6/78.2)	34.5	28.6
LLaMA-7B	1	0%	T	2.2 (59.4/1.1)	0.1 (100.0/0.1)	15.1 (60.9/8.6)	0.0 (0.0/0.0)	41.6 (26.5/97.2)	27.3	21.4
LLaMA-7B	3	0%	T	11.1 (22.4/7.4)	0.2 (16.7/0.1)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	41.2 (26.5/92.3)	26.2	20.0
LLaMA-7B	0	25%	T	51.2 (48.1/54.9)	74.3 (72.4/76.3)	49.3 (54.6/44.9)	62.5 (63.3/61.6)	49.9 (51.3/48.6)	57.4	57.3
LLaMA-7B	0	50%	T	57.9 (54.5/61.7)	79.0 (77.1/80.9)	54.7 (57.7/52.0)	69.3 (73.8/65.3)	55.6 (56.7/54.6)	63.1	62.9
LLaMA-7B	0	75%	T	60.8 (58.2/63.7)	82.3 (82.2/82.5)	56.4 (60.9/52.5)	72.9 (68.9/77.4)	57.3 (59.6/55.2)	65.9	66.3
LLaMA-7B	0	100%	T	53.9 (57.5/50.6)	80.9 (77.6/84.5)	55.8 (63.6/49.8)	72.1 (69.8/74.5)	57.4 (54.5/60.7)	64.0	64.0
Alpaca-7B	0	0%	T	18.1 (42.5/11.5)	52.6 (78.3/39.6)	29.2 (34.9/25.0)	29.2 (69.4/18.5)	48.4 (33.6/86.6)	40.5	36.2
Alpaca-7B	0	0%	A	15.8 (41.3/9.8)	49.5 (74.3/37.1)	23.4 (31.5/18.6)	21.8 (72.8/12.8)	48.2 (32.9/90.0)	38.7	33.6
Alpaca-7B	1	0%	T	0.2 (28.6/0.1)	0.5 (100.0/0.2)	5.1 (74.4/2.6)	0.0 (0.0/0.0)	42.0 (26.6/99.5)	26.9	20.5
Alpaca-7B	3	0%	T	9.4 (17.1/6.4)	3.7 (24.3/2.0)	3.1 (18.1/1.7)	5.7 (10.4/3.9)	39.6 (26.1/81.5)	24.4	19.1
Alpaca-7B	0	25%	T	48.9 (53.2/45.3)	73.8 (68.0/80.8)	51.2 (56.4/46.9)	63.4 (60.5/66.5)	52.1 (51.8/52.4)	58.2	58.4
Alpaca-7B	0	50%	T	56.4 (54.0/59.1)	78.7 (76.8/80.6)	55.5 (58.9/52.4)	68.8 (74.6/63.8)	57.4 (57.0/57.9)	63.2	62.8
Alpaca-7B	0	75%	T	57.9 (59.5/56.4)	81.0 (77.5/84.9)	58.9 (62.5/55.7)	70.2 (64.8/76.7)	56.9 (59.2/54.9)	65.0	65.7
Alpaca-7B	0	100%	T	54.7 (55.6/53.8)	81.1 (77.7/84.8)	58.0 (60.3/55.9)	69.8 (65.3/75.0)	55.1 (57.8/52.6)	63.6	64.4
LLaMA-3-8B	0	0%	T	1.6 (46.7/0.8)	44.5 (30.1/85.3)	38.5 (27.3/65.1)	12.4 (40.9/7.3)	0.4 (10.3/0.2)	29.4	31.7
LLaMA-3-8B	0	0%	A	0.9 (44.4/0.5)	41.0 (26.5/90.7)	39.5 (32.4/50.6)	9.2 (47.8/5.1)	1.6 (23.6/0.8)	28.2	29.5
LLaMA-3-8B	1	0%	T	44.1 (43.5/44.8)	68.5 (66.7/70.4)	39.8 (25.8/87.5)	22.8 (85.1/13.2)	2.7 (21.6/1.4)	40.5	43.5
LLaMA-3-8B	3	0%	T	23.1 (63.9/14.1)	64.2 (70.8/58.7)	44.7 (40.1/50.5)	37.6 (23.5/94.7)	0.7 (21.2/0.3)	37.1	43.7
LLaMA-3-8B	0	25%	T	54.6 (54.9/54.4)	76.6 (74.4/78.9)	46.9 (62.4/37.6)	67.4 (65.9/69.1)	54.4 (50.8/58.5)	60.5	59.7
LLaMA-3-8B	0	50%	T	58.4 (60.3/56.7)	80.4 (78.3/82.7)	54.4 (65.5/46.5)	72.8 (70.5/75.3)	60.7 (57.2/64.6)	65.6	65.1
LLaMA-3-8B	0	75%	T	60.2 (61.5/59.0)	81.9 (78.6/85.4)	56.9 (68.2/48.9)	74.6 (73.8/75.4)	61.8 (58.7/65.1)	67.2	66.8
LLaMA-3-8B	0	100%	T	63.9 (66.1/61.9)	83.2 (81.4/85.1)	59.5 (66.0/54.1)	75.8 (76.7/74.9)	63.4 (59.8/67.5)	69.5	68.7
GPT-3.5	0	0%	T	43.6 (28.7/91.2)	29.2 (87.0/17.5)	29.4 (63.4/19.1)	26.3 (72.2/16.1)	39.2 (52.5/31.3)	37.9	35.1
GPT-3.5	0	0%	A	43.7 (28.8/90.9)	32.8 (87.3/20.2)	29.2 (61.1/19.2)	24.6 (75.7/14.7)	38.4 (51.1/30.7)	38.1	35.1
GPT-3.5	1	0%	T	45.9 (36.3/62.5)	63.3 (74.8/54.9)	49.8 (46.7/53.4)	48.8 (65.4/38.9)	38.0 (44.7/33.1)	48.3	48.6
GPT-3.5	3	0%	T	47.4 (43.1/52.7)	67.6 (69.8/65.6)	49.5 (40.4/63.9)	54.1 (45.4/66.8)	18.3 (41.4/11.7)	48.3	52.1
GPT-4	0	0%	T	43.1 (28.0/93.7)	28.1 (94.1/16.5)	27.4 (82.6/16.4)	29.4 (85.7/17.7)	37.8 (53.8/29.1)	37.5	34.7
GPT-4	0	0%	A	42.9 (27.7/95.3)	27.6 (94.3/16.1)	30.7 (79.2/19.0)	27.2 (79.5/16.4)	31.8 (51.5/23.0)	36.4	34.0
GPT-4	1	0%	T	51.1 (37.9/78.2)	58.9 (80.4/46.5)	55.3 (49.9/61.9)	54.5 (62.0/48.7)	27.0 (45.6/19.1)	49.0	50.9
GPT-4	3	0%	T	49.6 (39.8/65.8)	59.9 (81.7/47.2)	54.7 (53.5/56.1)	58.3 (52.5/65.4)	30.9 (40.5/24.9)	49.4	51.9

Table B2: F1(precision/recall), UA, and WA of LLMs on IEMOCAP under the 5-Way classification set-up. In table headers, “N” stands for the number of ICL samples in the prompt; “P” stands for the proportion of training data used for fine-tuning; “M” stands for the modality of input, either transcription (T) or ASR hypothesis (A).

Model	N	P	M	Neutral	Fearful	Dissatisfied	Apogetic	Abusive	Excited	Satisfied	MF1	WF1
GPT-2	0	0%	T	0.1 (100.0/0.0)	0.0 (0.0/0.0)	9.3 (5.6/27.8)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	2.8 (1.4/64.8)	31.4 (35.7/28.1)	7.3	24.0
GPT-2	1	0%	T	81.2 (69.8/97.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	23.3 (14.8/54.8)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	3.9	0.6
GPT-2	0	25%	T	82.4 (70.1/99.8)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	69.9 (71.4/68.5)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	11.7	1.9
GPT-2	0	50%	T	82.3 (70.0/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	69.6 (95.2/54.8)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	11.6	1.9
GPT-2	0	75%	T	82.4 (70.3/99.5)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	58.9 (47.9/76.7)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	9.8	1.6
GPT-2	0	100%	T	82.3 (70.0/99.8)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	68.1 (74.2/63.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	11.4	1.9
LLaMA-7B	0	0%	T	82.1 (69.7/100.0)	0.0 (0.0/0.0)	0.3 (33.3/0.2)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	6.3 (75.0/3.3)	0.0 (0.0/0.0)	1.1	0.3
LLaMA-7B	0	0%	A	82.1 (69.7/100.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0	0.0
LLaMA-7B	1	0%	T	83.0 (78.1/88.5)	26.1 (60.0/16.7)	2.6 (47.1/1.3)	0.0 (0.0/0.0)	57.9 (52.4/64.7)	16.0 (9.2/58.2)	59.0 (74.1/49.0)	26.9	42.6
LLaMA-7B	3	0%	T	27.9 (81.2/16.9)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	0.0 (0.0/0.0)	39.2 (24.4/99.2)	6.5	27.2
LLaMA-7B	0	25%	T	93.9 (91.5/96.4)	26.1 (60.0/16.7)	55.2 (81.6/41.7)	72.3 (93.5/58.9)	11.1 (100.0/5.9)	43.6 (69.0/31.9)	90.9 (89.1/92.7)	49.9	79.5
LLaMA-7B	0	50%	T	94.4 (93.1/95.8)	41.7 (83.3/27.8)	68.6 (80.0/60.1)	75.8 (92.2/64.4)	64.0 (100.0/47.1)	51.4 (69.8/40.7)	91.1 (89.8/92.5)	65.4	83.6
LLaMA-7B	0	75%	T	93.8 (93.0/94.5)	35.3 (37.5/33.3)	61.7 (84.9/48.5)	57.4 (41.5/93.2)	69.2 (100.0/52.9)	50.6 (54.4/47.3)	90.8 (88.7/93.1)	60.8	81.3
LLaMA-7B	0	100%	T	94.2 (93.3/95.2)	43.8 (50.0/38.9)	68.1 (78.3/60.3)	75.4 (93.9/63.0)	69.2 (100.0/52.9)	51.3 (63.9/42.9)	90.7 (88.8/92.6)	66.4	83.2
Alpaca-7B	0	0%	T	65.4 (85.3/53.1)	1.9 (1.1/11.1)	24.8 (28.5/22.0)	46.0 (85.2/31.5)	0.0 (0.0/0.0)	18.1 (13.8/26.4)	53.3 (38.7/85.7)	24.0	44.6
Alpaca-7B	0	0%	A	65.5 (83.0/54.1)	2.0 (1.1/11.1)	22.9 (26.2/20.4)	52.9 (93.1/37.0)	0.0 (0.0/0.0)	17.3 (15.4/19.8)	50.7 (37.2/80.0)	24.3	42.6
Alpaca-7B	1	0%	T	75.4 (81.4/70.2)	3.4 (1.9/22.2)	13.8 (23.8/9.8)	26.2 (100.0/15.1)	30.0 (100.0/17.6)	8.9 (4.8/65.9)	67.4 (69.4/65.6)	25.0	51.2
Alpaca-7B	3	0%	T	65.4 (85.3/53.1)	1.9 (1.1/11.1)	24.8 (28.5/22.0)	46.0 (85.2/31.5)	0.0 (0.0/0.0)	18.1 (13.8/26.4)	53.3 (38.7/85.7)	24.0	44.6
Alpaca-7B	0	25%	T	93.3 (91.9/94.7)	17.4 (40.0/11.1)	53.5 (72.5/42.4)	74.8 (92.0/63.0)	0.0 (0.0/0.0)	45.4 (64.0/35.2)	90.2 (86.6/94.2)	46.9	78.7
Alpaca-7B	0	50%	T	94.4 (93.1/95.8)	43.5 (100.0/27.8)	68.1 (79.9/59.3)	74.2 (83.1/67.1)	64.0 (100.0/47.1)	46.0 (66.7/35.2)	91.0 (89.6/92.6)	64.5	83.2
Alpaca-7B	0	75%	T	93.6 (91.0/96.4)	35.7 (50.0/27.8)	45.6 (90.6/30.5)	75.7 (79.1/72.6)	38.1 (100.0/23.5)	50.7 (67.3/40.7)	91.1 (88.7/93.6)	56.1	78.0
Alpaca-7B	0	100%	T	94.0 (92.0/96.1)	10.5 (100.0/5.6)	62.2 (76.4/52.5)	71.8 (73.9/69.9)	0.0 (0.0/0.0)	39.3 (77.4/26.4)	90.7 (90.2/91.1)	45.8	80.7
LLaMA-3-8B	0	0%	T	44.2 (79.5/30.6)	1.3 (0.6/55.6)	1.0 (13.0/0.5)	24.0 (14.6/67.1)	26.8 (15.8/88.2)	5.7 (2.9/81.3)	59.3 (59.6/59.1)	19.7	42.4
LLaMA-3-8B	0	0%	A	47.0 (80.7/33.1)	1.2 (0.6/50.0)	1.0 (12.5/0.5)	24.1 (14.7/67.1)	27.5 (16.3/88.2)	5.7 (3.0/78.0)	59.2 (59.3/59.1)	19.8	42.3
LLaMA-3-8B	1	0%	T	83.5 (76.3/92.3)	4.5 (3.8/5.6)	2.4 (16.0/1.3)	39.4 (30.8/54.8)	9.0 (4.7/100.0)	35.4 (46.4/28.6)	50.5 (87.4/35.5)	23.5	38.0
LLaMA-3-8B	3	0%	T	55.5 (85.1/41.1)	0.0 (0.0/0.0)	0.3 (20.0/0.2)	35.3 (62.1/24.7)	2.3 (1.2/100.0)	24.6 (60.9/15.4)	54.8 (39.2/90.9)	19.5	39.9
LLaMA-3-8B	0	25%	T	93.6 (90.7/96.7)	27.3 (75.0/16.7)	52.4 (87.6/37.4)	74.6 (97.8/60.3)	74.1 (100.0/58.8)	49.0 (64.3/39.6)	90.4 (89.5/91.3)	61.3	79.2
LLaMA-3-8B	0	50%	T	93.8 (91.1/96.7)	26.1 (60.0/16.7)	56.5 (86.1/42.1)	76.4 (94.0/64.4)	74.1 (100.0/58.8)	49.3 (62.7/40.7)	90.5 (90.1/91.0)	62.2	80.3
LLaMA-3-8B	0	75%	T	94.3 (92.4/96.3)	38.5 (62.5/27.8)	64.0 (85.2/51.3)	75.2 (83.3/68.5)	64.0 (100.0/47.1)	48.6 (66.0/38.5)	90.8 (89.5/92.2)	63.5	82.2
LLaMA-3-8B	0	100%	T	94.5 (92.4/96.7)	50.0 (100.0/33.3)	66.2 (85.2/54.1)	74.4 (85.7/65.8)	78.6 (100.0/64.7)	52.1 (69.1/41.8)	90.8 (90.5/91.0)	68.7	82.9
GPT-3.5	0	0%	T	82.8 (76.9/89.8)	20.7 (27.3/16.7)	8.2 (35.0/4.6)	61.9 (87.5/47.9)	61.5 (88.9/47.1)	31.6 (27.4/37.4)	50.0 (58.9/43.5)	39.0	40.0
GPT-3.5	0	0%	A	82.9 (76.8/90.0)	28.6 (100.0/16.7)	8.2 (35.4/4.6)	61.9 (87.5/47.9)	61.5 (88.9/47.1)	31.5 (27.2/37.4)	49.7 (58.8/43.1)	40.2	39.8
GPT-3.5	1	0%	T	66.0 (82.3/55.1)	36.4 (100.0/22.2)	13.5 (33.1/8.4)	7.9 (100.0/4.1)	75.7 (70.0/82.4)	10.7 (5.9/61.5)	56.0 (42.6/81.6)	33.3	43.2
GPT-3.5	3	0%	T	57.9 (82.1/44.7)	10.0 (50.0/5.6)	16.6 (34.4/10.9)	36.0 (100.0/21.9)	69.0 (83.3/58.8)	6.4 (3.4/71.4)	59.3 (46.5/81.8)	32.9	46.7
GPT-4	0	0%	T	88.3 (86.0/90.8)	50.0 (100.0/33.3)	16.4 (47.2/9.9)	52.5 (37.4/87.7)	74.1 (100.0/58.8)	42.2 (36.2/50.5)	79.0 (78.5/79.6)	52.4	62.3
GPT-4	0	0%	A	88.3 (82.7/94.6)	41.7 (83.3/27.8)	47.9 (70.1/36.4)	47.9 (33.2/86.3)	75.9 (91.7/64.7)	39.8 (31.6/53.8)	66.6 (89.4/53.1)	53.3	60.8
GPT-4	1	0%	T	78.8 (93.8/68.0)	41.7 (83.3/27.8)	52.5 (46.7/60.1)	42.7 (28.2/87.7)	83.3 (78.9/88.2)	14.6 (8.0/80.2)	71.8 (63.8/82.1)	51.1	64.4
GPT-4	3	0%	T	83.2 (91.9/76.0)	26.1 (60.0/16.7)	51.1 (48.0/54.6)	55.0 (42.0/79.5)	77.8 (73.7/82.4)	28.0 (20.3/45.1)	74.2 (63.0/90.4)	52.0	66.5

Table B3: F1(precision/recall), MF1 and WF1 of LLMs on EmoWOZ. In table headers, “N” stands for the number of ICL samples in the prompt; “P” stands for the proportion of training data used for fine-tuning; “M” stands for the modality of input, either transcription (T) or ASR hypothesis (A).

Model	N	P	M	Development Set		Test Set	
				Depressed	Not Depressed	Depressed	Not Depressed
GPT-2	0	0%	T	0.0 (0.0/0.0)	82.5 (70.2/100.0)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
LLaMA-7B	0	0%	T	47.5 (31.1/100.0)	11.4 (100.0/6.1)	52.2 (35.3/100.0)	8.3 (100.0/4.3)
LLaMA-7B	0	0%	A	45.9 (29.8/100.0)	0.0 (0.0/0.0)	51.1 (34.3/100.0)	0.0 (0.0/0.0)
LLaMA-7B	1	0%	T	0.0 (0.0/0.0)	82.5 (70.2/100.0)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
LLaMA-7B	3	0%	T	42.1 (27.9/85.7)	10.8 (50.0/6.1)	48.9 (33.3/91.7)	8.0 (50.0/4.3)
LLaMA-7B	0	25%	T	0.0 (0.0/0.0)	81.0 (69.6/97.0)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
LLaMA-7B	0	50%	T	0.0 (0.0/0.0)	81.0 (69.6/97.0)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
LLaMA-7B	0	75%	T	0.0 (0.0/0.0)	79.5 (68.9/93.9)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
LLaMA-7B	0	100%	T	0.0 (0.0/0.0)	76.3 (67.4/87.9)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
Alpaca-7B	0	0%	T	47.5 (31.1/100.0)	11.4 (100.0/6.1)	53.3 (36.4/100.0)	16.0 (100.0/8.7)
Alpaca-7B	0	0%	A	45.9 (29.8/100.0)	0.0 (0.0/0.0)	53.3 (36.4/100.0)	16.0 (100.0/8.7)
Alpaca-7B	1	0%	T	0.0 (0.0/0.0)	82.5 (70.2/100.0)	15.4 (100.0/8.3)	80.7 (67.6/100.0)
Alpaca-7B	3	0%	T	45.9 (29.8/100.0)	0.0 (0.0/0.0)	51.1 (34.3/100.0)	0.0 (0.0/0.0)
Alpaca-7B	0	25%	T	12.5 (50.0/7.1)	82.1 (71.1/97.0)	0.0 (0.0/0.0)	77.2 (64.7/95.7)
Alpaca-7B	0	50%	T	11.8 (33.3/7.1)	80.5 (70.5/93.9)	0.0 (0.0/0.0)	77.2 (64.7/95.7)
Alpaca-7B	0	75%	T	10.5 (20.0/7.1)	77.3 (69.0/87.9)	13.3 (33.3/8.3)	76.4 (65.6/91.3)
Alpaca-7B	0	100%	T	18.2 (25.0/14.3)	75.0 (69.2/81.8)	0.0 (0.0/0.0)	75.0 (63.6/91.3)
LLaMA-3-8B	0	0%	T	47.1 (32.4/85.7)	37.2 (80.0/24.2)	43.2 (32.0/66.7)	36.4 (60.0/26.1)
LLaMA-3-8B	0	0%	A	46.4 (31.0/92.9)	21.1 (80.0/12.1)	42.9 (30.0/75.0)	14.3 (40.0/8.7)
LLaMA-3-8B	1	0%	T	0.0 (0.0/0.0)	82.5 (70.2/100.0)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
LLaMA-3-8B	3	0%	T	0.0 (0.0/0.0)	82.5 (70.2/100.0)	0.0 (0.0/0.0)	79.3 (65.7/100.0)
LLaMA-3-8B	0	25%	T	27.0 (21.7/35.7)	52.6 (62.5/45.5)	29.6 (26.7/33.3)	55.8 (60.0/52.2)
LLaMA-3-8B	0	50%	T	0.0 (0.0/0.0)	82.5 (70.2/100.0)	0.0 (0.0/0.0)	77.2 (64.7/95.7)
LLaMA-3-8B	0	75%	T	0.0 (0.0/0.0)	82.5 (70.2/100.0)	47.1 (36.4/66.7)	50.0 (69.2/39.1)
LLaMA-3-8B	0	100%	T	20.0 (33.3/14.3)	78.4 (70.7/87.9)	47.1 (36.4/66.7)	50.0 (69.2/39.1)
GPT-3.5	0	0%	T	54.5 (60.0/50.0)	79.2 (76.0/82.6)	64.3 (64.3/64.3)	84.8 (84.8/84.8)
GPT-3.5	0	0%	A	37.5 (75.0/25.0)	81.5 (71.0/95.7)	56.0 (63.6/50.0)	84.1 (80.6/87.9)
GPT-3.5	1	0%	T	13.3 (33.3/8.3)	76.4 (65.6/91.3)	40.0 (45.5/35.7)	78.3 (75.0/81.8)
GPT-3.5	3	0%	T	37.5 (75.0/25.0)	81.5 (71.0/95.7)	56.0 (63.6/50.0)	84.1 (80.6/87.9)
GPT-4	0	0%	T	63.6 (70.0/58.3)	83.3 (80.0/87.0)	59.3 (61.5/57.1)	83.6 (82.4/84.8)
GPT-4	0	0%	A	44.4 (66.7/33.3)	80.8 (72.4/91.3)	41.7 (50.0/35.7)	80.0 (75.7/84.8)
GPT-4	1	0%	T	80.0 (76.9/83.3)	88.9 (90.9/87.0)	55.6 (45.5/71.4)	72.4 (84.0/63.6)
GPT-4	3	0%	T	74.1 (66.7/83.3)	83.7 (90.0/78.3)	58.5 (44.4/85.7)	58.5 (44.4/85.7)

Table B4: F1(precision/recall) of LLMs on DAIC-WOZ. In table headers, “N” stands for the number of ICL samples in the prompt; “P” stands for the proportion of training data used for fine-tuning; “M” stands for the modality of input, either transcription (T) or ASR hypothesis (A).

# Sentiment-Aware Dialogue Flow Discovery for Interpreting Communication Trends

Patrícia Ferreira<sup>1,2</sup>

Isabel Carvalho<sup>1,2</sup>

Ana Alves<sup>1,3</sup>

Catarina Silva<sup>1,2</sup>

Hugo Gonçalo Oliveira<sup>1,2</sup>

<sup>1</sup> CISUC, LASI, <sup>2</sup> DEI, University of Coimbra, Portugal

<sup>3</sup> Polytechnic Institute of Coimbra, Portugal  
{patriciaf,isabelc,ana,catarina,hroliv}@dei.uc.pt

## Abstract

Customer-support services increasingly rely on automation, whether full or with human intervention. Despite optimising resources, this may result in mechanical protocols and lack of human interaction, thus reducing customer loyalty. Our goal is to enhance interpretability and provide guidance in communication through novel tools for easier analysis of message trends and sentiment variations. Monitoring these contributes to more informed decision-making, enabling proactive mitigation of potential issues, such as protocol deviations or customer dissatisfaction. We propose a generic approach for dialogue flow discovery that leverages clustering techniques to identify dialogue states, represented by related utterances. State transitions are further analyzed to detect prevailing sentiments. Hence, we discover sentiment-aware dialogue flows that offer an interpretability layer to artificial agents, even those based on black-boxes, ultimately increasing trustworthiness. Experimental results demonstrate the effectiveness of our approach across different dialogue datasets, covering both human-human and human-machine exchanges, applicable in task-oriented contexts but also to social media, highlighting its potential impact across various customer-support settings.

## 1 Introduction

Dialogue systems are increasingly pervasive, playing a crucial role in communication with customers in many companies. Monitoring and visualizing conversations produced by such systems offers a deeper comprehension of dialogue interactions, unveiling communication patterns, and providing valuable insights into the user experience. It is thus essential to ensure high-quality service. Here, the analysis of frequent dialogue flows plays an important role, as they will depict the organic evolution of interactions, enhancing human interpretability.

Obtaining dialogue flows from black-box systems, such as chatbots based on Large Language

Models (LLMs) or other encoder-decoder frameworks, can be challenging due to their generative and open-domain nature. Nonetheless, the ability to represent the conversation progression and consider emotional aspects such as the sentiment of the speakers is valuable, especially in activities requiring real-time assistance from responsible agents.

We propose an approach for automatic dialogue flow discovery from a history of written dialogues, and their representation in a transition graph. We begin by grouping similar utterances into clusters, which may be seen as dialogue states. Then we represent possible paths with their respective probabilities from the beginning to the end of the dialogue.

Furthermore, we enrich the states with the average sentiment of the included utterances. This has applications in a wide range of services and products involving dialogue or customer support, including call centers, emergency services, and virtual assistants. It also serves as an assessment tool, offering stakeholders a way to compare dialogue systems based on how they handle client requests while maintaining or improving their sentiment. Moreover, this approach can potentially identify topics that often result in negative sentiment. The main contributions of this work are summarized as:

- The proposal of a solution for the automatic discovery of dialogue flows that are adaptable to any language and domain, offering an interpretability layer to dialogue systems;
- The integration of sentiment analysis into existing/automatically generated flows, enriching interpretability with sentiment variations;
- The proposal of flow metrics for assessing (i) agents' performance based on sentiment variation, (ii) effectiveness in capturing common states, and (iii) sentiment and cluster cohesion within flows;
- A visual analysis of flows discovered from

diverse dialogue datasets, spanning various services and types, complemented by the proposed metrics, while showcasing the proposed approach and confirming its benefits;

- A proposal for an advanced analysis layer that includes sentiment variation representation within each cluster, offering valuable insights for assessing agent performance and identifying sentiment-associated states.

The remainder of the paper is structured as follows: Section 2 reviews work related to dialogue flow discovery and sentiment analysis; Section 3 describes the proposed approach for sentiment-aware dialogue flow discovery; Section 4 clarifies the meaning of each element in the sentiment-aware dialogue flows, and describes the experimental setup, the used datasets, and the flow metrics proposed; Section 5 presents and analyses the resulting flows; Finally, Section 6 concludes the paper and provides cues for future work.

## 2 Related Work

The categorization of utterances in dialogue systems may help in understanding user intentions and facilitating effective interactions (Deng et al., 2023; Gonalo Oliveira et al., 2022). Generally, utterances are classified according to user intentions (Vedula et al., 2020; Mou et al., 2022) or dialogue acts (Ribeiro et al., 2019; Liu et al., 2017), both providing valuable insights for task-oriented systems. However, the automatic classification of utterances is typically supervised and thus relies on annotated data, which is not always available. On the other hand, encoder-decoder systems, including those based on LLMs (e.g., ChatGPT<sup>1</sup>), do not rely on such classifications, but their flexibility comes at the cost of higher data demand and less control.

Traditional task-oriented dialogue systems are sustained by the design of flows to guide conversations towards specific goals. This entails defining specific user intentions and training phrases, and can be facilitated by tools like Google’s DialogFlow<sup>2</sup>, Microsoft Luis<sup>3</sup>, or Rasa<sup>4</sup>. Automating this process involves grouping semantically-similar utterances and representing them in a vector space, towards efficient intent discovery (Hashemi et al., 2016; Park et al., 2022; Liu et al., 2021).

<sup>1</sup><https://chat.openai.com/>

<sup>2</sup><https://cloud.google.com/dialogflow>

<sup>3</sup><https://www.luis.ai/>

<sup>4</sup><https://rasa.com/>

Representing dialogue flows as transition graphs offers insights on topics and other trends (Bouraoui and Lemaire, 2017). An earlier approach (Ritter et al., 2010) for flow discovery uses Hidden Markov Models on Twitter conversations. It introduces features like clustering similar utterances, vertices for marking the beginning and end of dialogues, as well as a threshold for ignoring low-probability transitions. Towards interpretability, clusters were labelled manually. Ferreira et al. (2024) developed a similar approach with automatic labelling.

By analysing communication trends, flow discovery may assist in the design of dialogue systems. This is the main goal of Graph2Bots (Bouraoui et al., 2019), which adopts co-clustering for discovering dialogue states and transitions in human-human conversations. An alternative approach (Sastre Martinez and Nugent, 2022) clusters utterances with DBSCAN and relies on finite-state automata for discovering ranked flows, based on the frequency of question-response sequences.

Sentiment Analysis (SA) (Liu, 2015) aims to extract sentiments from texts. In dialogues, it may help in identifying situations of sentiment degradation, which may then be acted upon, e.g., through a fallback system that replaces an artificial agent by a human; or by collecting information for later retraining the human or artificial agent.

SA has been combined with other tasks, such as dialogue act recognition, which reinforce one another. For instance, detecting agreement often corresponds with the expression of the same sentiment, while transitions from negative to neutral tend to coincide with changing to a statement. Works that tackled these tasks jointly (Xu et al., 2023; Qin et al., 2020; Li et al., 2020) benefited from it, and achieved high or state-of-the-art (SOTA) performances in datasets like Mastodon (Cerisara et al., 2018). Moreover, Song et al. (2023) outperformed several SOTA methods for user satisfaction estimation in task-oriented dialogue systems by exploiting SA in a multi-task adversarial strategy.

The seemingly symbiotic relationship between SA and other tasks motivated its application to dialogue flow discovery. Yet, to the best of our knowledge, no other work has combined these tasks.

## 3 Proposed Approach

We propose a generic approach for automatically discovering the most common flows in a history of dialogues, while simultaneously associating sen-

timents with their transitions. It comprises three distinct steps, outlined in Figure 1:

1. **Utterance Clustering** clusters semantically similar utterances, represented by their embedding. Discovered clusters may be seen as approximations to dialogue states.
2. **Flow Discovery** computes the most frequent paths. The result is a transition graph  $G(C, T)$ , where nodes  $c \in C$  represent dialogue states and edges  $t(c_i, c_j, p_{ij}) \in T$  represent transitions. The latter are weighted according to their probability, computed as in Equation 1, where  $|t(c_i, c_j)|$  represents the number of utterances in  $c_j$  that immediately follow a utterance in  $c_i$ .

$$p_{ab} = \frac{|t(c_a, c_b)|}{\sum_{x \in C} |t(c_a, c_x)|} \quad (1)$$

3. **Sentiment Classification** enriches  $G$  with sentiment information. When sentiment is not available with the data, each utterance’s sentiment can be determined by an external tool for the purpose. This enables the computation of the most predominant sentiment in each state and transition.

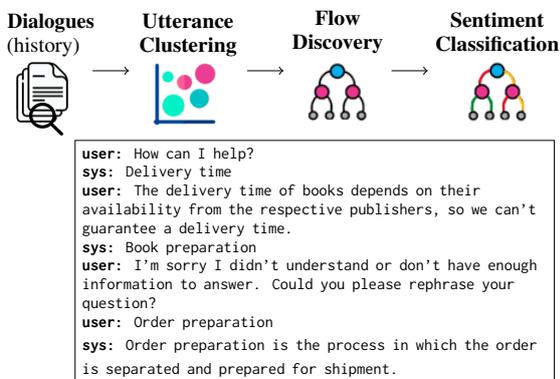


Figure 1: Overview of the proposed approach and an illustrative dialogue between a customer (user) and an artificial agent (sys).

The proposed approach can be applied to any collection of written dialogues, ideally with two speakers, but of any type, in any language, or domain, as long as utterances are provided in sequence and speakers are identified. The graph  $G$  can be visually represented, thereby enhancing human interpretation. The sentiment may be visually represented with different colours for each edge, such as green, yellow, and red, representing predominantly positive, neutral, or negative sentiments, respectively (as shown on the right-hand side of Figure

1). Sentiment-aware dialogue flows can be useful in various scenarios, including:

**Identifying communication trends** i.e., the discovery of flows from any type of dialogue promotes the identification of common and/or undesired topics or transitions, which can be used to improve the agent, e.g., by changing intents, reviewing protocols, or adjusting human resources;

**Interpreting black-box dialogue systems** i.e., the discovery of flows in human-machine dialogues adds an interpretability layer that increases understanding of the agent and promotes the identification of issues. Potential strategies for addressing such issues may include retraining the agent or implementing additional rules;

**Planning and developing dialogue systems** i.e., the analysis of human-human dialogues towards the identification of potential dialogue states and representative words or sentences, valuable to the agent’s development process.

In any scenario, the dialogue collection should be as comprehensive as possible and, ideally, cover all relevant intents. The set of applications attests to the versatility of the approach. Still, in this paper, we focus on the interaction between a customer and an agent, where the ability to understand and efficiently manage interactions is essential for improving the quality of service and, consequently, customer satisfaction.

## 4 Experimentation

In order to confirm the applicability of the sentiment-aware dialogue flows, extensive experimentation was conducted. This involved the implementation of each step of the proposed approach, introduced in Figure 1, with adequate tools, as well as the application to a range of dialogue datasets. This section details the implementation of the underlying processes but, before delving into the previous steps, we provide some clarifications on the visual notation used throughout the paper, aided by the illustrative diagram in Figure 2.

The diagram ( $G$ ) showcases the ideal scenario, in which an agent successfully manages to switch the customer’s sentiment from negative, at the Start Of the Dialogue (SOD), to positive, by the End Of the Dialogue (EOD). SOD and EOD are represented by specific nodes, which can be seen as states, represented as yellow boxes. The others

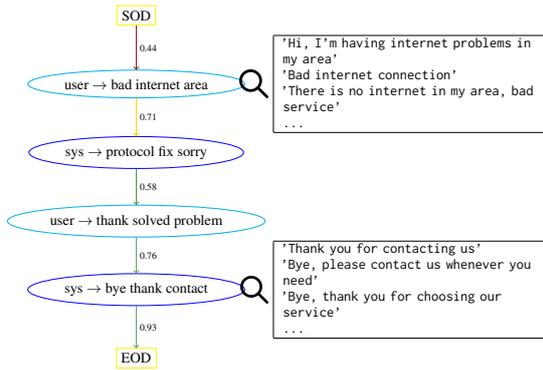


Figure 2: Example of a sentiment-aware dialogue flow showcasing an ideal scenario.

correspond to the discovered clusters ( $C$ ) and are represented by ellipses labelled with representative words in their utterances. States representing clusters by different speakers are also differentiated. In  $G$ , user clusters are coloured in light blue and agent clusters in dark blue. The diagram is complemented with examples of clustered utterances, on the right-hand side.

Edges represent transitions ( $T$ ) between clusters and have an associated weight, corresponding to their probability. For instance,  $G$  shows a 58% probability of moving from `sys→protocol fix sorry` to `user→thank solved problem`. The sum of all probabilities of  $T$  originating from the same cluster is 1. Nevertheless, in order to simplify the flow, a threshold can be applied for ignoring low-probability transitions, as carried out in this example. The colour of each transition represents the average sentiment within the destination cluster. Red corresponds to a more negative sentiment, green to a positive, and yellow to a neutral sentiment. For EOD, however, transitions represent the sentiment of the origin cluster, thus making the sentiment in the final interactions clearer and contributing to more immediate conclusions.

#### 4.1 Datasets

The proposed approach was applied to five different dialogue datasets, covering different channels (social media, chat, telephone), types of dialogue (task-oriented, open) and agent (human, machine), domains (tourism, telecommunications, retail, open) and languages (English and Portuguese). Specifically, the following datasets were used:

**EmoWOZ (Feng et al., 2022)** a public dataset of task-oriented dialogues that extends MultiWOZ (Budzianowski et al., 2018), thus covering

multi-domains related to tourism. EmoWOZ’s additionally has emotions assigned to utterances, including valence, translatable to a polarity (positive, neutral, negative).

**TwitterDialogueSAPT (TDSAPT) (Carvalho et al., 2023)** a public dataset of customer-support dialogues in Portuguese, extracted from Twitter, with entities (i.e., accounts) in the domains of Telecommunications, Television, Healthcare, eCommerce, and Finance, where utterances have manually-annotated sentiment. We adopted the original approach for extending this dataset with more dialogues from the same entities in the same timeline (April–May, November–December 2022).

**TelecomSAPT** transcriptions of customer-support dialogues, sampled from two months in the call center of a Portuguese Telecommunications company, with manually-labelled sentiment.

**RetailPT** a collection of customer-support dialogues of a Portuguese retail company, collected during a seasonal campaign that lasted 2.5 months (July–September 2023). Dialogues are between human customers and a proprietary Retrieval Augmented Generation system based on fine-tuning an optimised version of Quokka<sup>5</sup>.

**Mastodon (Cerisara et al., 2018)** a public dataset of dialogues extracted from the Mastodon social network, particularly from the octodon.social instance, with manually-annotated sentiment. These are open-domain conversations between two users and, as such, do not involve a service.

For some datasets, we could get the polarity of the utterances from available annotations. This was, however, not the case of RetailPT and the extension of TDSAPT, which employed a classifier fine-tuned in similar data (see Section 4.2).

Since the labels in TDSAPT were binary (negative and non-negative), we binarised the labels of all datasets, which still enabled the identification of negative transitions, the most problematic.

Table 1 describes the datasets according to channel (Chat, Phone, Social Media - SM) type of dialogue (Task Oriented - TO; Open) type of agents, domain, language (English - EN; Portuguese - PT), and number of dialogues.

<sup>5</sup>[hf.co/automaise/quokka-7b](https://hf.co/automaise/quokka-7b)

<sup>6</sup>Both TelecomSAPT and RetailPT were gently transferred to our team in the scope of projects with the industry, but are proprietary and cannot be publicly released.

Dataset	Channel	Type	Agent	Domain	Lang	#Dialogs
EmoWOZ	Chat	TO	Human	Tourism	EN	10,253
RetailPT	Chat	TO	Machine	Retail	PT	3,317
TelecomSAPT	Phone	TO	Machine	Telecom	PT	1,000
TDSAPT	SM	TO	Human	Several	PT	2,575
Mastodon	SM	Open	Human	Open	EN	535

Table 1: Brief description of each dataset, including channel, type of dialogue, type of agents, domain, language, and number of dialogues.

Table 2 presents the number of utterances in each dataset, the sentiment distribution (negative and non-negative) and informs on how the sentiment labels were obtained: in the data (D), automatic (A) by a supervised model, converted (C).

Dataset	# Utterances	% Neg	% Non-Neg	Source
EmoWOZ	140,801	1.57	98.43	C
RetailPT	19,098	28.79	71.21	A
TelecomSAPT	5,312	18.39	81.61	D
TDSAPT	5,966	36.15	63.85	D+A
Mastodon	2,205	31.61	68.39	D

Table 2: Analysis of the sentiment distribution in each dataset, including the source of sentiment labelling.

Tables 1 and 2 confirm the diversity of covered scenarios. They encompass various channels, dialogue types, agents, domains, and languages, attesting to the generalisation potential of the proposed approach. Datasets differ in size and prevalence of negative sentiment, spanning from as low as 1.6% of utterances in EmoWOZ to 36% in TDSAPT.

## 4.2 Implementation

Utterance embeddings were obtained from sentence transformers available in the HuggingFace Hub. Different models were used for English<sup>7</sup> and Portuguese<sup>8</sup>, both representing textual sequences in 384-dimension vectors.

Clustering was performed with the K-means method, as available in scikit-learn<sup>9</sup>. For each dataset, the number of clusters was optimised for maximizing the Silhouette score (Rousseeuw, 1987), which evaluates the cohesion and separation of formed groups. This relied on Optuna<sup>10</sup>, considering a range of 3–10 clusters for each speaker.

For the textual labels of the clusters, a document is created for each cluster, with its concatenated utterances. Using the same models as in the clustering step, the label resulted from the most frequent keyword for each cluster, obtained with KeyBERT (Grootendorst, 2020), considering a range

<sup>7</sup>[hf.co/sentence-transformers/all-MiniLM-L6-v2](https://hf.co/sentence-transformers/all-MiniLM-L6-v2)

<sup>8</sup><https://tinyurl.com/2fcwpuz7>

<sup>9</sup><https://tinyurl.com/4ymet8ff>

<sup>10</sup>[optuna.org/](https://optuna.org/)

of [1–3]-grams, and after removing stopwords in the NLTK (Bird and Loper, 2004) lists.

The sentiment of unlabeled utterances in TDSAPT and RetailPT was classified with a BERT model pretrained for Portuguese (Souza et al., 2020), fine-tuned for identifying negative and non-negative sentiments in Portuguese dialogues, in a similar fashion to the best model in related work (Carvalho et al., 2023). The main difference was the fine-tuning datasets, selected for sharing more similarities with the data to classify: in the extension of TDSAPT, the model was fine-tuned in the original dialogues of TDSAPT, with a 75% F1-score on it, whereas in RetailPT it was fine-tuned in TelecomSAPT, with a 74% F1-score on the former.

Finally, for representing the sentiment in each cluster, we compute the average sentiment in all its utterances. If the average sentiment is low ( $<0.4$ ), high ( $>0.6$ ), or in-between, we colour the incoming transitions in red, green or yellow, respectively. As the range of values associated with green and red is larger, we further define colour gradients: if the average sentiment is closer to 0.0 or 1.0, the corresponding colour gets darker. We recall that, as an average, this value may not represent the sentiment of all the utterances in each cluster. Hence, we propose a second, more in-depth analysis that includes the standard deviation (STD) of the sentiment in each cluster. Specifically, we compute: (i) the average sentiment (AVG); (ii) the sentiment at the highest deviation point (AVG+STD and assigning the corresponding colour); and (iii) the sentiment at the lowest deviation point (AVG-STD). This is considered in the graphical visualisation by adding a three-layered box to each cluster, with a larger middle layer coloured with the average sentiment, and the others with the sentiment at the lowest (left) and highest (right) deviation points. Some resulting dialogue flows are presented in Section 5.

## 4.3 Flow metrics

The discovered flows contribute to faster analysis of trends in the underlying dialogue datasets, but comparing flows from different datasets can still be subjective. To complement the analysis and make the comparison more straightforward, we designed objective metrics, computed directly from the flows. They capture the following aspects: (i) the agents’ performance based on the sentiment throughout the dialogue flow; (ii) how well the clusters represent the dataset based on the proportion of dismissed utterances; (iii) the flow’s cohesion regarding sen-

timent and its clusters. The computed metrics are described ahead, with customer support in mind.

*EOD<sub>-</sub>*: proportion of utterances that reach EOD with a negative sentiment. This can be applied to other sentiment levels but negativity is the one that should be mitigated;

$\Delta Sentiment$ : difference between the sentiment of the utterances at the end of the dialogue and of those at the start.

$SOD \rightarrow EOD$ : proportion of utterances from all speakers that were dismissed in the flow, i.e., those that took paths with a probability lower than the set threshold, ending up not represented;

**Flow Cluster Cohesion (FCC)**: average Silhouette score of the clusters;

**Flow Sentiment Cohesion (FSC)**: average standard deviation of the sentiment at each cluster;

**Average Initial Sentiment (AIS)**: average sentiment of each cluster with an incoming transition from SOD. As opposed to the values considered in  $\Delta Sentiment$ , this is calculated by cluster (i.e., each contains the average sentiment of the utterances within) and not by utterance.

**Average Final Sentiment (AFS)**: average sentiment of each cluster with outgoing transitions to EOD. As opposed to  $\Delta Sentiment$ , this is computed by cluster, not by utterance.

An analysis of these metrics should be enough to get insights on the performance of the agent(s). Ideally, it would present (i) a low *EOD<sub>-</sub>*, i.e., managed to avoid negative sentiment, (ii) a positive  $\Delta Sentiment$ , i.e., sentiment improved throughout the flow, (iii) a low  $SOD \rightarrow EOD$ , i.e., most utterances were represented, (iv) a high FCC, i.e., data fits the clusters well, and (v) a low FSC, i.e., sentiment at each cluster does not deviate much. Finally, AIS and AFS should be analysed together as the latter should be higher than the former, i.e., sentiment at the cluster level should improve throughout the flow. The next section reports on applying these metrics to the considered datasets.

## 5 Results and Discussion

Flows were discovered from every considered dataset and the designed metrics were computed as well. Together, they provide insights into inter-pretability, communication trends, and limitations

of the agents, among others. This section discusses some of the discovered flows and reports on the metrics computed for all. Due to lack of space, we do not present the flows for all datasets, but include them in the Appendix A.

The dialogue flow for RetailPT data is presented in Figure 3<sup>11</sup>. Various interactions between the user and the (artificial) agent can be observed. We immediately note that the first interaction of the agent (SOD’s outgoing transition), is always the same, with probability 1.0. The label of the initial cluster suggests an offer of assistance, which is confirmed by the data: in fact, all dialogues start with the How can I help? utterance.

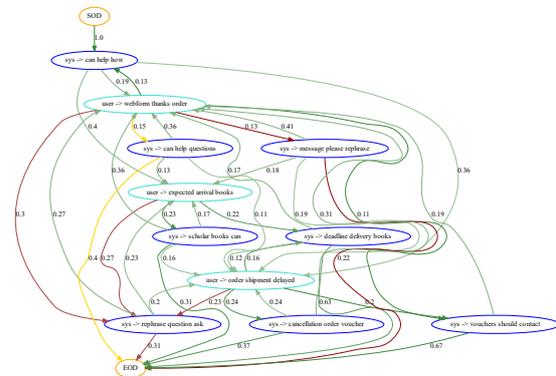


Figure 3: Sentiment-aware dialogue flow discovered for RetailPT.

With the help of the labels, we see that this interaction is followed by a user message: thanking for the order; querying about the arrival of the books; or informing on a shipment delay. As the probabilities of each transition from the can help how state do not sum up to 1.0, there is at least one low-probability transition (i.e.,  $p < 0.1$ ) not represented. Afterwards, the agent replies and, in some cases, asks the user to rephrase the question. Interactions continue until the EOD, which marks the end of the conversation.

In this case, non-negative sentiments (i.e., shades of green) predominate. Additionally, there are edges with neutral colours leading to the can help questions state and red edges associated with negative sentiments in the rephrase question ask and message please rephrase states. We may deduce that the agent is failing to process the requests, which could potentially increase the user’s frustration as some end the conversation afterwards.

<sup>11</sup>RetailPT data is in Portuguese. For an easier interpretation by the readers, cluster labels were translated to English.

This negative sentiment could potentially be mitigated by retraining the agent to better handle the queries that lead to those clusters.

It is important to note that the colour of the edges represents the average sentiment of utterances in each cluster, which may not fully capture the sentiment of the entire cluster. Hence, we created an advanced analysis layer, shown in Figure 4, which considers the sentiment’s standard deviation for each cluster via a three-layered cluster.

In the can help questions state, the average sentiment is represented by the colour yellow (i.e., neutral). Its left layer (red) represents the sentiment at its lower deviation value and its right layer (green) represents the sentiment at its highest deviation value. In this case, the average does not accurately represent the sentiment within that cluster as it also includes strong negative and positive values (i.e., deep shades of red and green).

In states such as scholar book can or vouchers should contact, there is minimal sentiment deviation, as each layer of the node appears uniformly green, suggesting that the sentiment within the utterances of underlying clusters is accurately represented by their average.

Table 3 reports on metrics computed for the utterances’ transitions and their sentiment. We recall that these can be used to evaluate an agent’s performance and how well the flow captures common states, i.e., represents most utterances.

EmoWOZ has  $EOD_- = 0$ , meaning no dialogue ends with negative sentiment. Moreover, it has the highest  $SOD \rightarrow EOD$ , meaning that, with the applied threshold (0.1), most utterances are lost along the way. As this is the largest dataset (seven times larger than RetailPT) it makes sense that it would be challenging to represent each utterance in it. Sentiment variation is the lowest for this dataset.

TelecomSAPT has the highest  $EOD_-$ , meaning it is the dataset that mostly finished with negative sentiment, followed by RetailPT. This means that the involved (artificial) agents could benefit from an in-depth analysis, possibly culminating in reviewing and/or retraining. These are also the only datasets with a negative sentiment variation, i.e., by the end of the dialogue, sentiment gets lower. They also show high  $SOD \rightarrow EOD$ , as does Mastodon, meaning these three datasets lose over half of their utterances throughout the flow.

Mastodon and TDSAPT show intermediate values overall and the latter has the lowest  $SOD \rightarrow EOD$ , meaning that more than half the utterances

are represented in the flow. Both datasets have a positive sentiment variation, suggesting an improvement by the end of the conversation.

In both cases, it is not easy to speculate more. Mastodon has social media dialogues, where sentiment can flow, without clear negative consequences as in customer-support. Moreover, TDSAPT includes dialogues with a broad range of entities, and would benefit from a future analysis of the flows for each, independently.

Dataset	$EOD_-$	$\Delta Sentiment$	$SOD \rightarrow EOD$
EmoWOZ	0.0	0.02	0.83
RetailPT	0.25	-0.28	0.63
TelecomSAPT	0.34	-0.06	0.55
Mastodon	0.08	0.18	0.65
TDSAPT	0.12	0.06	0.43

Table 3: Evaluation metrics for assessing agents’ performance and flow’s ability to capture common states.

Table 4 presents metrics for assessing the cohesion of flows regarding sentiment and clusters. In EmoWOZ no dialogue ends with a negative sentiment (1.00 AFS). It has also the lowest FSC, i.e., sentiment does not vary much within each cluster.

RetailPT has the highest AIS, however, AFS suggests that sentiment gets worse by the end of the dialogues. It has also the highest FCC, meaning that the data is well-fitted to the clusters.

Mastodon has the highest FSC, meaning that, contrary to EmoWOZ, sentiment diverges considerably within each cluster. However, AIS and AFS suggest that it increases by the end of the dialogue. It also presents the lowest FCC, meaning that data may not be well-fitted to the clusters, which aligns with the high divergence of sentiment within them.

TelecomSAPT and TDSAPT display intermediate results in flow cohesion and variation of sentiment within clusters. However, whereas the former’s AIS and AFS suggest sentiment across the dialogues is predominantly positive, for the latter, AIS and AFS have the lowest values, suggesting a more neutral sentiment. For TDSAPT, the difference between AIS and AFS is low, as is the  $\Delta Sentiment$ , but in different directions. The former value should be more accurate as it computes the variation by utterance instead of cluster.

Dataset	FCC	AIS	AFS	FSC
EmoWOZ	0.11	0.96	1.00	0.12
RetailPT	0.46	1.00	0.67	0.29
TelecomSAPT	0.26	0.82	0.71	0.26
Mastodon	0.04	0.68	0.71	0.41
TDSAPT	0.14	0.57	0.55	0.31

Table 4: Flow cohesion metrics for considered datasets.

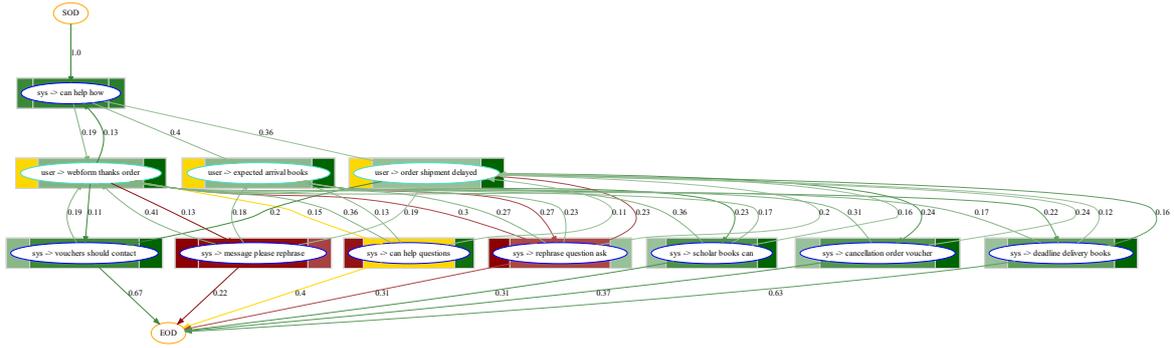


Figure 4: Sentiment-aware dialogue flow with standard deviation discovered for RetailPT data.

Finally, two factors could impact the discovery of sentiment-aware dialogue flows and, thus, their analyses: (i) the sentiment classifier, and (ii) the flow discovery process, including the clustering and labelling methods. The low performance of any of these can cause a chain reaction, decreasing the quality of the final analysis. As such, performance at each level should always be considered.

## 6 Conclusion

Technological advances have increased reliance on Artificial Intelligence, including for customer-support services. While efficient cost-wise, customers can tell they are interacting with an artificial agent or a human following a mechanical protocol, and this degrades their interaction and deteriorates the customers’ loyalty. Our goal is to mitigate that by providing additional interpretability, also contributing to increased trustworthiness.

We proposed a novel approach for automatically discovering the most common flows in a history of dialogues, while considering the sentiment. These are useful for various applications, from identifying communication trends to interpreting black-box dialogue systems, and contribute to uncovering the triggers of problematic situations.

Our solution is independent of domain and language, and does not require dialogues labelled with intents or acts. Its implementation enabled the discovery of flows from a diverse set of dialogue datasets, out of which interesting insights were gathered, also with the help of computed metrics. For instance, in dialogues with artificial agents (RetailPT, TelecomSAPT), sentiment gets worse throughout the flow. The automation of such agents results in more mechanical answers and, thus, more cohesive clusters (FCC), when compared to other datasets. Mastodon and TDSAPT were collected from social media and cover multi-

ple domains, which contributes to a higher variation of sentiment (FSC). Metrics also reveal that, with the parameters set (i.e., probability threshold of 0.1 and maximum 10 clusters for speaker), a large portion of utterances is lost in the flow discovery process. These regard low-probability transitions, but may degenerate interpretation, especially for large datasets as EmoWOZ. Yet, the alternative would be either to: reduce the number of clusters, with an impact on cohesion; or increase both the number of clusters and the threshold, with a negative impact in interpretability. Therefore, we plan to test alternative implementations and analyze their impact on the previous, including clustering and labelling methods, and sentiment classification, where new trends (Zhang et al., 2023) can be explored. The computation of more metrics should also be considered, e.g., for assessing the coverage of discovered flows in unseen dialogues from the same domain. Finally, towards stronger conclusions, flows should be discovered from additional datasets.

Another focus will be on flow visualization. While moving away from a graph-based model is unlikely, we consider integrating additional elements (e.g., reflecting the number of utterances in the node’s size) and interactivity, towards improved interpretability (e.g., selecting the best threshold; highlighting the path taken in a specific dialogue).

## Acknowledgements

This work was supported by: the project POWER (POCI-01-0247-FEDER-070365), co-financed by the European Regional Development Fund (FEDER), through Portugal 2020 (PT2020), and by the Competitiveness and Internationalization Operational Programme (COMPETE 2020); the Portuguese Recovery and Resilience Plan through project C645008882-00000055, Center for Responsible AI; and by national funds through

FCT, within the scope of the project CISUC (UID/CEC/00326/2020).

## References

- 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.
- Jean Léon Bouraoui, Sonia Le Meitour, Romain Carbou, Lina M Rojas Barahona, and Vincent Lemaire. 2019. Graph2bots, unsupervised assistance for designing chatbots. In *Procs. 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 114–117. ACL.
- Jean-Leon Bouraoui and Vincent Lemaire. 2017. Cluster-based graphs for conceiving dialog systems. In *Procs ECML-PKDD 2017 Workshop on Interactions between Data Mining and Natural Language Processing*. CEUR-WS.org.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Isabel Carvalho, Hugo Gonçalo Oliveira, and Catarina Silva. 2023. The importance of context for sentiment analysis in dialogues. *IEEE Access*, 11:86088–86103.
- Christophe Cerisara, Somayeh Jafaritazehjani, Adedayo Oluokun, and Hoa T Le. 2018. Multi-task dialog act and sentiment recognition on Mastodon. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 745–754.
- Yang Deng, Wenqiang Lei, Wai Lam, and Tat-Seng Chua. 2023. [A survey on proactive dialogue systems: Problems, methods, and prospects](#). In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pages 6583–6591. International Joint Conferences on Artificial Intelligence Organization. Survey Track.
- Shutong Feng, Nurul Lubis, Christian Geishauer, Hsien-chin Lin, Michael Heck, Carel van Niekerk, and Milica Gasic. 2022. [EmoWOZ: A large-scale corpus and labelling scheme for emotion recognition in task-oriented dialogue systems](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4096–4113, Marseille, France. European Language Resources Association.
- Patrícia Ferreira, Daniel Martins, Ana Alves, Catarina Silva, and Hugo Gonçalo Oliveira. 2024. Unsupervised flow discovery from task-oriented dialogues. *arXiv preprint arXiv:2405.01403*.
- Hugo Gonçalo Oliveira, Patrícia Ferreira, Daniel Martins, Catarina Silva, and Ana Alves. 2022. A Brief Survey of Textual Dialogue Corpora. In *Proceedings of the 13th Language Resources and Evaluation Conference, LREC 2022*, pages 1264–1274, Marseille, France. ELRA.
- Maarten Grootendorst. 2020. [KeyBERT: Minimal keyword extraction with BERT](#). 10.5281/zenodo.4461265.
- Homa B Hashemi, Amir Asiaee, and Reiner Kraft. 2016. Query intent detection using convolutional neural networks. In *International conference on web search and data mining, workshop on query understanding*.
- Jingye Li, Hao Fei, and Donghong Ji. 2020. [Modeling Local Contexts for Joint Dialogue Act Recognition and Sentiment Classification with Bi-channel Dynamic Convolutions](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 616–626, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Bing Liu. 2015. *Sentiment Analysis Mining Opinions, Sentiments, and Emotions*. Cambridge University Press.
- Pengfei Liu, Youzhang Ning, King Keung Wu, Kun Li, and Helen Meng. 2021. Open intent discovery through unsupervised semantic clustering and dependency parsing. *arXiv preprint arXiv:2104.12114*.
- Yang Liu, Kun Han, Zhao Tan, and Yun Lei. 2017. Using context information for dialog act classification in dnn framework. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2170–2178.
- Yutao Mou, Keqing He, Yanan Wu, Zhiyuan Zeng, Hong Xu, Huixing Jiang, Wei Wu, and Weiran Xu. 2022. Disentangled knowledge transfer for ood intent discovery with unified contrastive learning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 46–53.
- Jeiyeon Park, Yoonna Jang, Chanhee Lee, and Heuiseok Lim. 2022. Analysis of utterance embeddings and clustering methods related to intent induction for task-oriented dialogue. *arXiv preprint arXiv:2212.02021*.
- Libo Qin, Zhouyang Li, Wanxiang Che, Minheng Ni, and Ting Liu. 2020. [Co-GAT: A Co-Interactive Graph Attention Network for Joint Dialog Act Recognition and Sentiment Classification](#). *CoRR*, abs/2012.13260.
- Eugénio Ribeiro, Ricardo Ribeiro, and David Martins de Matos. 2019. Deep dialog act recognition using multiple token, segment, and context information representations. *Journal of Artificial Intelligence Research*, 66:861–899.

Alan Ritter, Colin Cherry, and Bill Dolan. 2010. [Un-supervised modeling of twitter conversations](#). In *Human Language Technologies - North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*.

Peter J Rousseeuw. 1987. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65.

Javier Miguel Sastre Martinez and Aisling Nugent. 2022. Inferring ranked dialog flows from human-to-human conversations. In *Procs. 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 312–324, Edinburgh, UK. ACL.

Kaisong Song, Yangyang Kang, Jiawei Liu, Xurui Li, Changlong Sun, and Xiaozhong Liu. 2023. A speaker turn-aware multi-task adversarial network for joint user satisfaction estimation and sentiment analysis. In *Proceedings of the 37th AAAI Conference on Artificial Intelligence*, pages 13582–13590.

Fábio Souza, Rodrigo Nogueira, and Roberto Lotufo. 2020. BERTimbau: Pretrained BERT models for Brazilian Portuguese. In *Proceedings of the Brazilian Conference on Intelligent Systems (BRACIS 2020)*, volume 12319 of *LNCS*, pages 403–417. Springer.

Nikhita Vedula, Nedim Lipka, Pranav Maneriker, and Srinivasan Parthasarathy. 2020. Open intent extraction from natural language interactions. In *Proceedings of The Web Conference 2020*, pages 2009–2020.

Yujun Xu, Enguang Yao, Chaoyue Liu, Qidong Liu, and Mingliang Xu. 2023. [A novel ensemble model with two-stage learning for joint dialog act recognition and sentiment classification](#). *Pattern Recognition Letters*, 165:77–83.

Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023. Sentiment analysis in the era of large language models: A reality check. *arXiv preprint arXiv:2305.15005*.

## A Application to diverse data

In the following sections, we showcase the application of our approach on the five datasets presented: EmoWOZ, RetailPT, TelecomSAPT, TwitterDialogueSAPT, and Mastodon.

### A.1 EmoWOZ

**EmoWOZ (Feng et al., 2022)** is a public dataset of task-oriented dialogues that extends MultiWOZ (Budzianowski et al., 2018), thus covering multi-domains related to tourism. It is the largest dataset covered in this work but also the one with the lowest percentage of negative utterances. It is also the only dataset where sentiment was converted as it is labelled for emotion. Figures 5 and 6

present the two sentiment-aware dialogue flows discovered for this dataset, with the latter presenting the sentiment standard deviation.

### A.2 RetailPT

**RetailPT** is a collection of customer-support dialogues of a Portuguese retail company. Dialogues are between human customers and a proprietary Retrieval Augmented Generation system. It is the second largest dataset covered in this work and the only one in which sentiment analysis was fully automatic, by a supervised model. Figures 7 and 8 present the two sentiment-aware dialogue flows discovered for this dataset, with the latter presenting the sentiment standard deviation.

### A.3 TelecomSAPT

**TelecomSAPT** contains transcriptions of customer-support dialogues, sampled from the call center of a Portuguese Telecommunications company, with manually-labelled sentiment. It is one of the smallest datasets covered in this work and the only one with a voice channel. Figures 9 and 10 present the two sentiment-aware dialogue flows discovered for this dataset, with the latter presenting the sentiment standard deviation.

### A.4 TwitterDialogueSAPT

**TwitterDialogueSAPT (TDSAPT) (Carvalho et al., 2023)** is a public dataset of customer-support dialogues in Portuguese, extracted from the social network Twitter, covering accounts in multiple domains, where utterances have manually-annotated sentiment. This dataset was extended for this work, and sentiment analysis was performed automatically by a supervised model for the new utterances. Figures 11 and 12 present the two sentiment-aware dialogue flows discovered for this dataset, with the latter presenting the sentiment standard deviation.

### A.5 Mastodon

**Mastodon (Cerisara et al., 2018)** is a public dataset of dialogues extracted from the Mastodon social network, particularly from the octodon.social instance, with manually-annotated sentiment. These are open-domain conversations between two users and, as such, do not involve a service. This is the smallest dataset covered by our work and the only one with a fully open domain and type of dialogue. Figures 13 and 14 present the two sentiment-aware dialogue flows discovered for

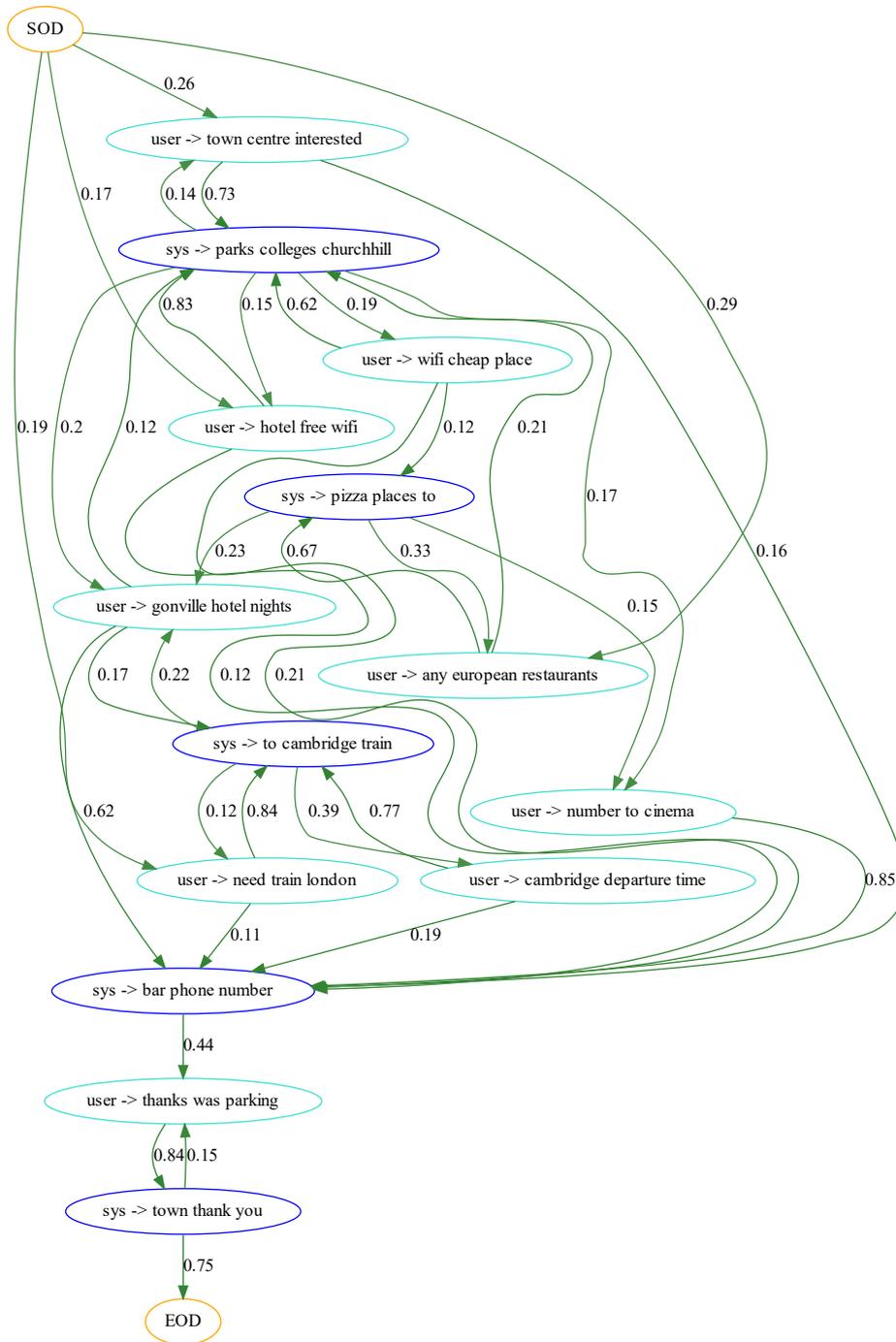


Figure 5: Sentiment-aware dialogue flow discovered for EmoWOZ data

this dataset, with the latter presenting the sentiment standard deviation.

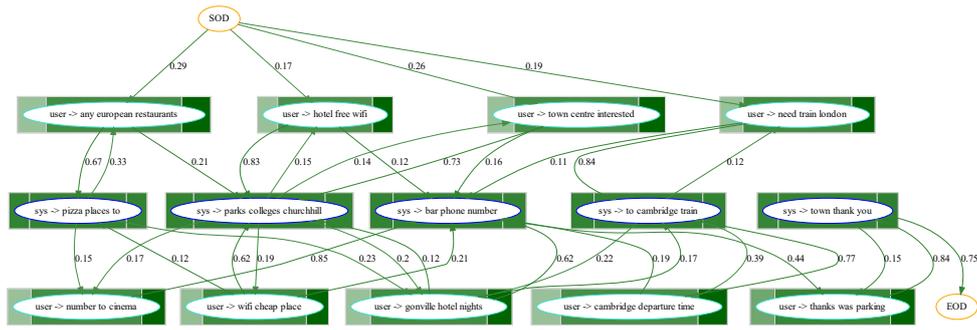


Figure 6: Sentiment-aware dialogue flow with standard deviation discovered for EmoWOZ data

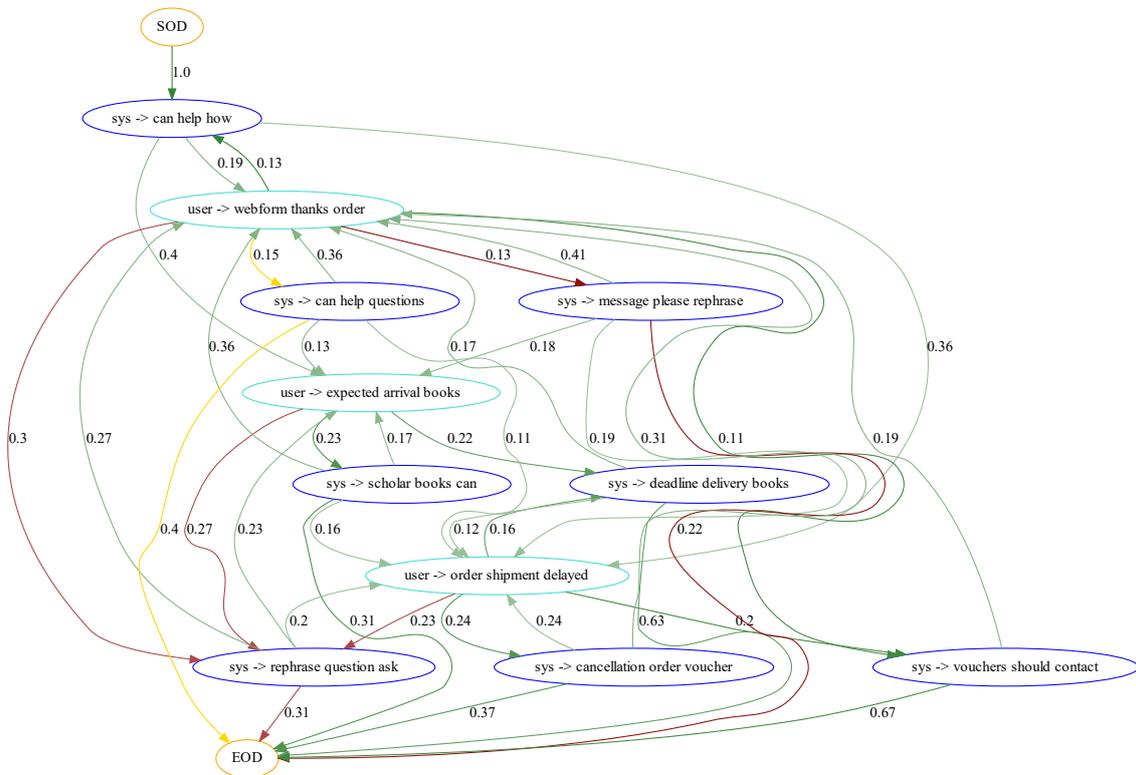


Figure 7: Sentiment-aware dialogue flow discovered for RetailPT data

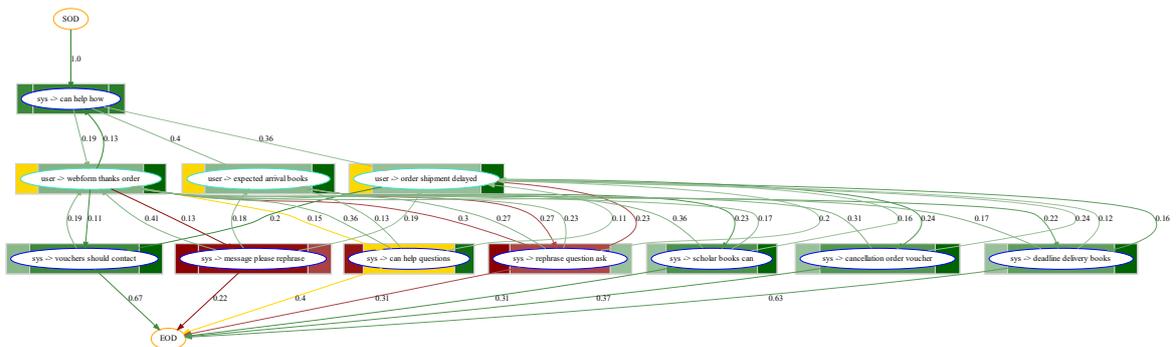


Figure 8: Sentiment-aware dialogue flow with standard deviation discovered for RetailPT data

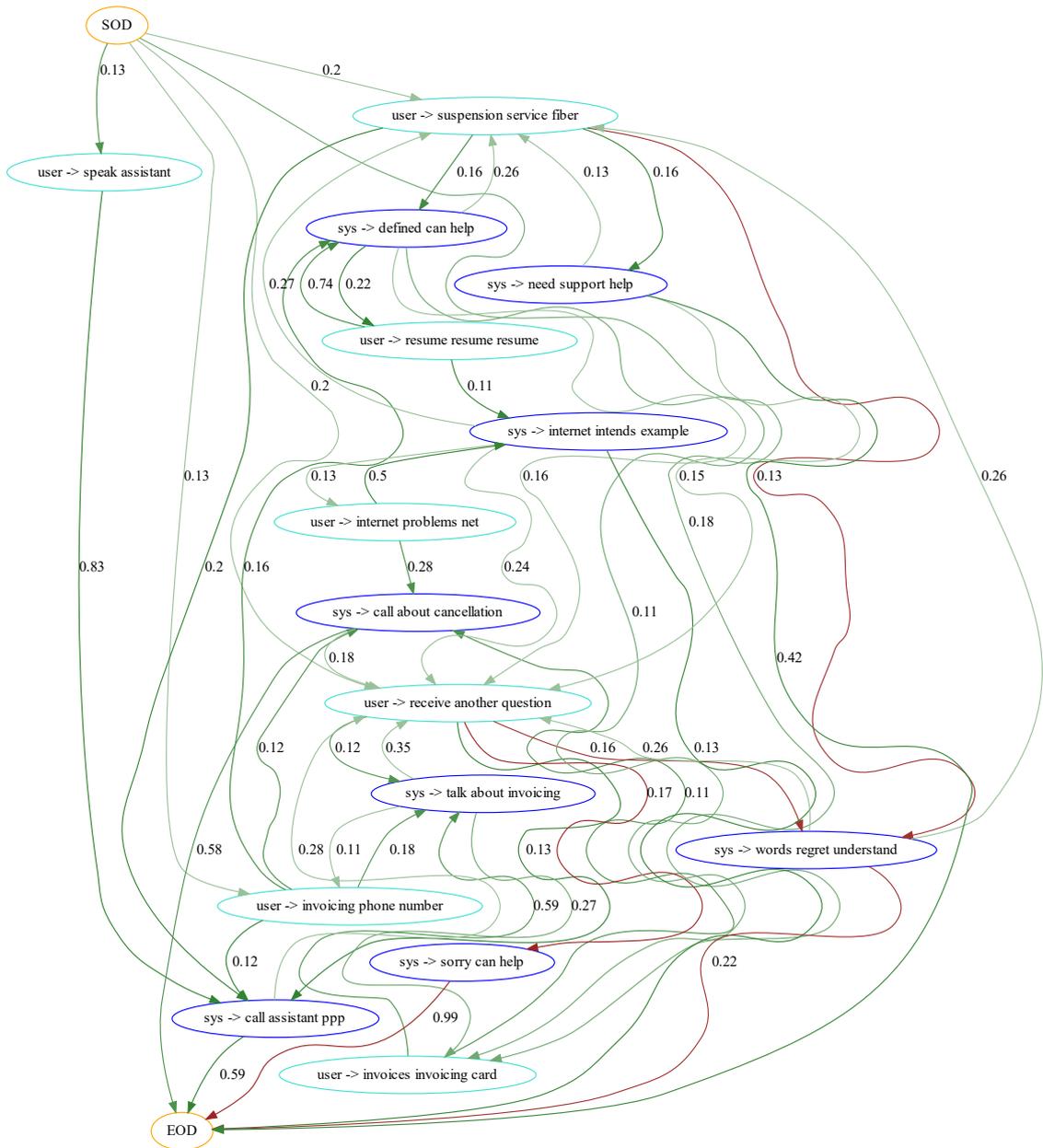


Figure 9: Sentiment-aware dialogue flow discovered for TelecomSAPT data

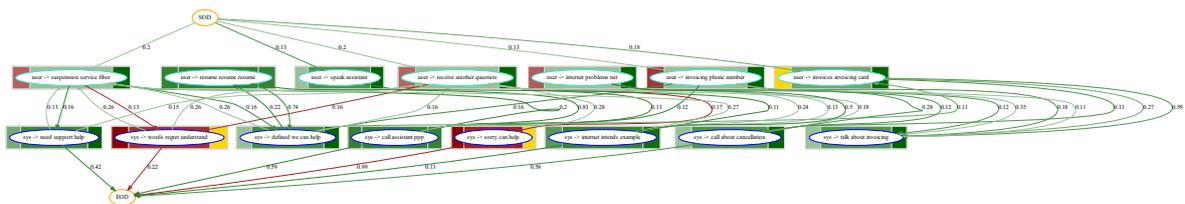


Figure 10: Sentiment-aware dialogue flow with standard deviation discovered for TelecomSAPT data

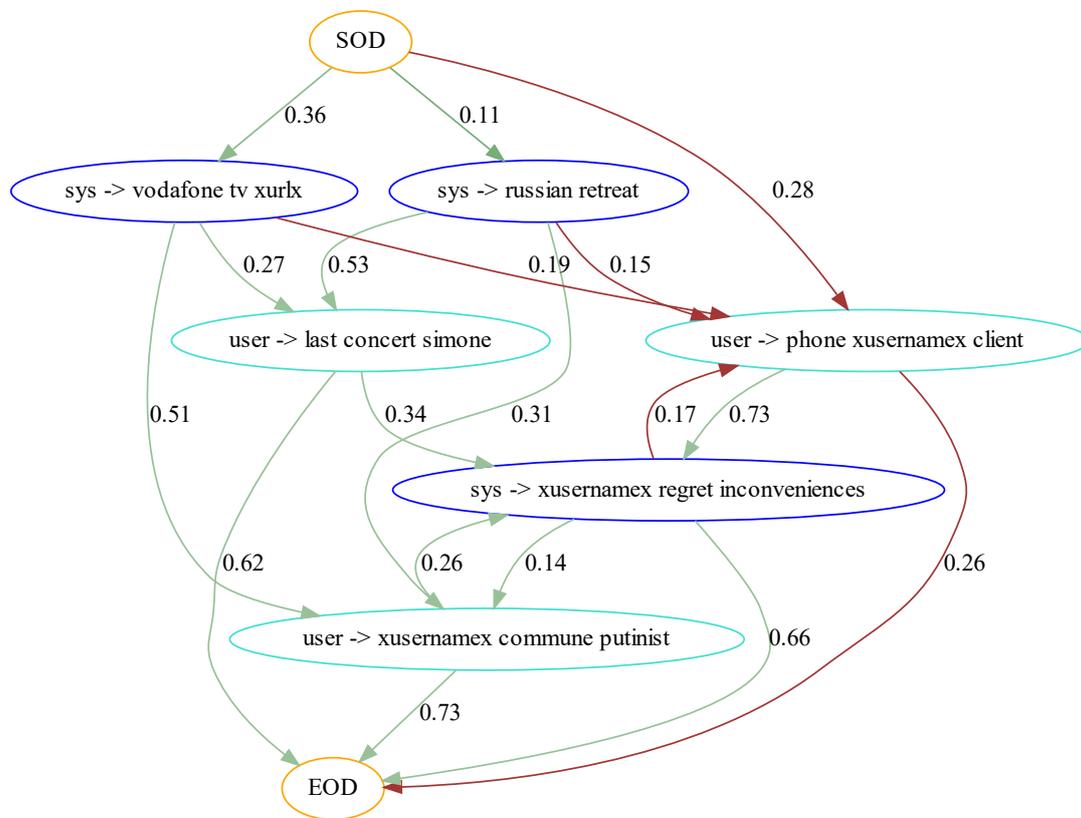


Figure 11: Sentiment-aware dialogue flow discovered for TDSAPT data

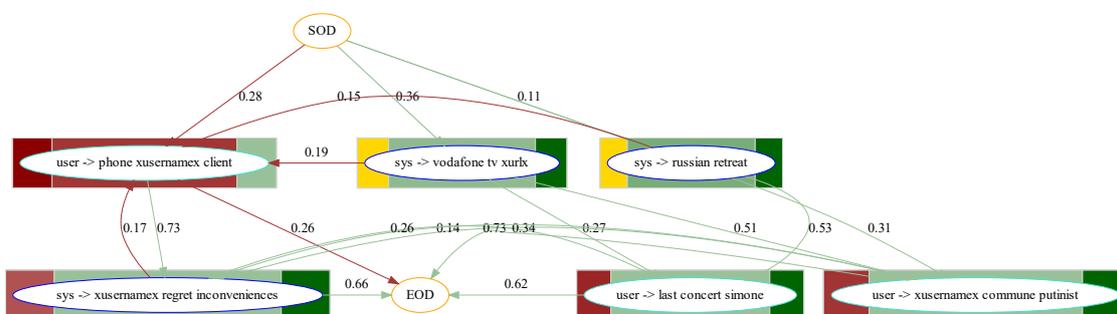


Figure 12: Sentiment-aware dialogue flow with standard deviation discovered for TDSAPT data

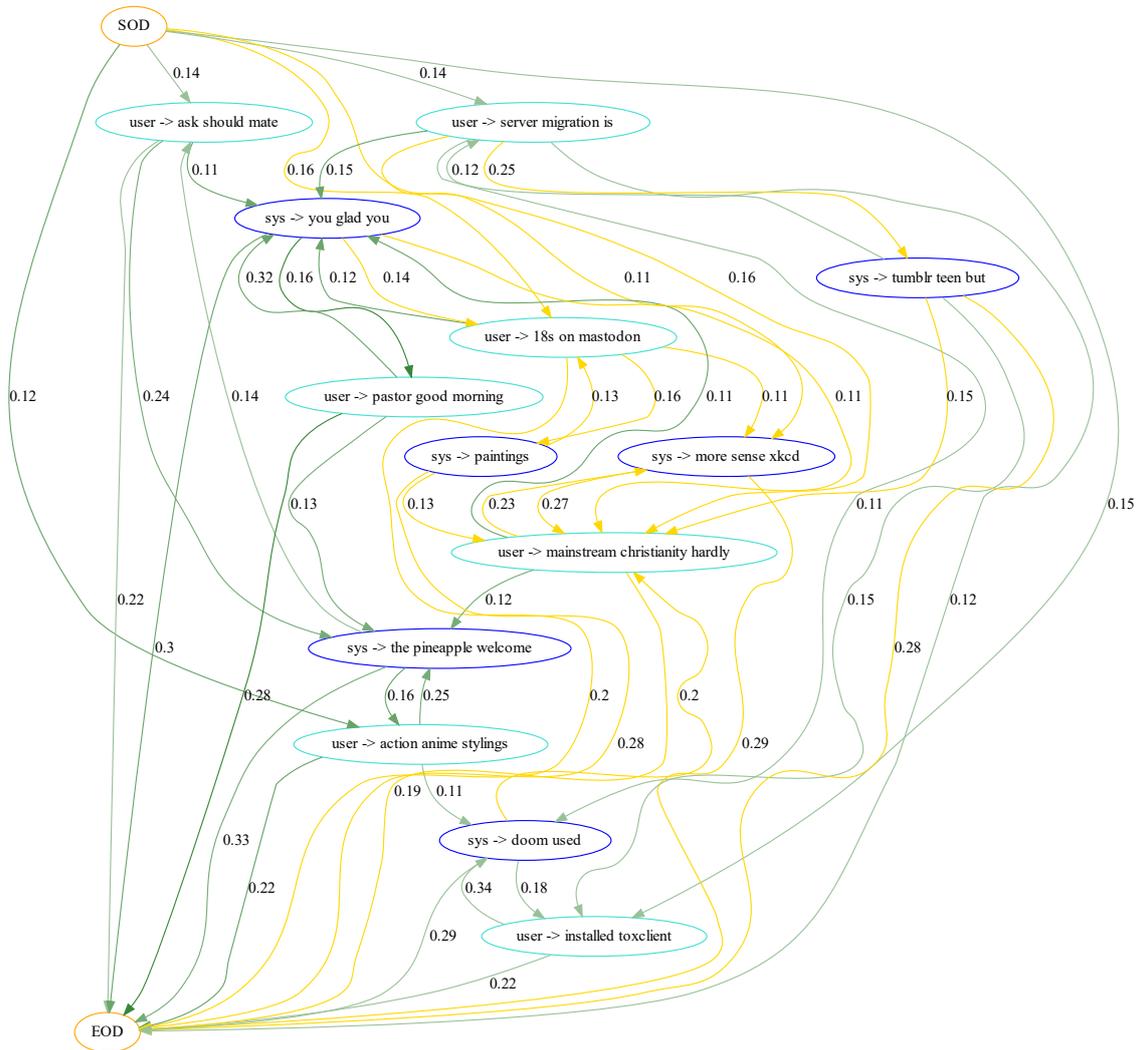


Figure 13: Sentiment-aware dialogue flow discovered for Mastodon data

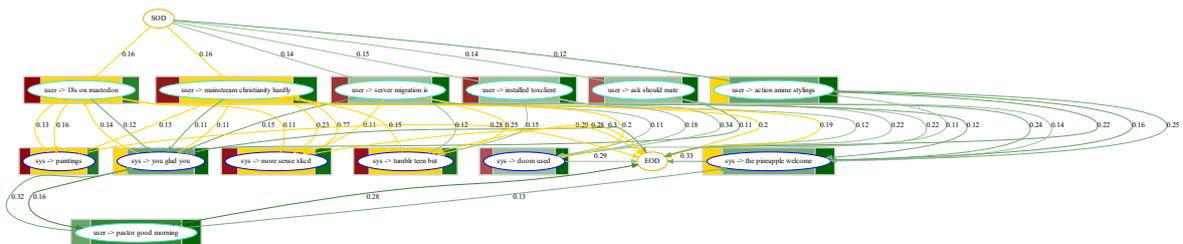


Figure 14: Sentiment-aware dialogue flow with standard deviation discovered for Mastodon data

# Analyzing and Enhancing Clarification Strategies for Ambiguous References in Consumer Service Interactions

Changling Li<sup>2\*</sup> Yujian Gan<sup>1\*</sup> Zhenrong Yang<sup>3</sup> Youyang Chen<sup>3</sup>  
Xinxuan Qiu<sup>2</sup> Yanni Lin<sup>2</sup> Matthew Purver<sup>1,4</sup> Massimo Poesio<sup>1,5</sup>

<sup>1</sup>Queen Mary University of London, UK <sup>2</sup>Guangxi Normal University, China  
<sup>3</sup>Guangxi University of Finance and Economics, China <sup>4</sup>Jožef Stefan Institute, Slovenia  
<sup>5</sup>University of Utrecht, Netherlands

Correspondence: [m.poesio@qmul.ac.uk](mailto:m.poesio@qmul.ac.uk)

## Abstract

When customers present ambiguous references, service staff typically need to clarify the customers' specific intentions. To advance research in this area, we collected 1,000 real-world consumer dialogues with ambiguous references. This dataset will be used for subsequent studies to identify ambiguous references and generate responses. Our analysis of the dataset revealed common strategies employed by service staff, including directly asking clarification questions (CQ) and listing possible options before asking a clarification question (LCQ). However, we found that merely using CQ often fails to fully satisfy customers. In contrast, using LCQ, as well as recommending specific products after listing possible options, proved more effective in resolving ambiguous references and enhancing customer satisfaction.<sup>1</sup>

## 1 Introduction

Clarification questions (CQ) have long been a focal point in dialogue research due to their various functions, with resolving ambiguities being one of the most crucial (Purver, 2004a; Boni and Manandhar, 2005; Ginzburg, 2012; Liu et al., 2014; Dhole, 2020; Lautraite et al., 2021; Testoni and Fernández, 2024). Previous studies have primarily examined whether models are capable of generating suitable clarification requests in response to ambiguities (Purver et al., 2001; Zhang and Choi, 2023; Deng et al., 2023). However, little attention has been paid to determining the most effective type of clarification request (Liu et al., 2014; Zhang and Choi, 2023). This gap in research prompts a significant question: What type of clarification request should intelligent customer service systems generate when addressing ambiguous references?

\*These two authors contributed equally to this work.

<sup>1</sup>You can find our data [here](#).

<b>Dialogue 1:</b>
A: I want the same pizza as last night.
B: What type of pizza would you like?
A: I want a Hawaiian pizza.
<b>Dialogue 2:</b>
A: I want a pizza.
C: What type of pizza would you like?
A: I want a Hawaiian pizza.

Table 1: Questions for general and specific references.

Before addressing this issue, it is necessary to clarify the definition of a CQ. Purver (2004b) defines a ‘clarification question/request’ in dialogue systems as a type of communicative action where one participant asks another to provide more information or to make their previous statement clearer. This typically occurs when the listener does not fully understand the speaker’s message due to ambiguity, vagueness, or missing information. In Dialogue 1 of Table 1, B provides an example of a CQ. However, Purver (2004b) believes that C in Dialogue 2 of Table 1 does not qualify as a CQ, but is merely an information request. This is because in Dialogue 2, C understands A’s message (a general, existentially quantified reference) but needs further information to clarify A’s needs. In Dialogue 1, B does not fully understand A’s message (a specific, definite reference), indicating ambiguity in A’s statement, and thus needs additional information to clarify. Unlike Purver (2004b), Rodríguez and Schlangen (2004) and Rieser and Moore (2005) consider both B and C to be CQs. Given that the subsequent dialogue content of B and C is similar, this paper adopts the same taxonomy and definition as Rodríguez and Schlangen (2004) and Rieser and Moore (2005).

Besides, we must define what an ambiguous reference is. An “ambiguous reference” in communication refers to a statement, word, or phrase whose

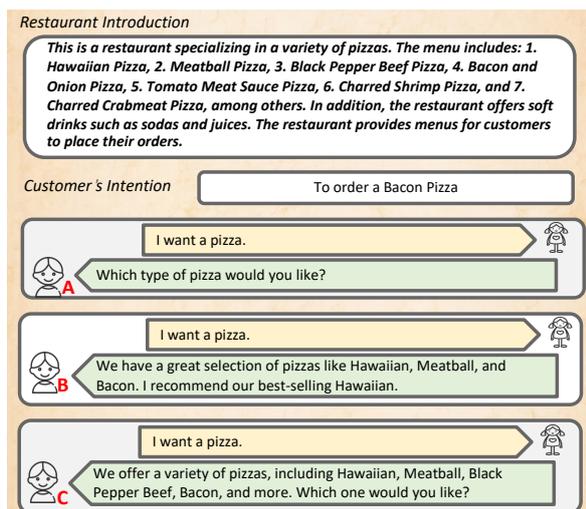


Figure 1: An example of a dialogue background with three different responses to an ambiguous reference. The girl represents the customer, and the boy represents the service staff.

meaning is not clear due to multiple possible interpretations (Eckert et al., 2003). In consumer service settings, an ambiguous reference can lead to confusion about product specifications, pricing, or customer intentions, which may hinder effective communication between the service staffs and customers. Effective clarification questions are essential in these scenarios to resolve ambiguities (Majumder et al., 2021).

Figure 1 illustrates an example of an ambiguous reference: a customer at a restaurant that serves various types of pizza orders a pizza but does not specify which kind. The figure also demonstrates several potential responses that the server might use to clarify this ambiguity. It is important to identify which response methods are commonly used by service staff and are favored by customers in practical scenarios. Addressing these questions is vital for the development of intelligent customer service systems.

To answer these questions, we collected 1,000 Chinese conversations from real-world consumer environments. Nearly every customer in these dialogues initiated at least one ambiguous reference. After organizing the data, we annotated each sentence to prepare the dataset for several uses: (1) training or evaluating a model’s capability to identify ambiguous references in conversations; (2) training or evaluating a model’s ability to resolve ambiguous references effectively through dialogue; (3) analyzing real-world service staff approaches to ambiguous references. Our analysis revealed that service staffs often use direct clarification questions

or list potential options before asking a clarification question to clear up any ambiguity, as demonstrated in responses A and C in Figure 1.

Our dataset highlights the response strategies typically used by service staff, yet these may not always align with what customers consider optimal. To gain deeper insights into customer preferences, we developed a questionnaire based on three response methods illustrated in Figure 1 and surveyed customers on their satisfaction with each response. The findings show that customers’ satisfaction levels with responses B and C are comparable and notably higher than with response A. This indicates that direct clarification questions are not the sole effective approach for addressing ambiguous references.

## 2 Dataset Construction

We gathered dialogue data from actual online and offline consumer interactions to explore how service staff addresses ambiguous references in customer inquiries. We compiled a dataset of 1,000 Chinese conversations, which were collected by four undergraduates over a period of three months, drawing on their personal shopping experiences and those of their peers.

### 2.1 Construction Principles

The dialogue dataset was constructed adhering to strict principles:

**Authenticity:** Every dialogue was directly drawn from actual consumer experiences, covering both successful and unsuccessful transactions.

**Completeness:** We ensured every conversation collected comprehensively included queries about products/services and responses from service staff.

**Diversity:** Dialogues span a range of online and offline scenarios, with offline scenarios including shops, restaurants, clothing stores, and other venues.

**Privacy Protection:** We rigorously anonymized all dialogues, removing or modifying any identifiable details, such as shop and brand names or personal identifiers.

### 2.2 Methodology for Dialogue Data Collection

Before we began data collection, we trained four data collectors to present requests with ambiguous references to service staff during their regular consumer activities, and to observe the responses. Once the transactions were complete—or if they

were terminated because the product or service was unavailable—the collectors exited the venues or ended the online sessions and reconstructed the dialogues from memory. It’s important to note that all dialogues recorded are reconstructions based on actual conversations, and any personally identifiable information has been removed.

### 2.3 Pre-Analysis of Dialogues and Consumer Scenario Classification

After gathering approximately 300 dialogue samples, we performed an initial manual summary analysis. This analysis showed that service staff respond to customers’ ambiguous requests using four main strategies: clarification questions, listing, listing followed by clarification questions, and information gathering, or they may choose to ignore the ambiguous reference. Specifically, clarification questions (*CQ*) directly address the ambiguity, as illustrated in response A of Figure 1. Listing (*LIST*) involves detailing potential options, as depicted in response B of Figure 1. Listing followed by clarification questions (*LCQ*) combines listing options with clarification questions, as seen in response C of Figure 1. Information gathering (*IG*) involves asking questions that do not directly relate to the ambiguity, such as inquiring about the customer’s preference for spicy or sweet flavors within the context of the ambiguous reference shown in Figure 1. Ignoring the ambiguous reference (*IAR*), like *IG*, overlooks the need for clarification; however, unlike *IG*, responses here are declarative rather than interrogative.

Additionally, we observed that different consumer environments may influence the responses. From the analysis of dialogue samples, we classified the consumer environments into five main categories: those with only a menu, only product displays, both a menu and product displays, neither menus nor product displays, and online shopping. The first category, labeled as ‘*MENU*’, includes scenarios found typically in restaurants where customers can see the menu but not the actual food. The second, ‘*PROD*’, refers to environments like supermarkets where only product displays are available. The third category, ‘*M&P*’, applies to fast food outlets where both menus and food are visible in display counters. The fourth, ‘*NO-M&P*’, includes service-oriented settings such as barber shops and mobile repair stores, where neither menus nor products are displayed. Lastly, the ‘*OL*’ category encompasses purely online shop-

ping. These first four categories are associated with offline consumer settings, while the last category specifically pertains to online shopping.

### 2.4 Dataset Annotation Steps

Documenting consumer dialogues is merely the initial step; they also require detailed annotation. This involves categorizing responses from service staff, briefly describing the consumer scenario as illustrated by the restaurant example in Figure 1, and identifying the type of consumer scenario. The steps for organizing and annotating this data are as follows. Step 1: Load the dialogue into a data annotation platform and record the time, city location, and specific consumer scenario, along with a concise description of it. Step 2: Meticulously annotate each sentence in the dialogue by category, including ‘*CQ*’, ‘*LIST*’, ‘*LCQ*’, ‘*IG*’, ‘*IAR*’, ‘ambiguity’ (if the customer raises an ambiguous reference), and ‘none’ (if it doesn’t fit into any of the previous categories), resulting in a total of seven categories. Step 3: Perform internal cross-validation within the team. Discuss any discrepancies in annotations during team meetings and make final decisions collaboratively.

### 3 What type of clarification question do service staff prefer to use?

Table 2 illustrates how service staff respond to requests with ambiguous references across various scenarios based on the dataset introduced in Section 2. The reason the total responses exceed the number of dialogues in the dataset is that an ambiguous reference can include multiple elements needing clarification. In most instances, service staff predominantly rely on clarification questions, including both *CQ* and *LCQ*, which constitute approximately 90% of all responses. Except in the *PROD* scenario, the frequency of using either *CQ* or *LCQ* is similar, indicating no clear preference among service staff. However, in locations where only products are displayed (*PROD*), it appears that service staff more frequently opt for *LCQ*. This approach may be necessary because similar items are not always on adjacent shelves, thus listing items from various locations helps staff better understand customer needs and guide them accurately. Furthermore, since customers lack menus and neither party may directly see the required items, *LCQ* could also improve the customer’s sensory experience.

Among the remaining response types—*LIST*,

Scenario	Total	<i>CQ</i>	<i>LIST</i>	<i>LCQ</i>	<i>IG</i>	<i>IAR</i>
<b>Offline</b>						
<i>MENU</i>	485	0.449	0.049	0.476	0.021	0.004
<i>PROD</i>	624	0.405	0.050	0.514	0.027	0.003
<i>M&amp;P</i>	404	0.505	0.010	0.483	0.002	0
<i>NO-M&amp;P</i>	58	0.483	0.017	0.448	0.051	0
<b>Online</b>						
<i>OL</i>	828	0.430	0.087	0.448	0.023	0.012

Table 2: Response Strategies by Service Staff Across Different Scenarios: Row headers distinguish online and offline scenarios as detailed in Section 2.3. The ‘Total’ column sums counts from five response strategies, each defined in the remaining column headers with explanations also in Section 2.3.

*IG*, and *IAR*—none directly involve clarification questions. *LIST* is the predominant method within these, and while it does not directly seek clarifications, it demonstrates that service staff have detected the ambiguous references in customer communications and are attempting to resolve the ambiguity in a non-questioning manner. Conversely, *IG* and *IAR* indicate a failure by service staff to accurately identify the ambiguity. Fortunately, occurrences of these latter two responses are infrequent in real-world scenarios.

#### 4 What type of clarification questions do consumers prefer to receive?

In the last section, we explored the preferred response types to ambiguous references from the perspective of service staff. This section shifts focus to customer preferences regarding the responses they receive from service staff. We conducted a hybrid online and offline survey to analyze these preferences, utilizing the Tencent Questionnaire mini-app for creation and distribution. The survey was primarily distributed in the Guangxi region of China. A total of 413 questionnaires were issued, and all were returned with valid responses.

##### 4.1 Questionnaire Design

The survey encompasses gathering basic information from participants and assessing their satisfaction with responses provided by service staff across various consumer settings. We designed 10 scenarios for this purpose, split evenly between online and offline, each offering three distinct responses from service staff for evaluation. This diverse scenario approach helps mitigate potential biases in ratings due to specific environmental or stylistic responses. The three response types assessed in-

	<i>LIST</i>		<i>CQ</i>		<i>LCQ</i>	
	Mean	Std	Mean	Std	Mean	Std
Both	4.048	0.94	3.492	1.04	4.059	0.93
Online	4.077	0.93	3.458	1.05	3.992	0.91
Offline	4.019	0.95	3.526	1.04	4.123	0.94

Table 3: Mean and Standard Deviation for Three Response Strategies: Detailed explanations of the strategies are provided in Section 2.3. ‘Both’ represent both online and offline.

clude *CQ*, *LIST*, and *LCQ*. While *LIST* is less frequently used, assessing *LIST* helps determine which aspects of the *LCQ* are most valuable to customers. Satisfaction ratings are captured on a 5-point Likert scale, ranging from 1 (strongly dislike) to 5 (strongly like), ensuring that preferences are accurately quantified. For detailed content of the questionnaire, see Appendix C.

##### 4.2 Questionnaire Data Analysis

Table 3 shows customer satisfaction rating with three distinct response types from service staff across various scenarios. A key takeaway from Table 3 is that satisfaction with mere clarification questions is the lowest, even less than the satisfaction with listing potential options, which are infrequently used by service staff. Furthermore, satisfaction levels for *LIST* and *LCQ* are similar, both substantially higher than for mere clarification. This suggests that in responses incorporating both listing and clarification, the listing aspect is deemed more crucial than clarification. Additional evidence comes from online scenarios, where satisfaction with listing alone marginally surpasses that with *LCQ*. Consequently, it is apparent that consumers prefer service staff to explicitly and exhaustively outline all options.

We next performed a one-way analysis of variance (ANOVA) on the data from Table 3 to delve deeper into the satisfaction differences across various response types in different scenarios. Initially, we analyzed dialogues from both online and offline. The analysis revealed that for comparisons between *LIST* versus *CQ* and *LCQ* versus *CQ*, the resulting *p*-values were nearly zero. This led us to reject the null hypothesis of no significant differences, demonstrating notable satisfaction disparities among these response types. In contrast, the *p*-value between *LIST* and *LCQ* was 0.85, which did not warrant rejecting the null hypothesis, indicating no significant satisfaction differences be-

tween these two types of responses. Moreover, when comparing online to offline data, the  $p$ -values for *LIST*, *CQ*, and *LCQ* were 0.03, 0.05, and 0, respectively. These findings highlight significant variations in satisfaction rating between online and offline, emphasizing the necessity for tailored customer service dialogue designs for each scenario. This implies that strategies effective in offline settings may not necessarily translate well to online interactions, and vice versa. For a more detailed analysis, Appendix B categorizes the data by age and educational level.

### 4.3 Key Takeaways

Building on the analysis, we offer the following three key insights:

- While simple clarification questions can resolve ambiguous references, they are not the most effective approach.
- Listing combined with clarification stands out as the best strategy for dealing with ambiguous references.
- Businesses can effectively resolve customer ambiguities by combining the listing of potential choices with actions like suggesting new releases, which typically maintains high levels of customer satisfaction.

## 5 Conclusion

This study analyzed service staff responses to ambiguous references using data from 1,000 customer interactions and feedback from 413 customer questionnaires. The results show that while simple clarification questions resolve ambiguities, they do not achieve high customer satisfaction. In contrast, strategies combining listing with clarification questions or others increase customer satisfaction. Future research should continue to analyze our dataset to develop more sophisticated responses that could outperform those by human service staff.

## 6 Limitation

This research is confined to a Chinese-language dialogue dataset, with the analysis restricted to surveys conducted within China. Consequently, the findings may not be directly applicable to other linguistic contexts. Furthermore, the relatively small sample of participants over the age of 65 in our questionnaire might not accurately reflect the broader opinions of this demographic.

## 7 Acknowledgements

We thank the anonymous reviewers for their helpful comments. Yujian Gan, Matthew Purver and Massimo Poesio acknowledge financial support from the UK EPSRC under grant EP/W001632/1, and Purver also from the EPSRC under grant EP/S033564/1 and from the Slovenian Research Agency (ARRS) core research programme Knowledge Technologies (P2-0103).

## References

- Saeid Amiri, Sujay Bajracharya, Cihangir Goktolga, Jesse Thomason, and Shiqi Zhang. 2019. [Augmenting knowledge through statistical, goal-oriented human-robot dialog](#).
- Negar Arabzadeh, Mahsa Seifkar, and Charles L. A. Clarke. 2022. [Unsupervised question clarity prediction through retrieved item coherency](#).
- Marco De Boni and Suresh Manandhar. 2005. [Implementing clarification dialogues in open domain question answering](#). *Natural Language Engineering*, 11:343–361.
- Yang Deng, Lizi Liao, Liang Chen, Hongru Wang, Wenqiang Lei, and Tat-Seng Chua. 2023. [Prompting and evaluating large language models for proactive dialogues: Clarification, target-guided, and non-collaboration](#).
- Kaustubh D Dhole. 2020. [Resolving intent ambiguities by retrieving discriminative clarifying questions](#).
- Claudia Eckert, Martin Stacey, and Christopher Earl. 2003. Ambiguity is a double-edged sword: similarity references in communication. In *Proceedings of the 14th international conference on engineering design*.
- Yue Feng, Hossein A. Rahmani, Aldo Lipani, and Emine Yilmaz. 2023. [Towards asking clarification questions for information seeking on task-oriented dialogues](#).
- Jonathan Ginzburg. 2012. *The Interactive Stance: Meaning for Conversation*.
- Ryan Blake Jackson and Tom Williams. 2020. [Enabling morally sensitive robotic clarification requests](#).
- Hadrien Lautreite, Nada Naji, Louis Marceau, Marc Queudot, and Eric Charton. 2021. [Multi-stage clarification in conversational ai: The case of question-answering dialogue systems](#).
- Alex Liu, Rose Sloan, Mei-Vern Then, Svetlana Stoyanchev, Julia Hirschberg, and Elizabeth Shriberg. 2014. [Detecting inappropriate clarification requests in spoken dialogue systems](#).

- Bodhisattwa Prasad Majumder, Sudha Rao, Michel Galley, and Julian Mcauley. 2021. [Ask what’s missing and what’s useful: Improving clarification question generation using global knowledge.](#)
- Matthew Marge and Alexander I Rudnicky. 2019. Miscommunication detection and recovery in situated human–robot dialogue. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 9(1):1–40.
- Matthew Purver. 2004a. The theory and use of clarification requests in dialogue.
- Matthew Purver, Jonathan Ginzburg, and Patrick G T Healey. 2001. [On the means for clarification in dialogue.](#) *Proceedings of the Second SIGdial Workshop on Discourse and Dialogue* -.
- Matthew Richard John Purver. 2004b. *The theory and use of clarification requests in dialogue.* Ph.D. thesis, University of London King’s College.
- Verena Rieser and Johanna Moore. 2005. Implications for generating clarification requests in task-oriented dialogues. In *Proceedings of the 43rd Annual Meeting of the ACL*, pages 239–246, Ann Arbor.
- Kepa Rodríguez and David Schlangen. 2004. Form, intonation and function of clarification requests in German task-oriented spoken dialogues. In *Proceedings of the 8th Workshop on the Semantics and Pragmatics of Dialogue (Catalog)*, pages 101–108, Barcelona.
- Alberto Testoni and Raquel Fernández. 2024. [Asking the right question at the right time: Human and model uncertainty guidance to ask clarification questions.](#)
- Tom Williams, Fereshta Yazdani, Prasanth Suresh, Matthias Scheutz, and Michael Beetz. 2019. Dempster-shafer theoretic resolution of referential ambiguity. *Autonomous Robots*, 43:389–414.
- J Wyatt. 2005. Planning clarification questions to resolve ambiguous references to objects. In *Proceedings of the 4th IJCAI workshop on knowledge and reasoning in practical dialogue systems, Edinburgh, Scotland*, pages 16–23.
- Michael J. Q. Zhang and Eunsol Choi. 2023. [Clarify when necessary: Resolving ambiguity through interaction with lms.](#)

## A Relate Work

### A.1 Ambiguous Reference

Researchers are interested in conversational agents facing the challenge of ambiguous reference (Eckert et al., 2003). For instance, Wyatt (2005) aims to design robots capable of engaging in task-directed conversations with humans about ambiguous references in tabletop scenes. Marge and Rudnicky (2019) presents a method for robustly handling miscommunication between people and robots in task-oriented dialogues, supporting the detection and

recovery from situated grounding problems related to referential ambiguity and impossible actions. Williams et al. (2019) initially provides recommendations for designers of robots that need to generate such requests and further demonstrates that a Dempster-Shafer reasoning component, when combined with probabilistic reference resolution, can address both pragmatic and referential uncertainties.

### A.2 Clarification Questions in Dialogues

The generation of clarification questions is vital in dialogue system research, enhancing system accuracy and user experience. Literature in this area covers various aspects: optimal timing for posing questions, task-oriented models for different scenarios, and ethical frameworks for clarification.

Arabzadeh et al. (2022) introduced an unsupervised learning method for predicting when to pose clarification questions based on retrieval item consistency and contextual similarity, showing superior generalization over neural network methods. Feng et al. (2023) developed a multi-attention sequence-to-sequence model that integrates contextual information and task knowledge to improve the specificity and accuracy of clarification questions in task-oriented dialogue systems. Further, Amiri et al. (2019) combined semantic parsing with probabilistic dialogue management to enhance knowledge base quality and human-robot interaction by generating goal-oriented clarification questions. Meanwhile, Jackson and Williams (2020) introduced a moral reasoning strategy in robots to ensure ethical responses when initiating clarification questions, integrating a moral assessment module into the robot architecture. In large language model applications, Deng et al. (2023) implemented the Proactive Chain-of-Thought (ProCoT) scheme to augment goal planning in reasoning chains, significantly improving the handling of clarification and goal-oriented questions.

## B Further Analyzing Questionnaire Data

In addition to overall lower satisfaction with *CQ* than with *LIST* and *LCQ*, *CQ* consistently ranks below both in satisfaction across all ten consumer scenarios in the questionnaire. Satisfaction levels between *LIST* and *LCQ* fluctuate. *LIST* occasionally achieves slightly higher satisfaction than *LCQ*, particularly in online scenarios, similar to the general trend in Table 3.

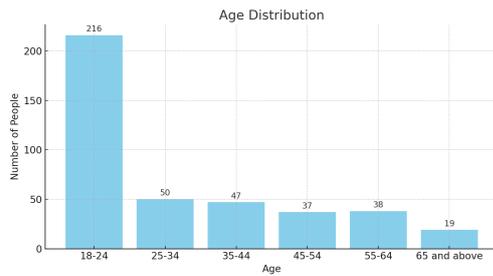


Figure 2: Age Distribution of Questionnaire Respondents

Figure 2 displays the age distribution of survey participants, predominantly ranging from 18 to 24 years old. Upon segmenting participants by age, it becomes apparent that all groups express the lowest satisfaction with *CQ*. Older participants increasingly favor *LCQ*, with satisfaction for *LIST* nearly as high. Distinctly, those aged 65 and above, while also least satisfied with *CQ*, perceive less difference between *CQ* and *LIST* responses compared to other age groups, showing almost no preference between these response types.

We performed analogous analyses based on gender and educational levels, revealing that irrespective of gender or educational classification, respondents showed a marked preference for *LIST* and *LCQ*, consistently *CQ* the lowest. Besides, females tended to rate responses higher than males across the board. Furthermore, there was a clear trend of increasing satisfaction with all three response types as educational attainment rose.

## C Questionnaire

## Questionnaire

01 Required Your age: (single choice)  
Under 18, 18-24, 25-34, 35-44, 45-54, 55-64, 65 and above

02 Required Your gender: (single choice) Male, Female

03 Required Your education level: (single choice)  
Middle school and below, High school or vocational school, Associate degree and undergraduate, Master's degree and above

04 Required Your occupation: (single choice)  
Student, Educator, IT Industry, Healthcare, Government Agency, Freelancer, Retired, Farmer, Public Institution, Enterprise, Other

Score 1 to 5: 1 - Very dissatisfied, 2 - Dissatisfied, 3 - Neutral, 4 - Satisfied, 5 - Very satisfied

05 Required: Assume you are at a breakfast shop and want to buy a char siu bao.  
You say, "I want to buy a bun." Please rate the following responses from different salespeople:

A: "What kind of bun would you like?"

B: "We have meat bun, char siu bao, red bean paste bun, and three delicacies bun. Which one would you like?"

C: "We currently have meat bun, char siu bao, red bean paste bun, and three delicacies bun."

06 Required: Assume you are at a cake shop and want to buy a cheese cake.  
You say, "I want to buy a cake." Please rate the following responses from different salespeople:

A: "Hello, we have three flavors: original, cheese, and taro. You can choose any."

B: "Hello, what kind of cake would you like?"

C: "Hello, we currently offer signature original cake, cheese cake, and taro cake. Which one would you like?"

07 Required: Assume you are at a department store and want to buy an oil-control shampoo.  
You say, "I want to buy shampoo." Please rate the following responses from different salespeople:

A: "Hello, what effect do you need from the shampoo?"

B: "Hello, we currently have shampoos with oil control, smoothing, and color protection effects. These three are very popular."

C: "Hello, we have shampoos mainly for oil control, smoothing, and color protection. Which one would you like?"

08 Required: Assume you are at a yogurt drink shop and want to buy a strawberry yogurt.  
You say, "I want to buy a cup of yogurt." Please rate the following responses from different salespeople:

A: "Our peach, strawberry, and avocado flavored yogurts are very popular. You can choose any."

B: "Our signature flavors are peach yogurt, strawberry yogurt, and avocado yogurt. Which flavor would you like?"

C: "What kind of yogurt would you like?"

The questions 9 and 10 are omitted here...

11 Required: Assume you are consulting an online customer service representative and want to buy a 20-inch suitcase.  
You say, "I want to buy a suitcase." Rate the responses from different CSRs:

A: "We have suitcases in various sizes: [18-inch Link] [20-inch Link] [22-inch Link] [24-inch Link].

Which one would you like?"

B: "Our store has 18-inch, 20-inch, 22-inch, and 24-inch suitcases.

For more details and to order, please click: [18-inch Link] [20-inch Link] [22-inch Link] [24-inch Link]."

C: "What size of suitcase would you like?"

12 Required: Assume you are consulting an online customer service representative and want to buy a double-door refrigerator.  
You say, "I want to buy a refrigerator." Rate the responses from different CSRs:

A: "Hello, we have French four-door, double-door, and T-type three-door refrigerators.

For more details and to order, please click: [French Four-door Link] [Double-door Link] [T-type Three-door Link]."

B: "Hello, we have these three types of refrigerators: [French Four-door Link] [Double-door Link] [T-type Three-door Link].

Which one would you like?"

Customer Service C: "Hello, what type of refrigerator would you like?"

13 Required: Assume you are consulting an online customer service representative and want to buy a Y brand facial cleanser.  
You say, "I want to buy a facial cleanser." Rate the responses from different CSRs:

A: "Hello, what brand of facial cleanser would you like?"

B: "Hello, our store has R brand, T brand, and Y brand facial cleansers.

For more details and to order, please click: [R Brand Link] [T Brand Link] [Y Brand Link]."

C: "Hello, our store has [R Brand Link] [T Brand Link] [Y Brand Link] facial cleansers.

Which one would you like?"

14 Required: Assume you are consulting an online customer service representative and want to buy an M brand hair dryer.  
You say, "I want to buy a hair dryer." Rate the responses from different CSRs:

A: "Hello, our store has [M Brand Link] [H Brand Link] [P Brand Link] hair dryers.

Which one would you like?"

B: "Hello, what brand of hair dryer would you like?"

C: "Hello, our store offers M brand, H brand, and P brand hair dryers.

For more details and to order, please click: [M Brand Link] [H Brand Link] [P Brand Link]."

Figure 3: Survey Questionnaire on Customer Preferences for Response Styles.

# Coherence-based Dialogue Discourse Structure Extraction using Open-Source Large Language Models

Gaetano Cimino<sup>1\*</sup>, Chuyuan Li<sup>2</sup>, Giuseppe Carenini<sup>2</sup>, Vincenzo Deufemia<sup>1</sup>

<sup>1</sup>University of Salerno, 84084, Fisciano, Salerno, Italy

<sup>2</sup>University of British Columbia, V6T 1Z4, Vancouver, BC, Canada

{gcimino, deufemia}@unisa.it

chuyuan.li@ubc.ca, carenini@cs.ubc.ca

## Abstract

Despite the challenges posed by data sparsity in discourse parsing for dialogues, unsupervised methods have been underexplored. Leveraging recent advances in Large Language Models (LLMs), in this paper we investigate an unsupervised coherence-based method to build discourse structures for multi-party dialogues using open-source LLMs fine-tuned on conversational data. Specifically, we propose two algorithms that extract dialogue structures by identifying their most coherent sub-dialogues: DS-DP employs a dynamic programming strategy, while DS-FLOW applies a greedy approach. Evaluation on the STAC corpus demonstrates a micro-F<sub>1</sub> score of 58.1%, surpassing prior unsupervised methods. Furthermore, on a cleaned subset of the Molweni corpus, the proposed method achieves a micro-F<sub>1</sub> score of 74.7%, highlighting its effectiveness across different corpora.

## 1 Introduction

Understanding multi-party dialogue structure is crucial for various natural language tasks like dialogue comprehension, summarization, and sentiment analysis (Joty et al., 2019; Li et al., 2019; He et al., 2021; Feng et al., 2022). The goal is to extract a coherent discourse structure from a dialogue transcript, wherein pairs of clause-like texts are linked through rhetorical relations. To obtain a good understanding of the coherent discourse structures in dialogues, the Segmented Discourse Representation Theory (SDRT) (Asher and Lascarides, 2003) framework proposes to annotate dialogues with dependency graphs, where edges link text spans labeled with semantic-pragmatic relations. An example dialogue derived from the Strategic Conversations corpus (STAC) (Asher et al., 2016) corpus is shown in Figure 1, with

\* This work was done during a visit to the University of British Columbia.

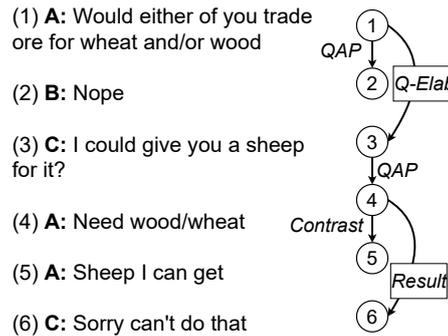


Figure 1: A dialogue instance from the STAC corpus (id *pilot04\_6*), illustrating user utterances on the left and the corresponding ground-truth dialogue structure and relations on the right. The graph reveals three distinct sub-dialogues: (1, 2), (1, 3, 4, 5), and (1, 3, 4, 6).

nodes denoting discourse units and edges relation types (i.e., *Question-Answer Pair (QAP)*, *Question-Elaboration (Q-Elab)*, *Contrast*, and *Result*).

Multi-party dialogues pose greater challenges compared to two-party dialogues, due to the involvement of numerous speakers, each contributing uniquely with more speech turn interactions and structural particularities (Asher et al., 2016). Nevertheless, this complexity allows for the segmentation of dialogues into independent conversational flows that share a common overarching topic. These conversational flows exhibit distinct internal progression and structure, thereby permitting them to be regarded as sub-dialogues (Fernández et al., 2008; Frampton et al., 2009; Sun et al., 2016). As a result, the discourse structure of multi-party dialogues can be predicted by decomposing dialogues into coherent sub-dialogues, where each sub-dialogue reflects the flow of conversation, starting with an initial utterance and concluding when no further elaboration occurs. For instance, the dialogue in Figure 1 comprises three sub-dialogues: (1, 2), (1, 3, 4, 5), and (1, 3, 4, 6). However, exploring all possible sub-dialogues to identify the coherent ones is unrealistic because it involves ana-

lyzing all possible ordered sequences of utterances within a dialogue, leading to exponential growth that makes exhaustive analysis impractical.

Supervised evaluation of dialogues is challenging due to data sparsity (Li et al., 2022). To address this issue, some studies have proposed unsupervised (Li et al., 2023) and semi-supervised (Badene et al., 2019b,a; Nishida and Matsumoto, 2022; Li et al., 2024a) methods. These methods typically predict the best discourse pair given a discourse unit, while overlooking the previous context. However, we advocate that identifying sub-dialogues can offer a broader context to better understand the thematic coherence within a dialogue, thus building a more accurate discourse structure.

In this paper, we propose an unsupervised, sub-dialogue-oriented method for extracting “naked” discourse structures without discourse relations in multi-party dialogues. Although without relations, discourse structures alone have been shown to be crucial features for tasks such as content selection (Louis et al., 2010) and summarization (Xiao et al., 2020; Xu et al., 2020). Precisely, we introduce two algorithms: Multi-Party Dialogue Structure Extraction based on Dynamic Programming (DS-DP) and Multi-Party Dialogue Structure Extraction based on Flow Conversation Analysis (DS-FLOW), designed to decompose dialogues into coherent sub-dialogues. DS-DP identifies the most coherent (partial) sub-dialogues ending in each discourse unit by applying a dynamic programming strategy. In contrast, DS-FLOW greedily predicts for each discourse unit the most likely coherent subsequent utterances, followed by a process that ensures the completeness of the resulting discourse structure. In both algorithms, we use perplexity as a metric to evaluate sub-dialogue coherence. To compute perplexity scores, we draw inspiration from work on Pre-trained Language Models (PLMs) and Large Language Models (LLMs) fine-tuned on conversational data implicitly capturing dialogue quality (Mehri and Eskénazi, 2020; Bruyn et al., 2022).

We utilize open-source models, as proprietary models are limited to text-based prompts and do not permit analysis of output probabilities. In practice, we compare the performance of two open-source LLMs: a chatbot trained by fine-tuning LLaMA on user-shared conversations Vicuna-13B (Chiang et al., 2023) and a general-purpose model Mistral-7B (Jiang et al., 2023). We evaluate our method on the STAC corpus (Asher et al., 2016) and a revised subset of the Molweni corpus (Li et al., 2020). The

results demonstrate the effectiveness of our solution, as it outperforms prior unsupervised methods. Specifically, we achieve a micro-F<sub>1</sub> score of 58.1% on STAC and 74.7% on Molweni, demonstrating its robustness across different corpora.

The contributions of this paper are twofold. First, we propose a fully unsupervised method for extracting graph structures of multi-party dialogues, which is the first of its kind to the best of our knowledge. Second, we introduce and evaluate two novel algorithms that leverage open-source LLMs to decompose dialogues into coherent sub-dialogues, enabling a more fine-grained analysis of discourse structures.

## 2 Related Work

**Multi-Party Dialogue Discourse Parsing** Various methodologies have been proposed for parsing multi-party dialogues. Perret et al. (2016) developed an Integer Linear Programming approach predicting non-tree structures by encoding linguistic principles as constraints. Wang et al. (2021) presented the Structure Self-Aware model, using an edge-centric graph neural network to learn representations of discourse unit pairs directly. Bennis et al. (2023) introduced BERTLine, a discourse parsing model leveraging a multi-task setup to jointly predict discourse attachments and relation labels, achieving state-of-the-art performance. Mao et al. (2024) proposed the Hierarchical Graph Fusion Network, using hierarchical graph neural networks to encode contextual levels like utterances, dialogue topics, and user preferences. While effective, these approaches rely on annotated data, posing challenges due to limited resources. To address data sparsity, recent studies have explored unsupervised and semi-supervised strategies using PLMs and LLMs. For instance, Li et al. (2023) proposed extracting dependency trees from PLM attention matrices using unsupervised metrics or semi-supervised strategies with small validation sets. Instead, Li et al. (2024a) designed a semi-supervised pipeline to predict structures and relations sequentially via self-training. In another study, Chan et al. (2023) used zero- and few-shot prompting techniques to assess ChatGPT on discourse parsing, but achieved abysmal results. In contrast, the method proposed in this paper requires no annotation, using fully unsupervised approaches to extract discourse structures in multi-party dialogues. Furthermore, unlike the unsupervised ap-

proach proposed by Li et al. (2023), our solution can also extract graph structures rather than being limited to dependency trees.

**LLMs for Dialogue Evaluation** Prior research has highlighted the inherent ability of PLMs and LLMs to implicitly capture dialogue quality, making them suitable for evaluating dialogues. Mehri and Eskénazi (2020) and Bruyn et al. (2022) introduced the FED and FULL metrics, respectively, to assess open-domain dialogue systems utilizing PLMs and LLMs without requiring ground-truth responses or supervised training data. These metrics evaluate dialogue quality by estimating the likelihood of a model generating follow-up utterances aligned with different dimensions of dialogue quality after a given system response. The strong correlation observed between metric scores and human judgments suggests that PLMs and LLMs have acquired meaningful representations of dialogue quality aligned with human perceptions. Similarly, Zhang et al. (2024) analyzed LLMs as automatic dialogue evaluators, inspired by the remarkable performance of LLMs fine-tuned using the instruction-tuning approach (Zhang et al., 2023). Their study involved multidimensional evaluation of proprietary and open-source LLMs for assessing dialogue quality across various dimensions. Results indicate that appropriately aligned and utilized LLMs can effectively serve as generalized automatic dialogue evaluators, complementing human judgments. Motivated by these findings, in this paper we evaluate sub-dialogues by leveraging LLMs’ capabilities in generating coherent dialogues and adhering to relevant instructions.

**Sub-Dialogue Detection** Sub-dialogues are extensively studied in computational tasks, notably within Dialogue State Tracking (DST) (Sun et al., 2016; Lee et al., 2021), aiming to understand decisions in multi-party conversations. In this framework, a dialogue session is deconstructed into a series of sub-dialogues, each consisting of consecutive multi-turn exchanges focused on a shared topic. Departing from conventional DST approaches, in this paper we adopt sub-dialogue delineation to extract the underlying structure of multi-party dialogues. Specifically, our method involves unsupervised evaluation of multi-party dialogue discourse units, leveraging insights from LLMs fine-tuned on conversational data.

### 3 Method

In this section, we first formally describe the dialogue parsing task. We then describe two algorithms for sub-dialogue extraction. The first relies on a dynamic programming strategy and is formally denoted as Multi-Party Dialogue Structure Extraction based on Dynamic Programming (DS-DP). Conversely, the second algorithm, grounded in the analysis of conversation flows, is formally named Multi-Party Dialogue Structure Extraction based on Flow Conversation Analysis (DS-FLOW).

#### 3.1 Problem Formulation

Let  $D = (e_1, e_2, \dots, e_n)$  be a dialogue consisting of  $n$  *Elementary Discourse Units* (EDUs), each representing the smallest unit of discourse. In the SDRT framework, a dialogue  $D$  can be represented as a *Directed Acyclic Graph* (DAG), denoted as  $DAG(D)$ , wherein EDUs are connected with directed edges. Dialogue discourse parsing aims to automatically derive a DAG that best represent the SDRT structure of a dialogue. In our proposal, the construction of  $DAG(D)$  involves creating  $m$  sub-dialogues  $Sub_1(D), \dots, Sub_m(D)$ , where each sub-dialogue possesses an intrinsic structure  $S_{Sub_j(D)} \subseteq DAG(D)$ . Thus, discourse structure extraction can be reframed as linking EDUs within  $D$  to form the most coherent sub-dialogues, such that  $DAG(D) = \bigcup_{j=1}^m S_{Sub_j(D)}$ . Note that the following properties hold:

- We assume the absence of backward links in the final DAG as an utterance cannot depend, either anaphorically or rhetorically, on subsequent utterances within a dialogue, as they are previously unknown (Afantenos et al., 2012; Li et al., 2023).
- Each sub-dialogue must include an edge originating from the initial EDU  $e_1$ . This means that:

$$\forall j \in \{1, \dots, m\} \exists (e_1, e_k) \in S_{Sub_j(D)}$$

This constraint is justified by the fact that all EDUs, except  $e_1$ , must have at least one incoming edge from a previous node, and recursively following these edges backward leads to  $e_1$ . In SDRT, the DAGs have unique roots, so that every single EDU is reachable from the first EDU, i.e., the axiom (Perret et al., 2016).

- Sub-dialogues can overlap, allowing certain edges to be part of multiple sub-dialogues, justified by the fact that speaker interventions may

contribute to different themes within one dialogue. For instance, the edge  $(e_3, e_4)$  in Figure 1 appears in two sub-dialogues. In sub-dialogue  $(e_1, e_3, e_4, e_5)$ , this edge leads to speaker A’s elaboration in  $e_5$  on their inquiry in  $e_4$ , which was prompted by speaker C’s question in  $e_3$ . In sub-dialogue  $(e_1, e_3, e_4, e_6)$ , the edge  $(e_3, e_4)$  leads to speaker C’s declination in  $e_6$  of speaker A’s inquiry in  $e_4$ .

### 3.2 DS-DP Algorithm

This algorithm uses dynamic programming to efficiently explore the space of all possible sub-dialogues within a dialogue. As a first step, given a dialogue  $D$  as input, the algorithm maps it into a fully-connected graph with only forward links  $G = (V, E)$ . In this graph,  $V$  represents the set of EDUs within the dialogue, and  $E$  includes all potential links in the dialogue’s structure. The DS-DP algorithm aims to extract from  $G$  the paths corresponding to the most coherent (partial) sub-dialogues starting from the initial EDU and ending in each subsequent EDU, based on a coherence metric denoted as *eval*. To this end, it defines two matrices of size  $(|V| - 1) \times (|V| - 1) \times (|V| - 2)$ , which we call  $M_{co}$  and  $M_{pred}$ . Here,  $M_{co}[i][j][k]$  denotes the maximum coherence of a sub-dialogue passing through node  $i$ , ending in node  $j$ , with  $k$  preceding nodes before node  $i$ . Similarly,  $M_{pred}[i][j][k]$  stores the previous node to achieve the maximum coherence value of the sub-dialogue ending in node  $j$ , passing through node  $i$ , and considering  $k$  preceding nodes before node  $i$ . Taking the unidirection property into account, only the upper right half of the matrices contain valid values; no values are stored in the lower left part of the matrices. For initialization, the assignment

$$\forall j \quad M_{co}[0][j][0] = eval(e_1, e_j)$$

is set, rooted in the recognition that the only sub-dialogues without preceding nodes are those progressing from the initial node to any subsequent node. As a result, for each EDU  $e_i$  ( $i > 1$ ), the algorithm computes the most coherent sub-dialogues starting from  $e_1$  and ending in  $e_i$ , with  $k$  intermediate nodes ( $k \in [1, i - 1]$ ). Specifically, each EDU  $e_i$  may either directly connect to  $e_1$  or include up to  $i - 1$  edges within its most coherent sub-dialogue.

The pseudo-code of the DS-DP algorithm for matrix construction is presented in Algorithm 1. It iterates through each  $k$  value within the range from 1 to  $|V| - 2$ . For each  $k$ , it systematically traverses each

---

#### Algorithm 1 DS-DP - Matrix Construction

---

**Input:**  $G = (V, E)$   
**Output:** Updated  $M_{co}$  matrix and  $M_{pred}$  matrix

```

1: for  $k \leftarrow 1$  to  $|V| - 2$  do
2:   for  $i$  in  $V$  do
3:     for  $j$  in  $V$  do
4:       if  $j > i$  then
5:         for each node  $u$  with an edge into  $i$  do
6:           if  $M_{co}[u][i][k - 1] \neq \text{NULL}$  then
7:              $val \leftarrow eval(k - 1 \text{ EDUs}, u, i, j)$ 
8:             if  $val$  better than  $M_{co}[i][j][k]$  then
9:                $M_{co}[i][j][k] \leftarrow val$ 
10:               $M_{pred}[i][j][k] \leftarrow u$ 
11: end of all loops and conditions
```

---

node  $i$  according to the topological order defined on  $G$ . Subsequently, for each node  $i$ , it explores all possible successor nodes  $j$ . During this traversal, it examines each node  $u$  that has an edge directed towards  $i$ . The condition  $M_{co}[u][i][k - 1] \neq \text{NULL}$  indicates that node  $u$  has been previously visited, implying the feasibility of reaching node  $i$  from  $u$  by considering  $k - 1$  preceding nodes along the path. This condition ensures the consideration of only those nodes  $u$  that are accessible from  $i$  and can therefore serve as intermediary nodes to reach  $j$  with  $k - 1$  previous nodes. Upon satisfying this condition, the algorithm evaluates the coherence of the sub-dialogue ending at  $j$ , including  $k - 1$  preceding nodes,  $u$ , and  $i$ . If this assessment yields a coherence value superior to the one currently stored in  $M_{co}[i][j][k]$ , the matrix is updated with the new coherence value, and the predecessor information is recorded in  $M_{pred}[i][j][k]$ . To identify a sub-dialogue starting from the initial node and ending in a specified node  $e_j$ , it is sufficient to examine all non-null entries in the  $M_{co}$  matrix while keeping  $j$  constant. Subsequently, the sub-dialogue characterized by the highest coherence value is considered. The final DAG is then constructed by combining the identified sub-dialogues. An illustration of the application of DS-DP to the dialogue depicted in Figure 1 is provided in Appendix E.

From a complexity analysis perspective, the DS-DP algorithm comprises four nested loops for matrix construction and two nested loops for structure prediction. The first three outermost loops in Algorithm 1 iterate  $O(|V|)$  times each, resulting in a time complexity of  $O(|V|^3)$ . The last innermost loop processes all incoming edges of the current node, which has a time complexity of  $O(|V|)$ . For structure prediction, each node  $j$  requires iteration over  $j - 1$  values (since each path ending in node  $j$  can have a maximum length of  $j - 1$ ), resulting in

a time complexity of  $O(|V|^2)$ . Consequently, the overall time complexity of DS-DP is dominated by the matrix construction process, which has a worst-case time complexity of  $O(|V|^4)$ .

**Coherence Evaluation** Coherence, defined by the seamless flow and logical progression inherent in conversational interactions, stands as a pivotal criterion for assessing dialogues. Within the text analysis context, *perplexity* emerges as a valuable metric for evaluating the coherence of textual constructs (Colla et al., 2022). Consequently, we adopt perplexity as the *eval* metric to quantitatively measure how effectively a sub-dialogue maintains its logical structure and natural progression. Drawing from earlier studies indicating that LLMs capture elements of dialogue quality (Mehri and Eskénazi, 2020; Bruyn et al., 2022), we employ them to estimate the joint probability of each sub-dialogue  $Sub_D = (e_1, \dots, e_l)$  of a dialogue  $D$ . The perplexity score is calculated as

$$Pe(Sub_D) = \exp\left(-\frac{1}{l} \sum_{i=1}^l \log P(e_i | e_{<i})\right)$$

and provides insights into the model’s level of certainty or uncertainty in predicting the unfolding discourse. Lower perplexity scores indicate higher coherence, demonstrating the model’s proficiency in comprehending the logical flow of conversation.

### 3.3 DS-FLOW Algorithm

While DS-DP constructs sub-dialogues by identifying the most likely antecedents of a given EDU, DS-FLOW mainly focuses on capturing the most fluent successive utterances of a given EDU. Specifically, it consists of three steps: (i) In the first step, for each EDU excluding the final one, the algorithm predicts the most probable subsequent EDU that elaborates upon it. Notably, previous incoming links to EDUs are utilized to inform these predictions. We evaluate sub-dialogue coherence using the perplexity metric, as discussed in DS-DP. (ii) In the second step, a filtering mechanism is applied to recognize the conclusion of conversational segments. This step addresses the issue that not all utterances are elaborated upon further, resulting in certain nodes lacking outgoing links. (iii) The third step involves a backward analysis to address potential *orphan* EDUs (i.e., EDUs without incoming edges) due to the filtering process or the lack of links predicted in the first step. For each orphan EDU  $e_i$ , the analysis selects a parent out of all

sub-dialogues ending in an EDU  $e_j$  where  $i > j$ .

An illustration of the application of DS-FLOW to the dialogue depicted in Figure 1 is provided in Figure 2. It elucidates the following steps: the initial identification of an outgoing link for each EDU, the subsequent filtration of links ( $e_2, e_4$ ) and ( $e_5, e_6$ ), and the selection of sub-dialogues ( $e_1, e_3$ ) and ( $e_1, e_3, e_4, e_5$ ) in the backward analysis due to the absence of incoming links for  $e_3$  and  $e_5$  in the first two steps, thereby augmenting the final DAG with edges ( $e_1, e_3$ ) and ( $e_4, e_5$ ).

From a complexity analysis perspective, the DS-FLOW algorithm constructs a DAG with  $|V| = n$  nodes from a dialogue containing  $n$  EDUs through three sequential steps. Initially, it predicts the subsequent EDU for each dialogue segment by leveraging prior connections, achieving a linear time complexity of  $O(|V|)$ . Following this, it filters segments that terminate without additional elaboration, also operating in linear time  $O(|V|)$ . Subsequently, in its third step, DS-FLOW undertakes a backward analysis to assign appropriate parents to orphan EDUs from previously identified sub-dialogues. In the worst-case scenario, this involves evaluating each orphan against all preceding EDUs, resulting in a time complexity of  $O(|V|^2)$ . Consequently, the overall time complexity of DS-FLOW is  $O(|V|^2)$ .

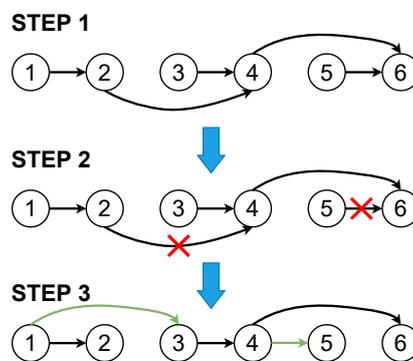


Figure 2: An example of DS-FLOW execution.

**Filtering Mechanism** An approach to implementing the filtering mechanism entails employing instruction-tuned LLMs as automatic dialogue evaluators (Zhang et al., 2024), prompting these models to generate responses for potential dialogue continuations. Specifically, given a pair of EDUs, the LLM is tasked with evaluating whether the second EDU (i.e., the next sequential EDU in the dialogue) builds upon the first one. However, as noted by Zhang et al. (2024) and confirmed through our own experimentation, the text generated by open-

source LLMs can become problematic, featuring content that is nonsensical or inaccurate (Rawte et al., 2023). Consequently, we follow the method outlined by Gupta et al. (2022), employing an implicit scoring mechanism. Specifically, when presented with an instruction prompt input<sup>1</sup>, we focus on the output probabilities associated with the words “Yes” and “No” as generated by the LLM. In this context, we compute the probability

$$P(e_i \rightarrow e_j) = \frac{P(\text{Finaltoken}=\text{“Yes”})}{P(\text{Finaltoken}=\text{“Yes”})+P(\text{Finaltoken}=\text{“No”})}$$

where we evaluate the likelihood of having an EDU  $e_j$  as a subsequent utterance following the EDU  $e_i$ . In the DS-FLOW algorithm, we discard outgoing links with a probability lower than 0.5. The evaluation of the filtering mechanism’s performance revealed that it filters a limited number of links with commendable reliability. Additional details are provided in Appendix C.

### 3.4 Additional Constraints

The proposed approaches for analyzing sub-dialogues within a dialogue face a challenge of preserving semantic coherence. Specifically, certain sub-dialogues may lack coherence, such as examining the link between the first and last EDUs in a long dialogue, which is unlikely to constitute a valid connection. To illustrate this challenge, consider the following dialogue excerpt:

( $e_1$ ) **A**: Did you enjoy the movie last night?  
 ( $e_2$ ) **B**: Yeah, the plot twist was unexpected.  
 . . .  
 ( $e_{p-1}$ ) **A**: What did you think about the ending?  
 ( $e_p$ ) **B**: Oh, it was great!

where  $p$  is a large number. In this scenario, a valid link exists between  $e_1$  and  $e_2$ . However, when examining individual pairs of EDUs,  $e_p$  may erroneously be deemed as coherent with  $e_2$  in relation to  $e_1$ , despite their temporal separation and semantic incongruity within the ongoing conversation. To mitigate the issue of incoherent sub-dialogues, we advocate for including a hard constraint on the distance between two EDUs under scrutiny. As done by Bennis et al. (2023), when assessing the potential linkage between an EDU  $e_j$  and one of its preceding  $e_i$ , we impose the condition  $j > i \geq j - 10$ . By analyzing the development sets from the STAC and Molweni corpora, we observed that fewer than

1.9% of the links fail to meet the specified condition. By limiting the distance between EDUs, we reduce computational complexity and enhance the likelihood of extracting relevant information from nearby EDUs, thereby improving the coherence of sub-dialogues. Additionally, we propose integrating a penalization factor  $P_{dist}(d)$ , where  $d$  represents the number of intervening speech turns between two EDUs. This factor increases the perplexity associated with a sub-dialogue as the temporal distance between the two EDUs to be linked increases. By prioritizing proximity between EDUs, the incorporation of the penalization factor aims to account for the potential degradation of coherence over time. Specifically, we adopt  $\sqrt{d}$  as the penalization factor for perplexity scores. This penalty is applied by multiplying the perplexity score of a sub-dialogue by the output of  $P_{dist}(d)$ .

## 4 Experimental Setup

**Corpora** We conduct experiments on two commonly used SDRT-annotated dialogue corpora: (*i*) **STAC** (Asher et al., 2016). This corpus contains 1161 multi-party dialogues arising from interactions within an online version of the game “The Settlers of Catan”. Given the unsupervised nature of our method, we evaluate it on the test set, which consists of 109 documents, amounting to 1129 EDU pairs. (*ii*) **Molweni** (Li et al., 2020). Derived from the Ubuntu Chat Corpus (Lowe et al., 2015), this corpus centers around technical discussions concerning the Ubuntu system. Due to quality issues with the original annotations (Li et al., 2023), we employ the “Molweni-clean” version proposed by Li et al. (2024a), which consists of 50 documents, encompassing 373 EDU pairs. Detailed corpus statistics are presented in Table 1.

**Evaluation Metrics** To assess the performance of the proposed approaches, we report the micro- $F_1$ , recall, and precision for the generated structures.

**Compared Methods** We contrast our method with the straightforward yet strong unsupervised LAST baseline (Schegloff, 2007), which links each EDU with the preceding one. Moreover, we compare it with the method proposed by Li et al. (2023), currently the only known unsupervised approach in the literature proficient at predicting discourse structure, albeit without explicitly extracting DAGs. Finally, to draw insights from modern LLMs, we present results from ChatGPT (*gpt-3.5-turbo ver-*

<sup>1</sup>The prompt is detailed in Appendix B.

Corpus	#Doc	#Turn/doc	#Tok/doc	#Spk/doc
STAC	109	10.6	42.5	3.0
Molweni-clean	50	8.5	91.1	3.2

Table 1: Key statistics of corpora: number of documents (#Doc), averaged speech turns, tokens, and speakers per document (#Turn/doc, #Tok/doc, #Spk/doc).

sion), Vicuna, and Mistral in a zero-shot setting.

**Implementation Details** We use the Vicuna-13b and Mistral-7b models from the Hugging-Face library (Wolf et al., 2020). We employ the *lm-evaluation-harness*<sup>2</sup> framework for computing perplexity scores. We replace speaker names with markers (e.g., John → “spk1”) to match the inference setup in the employed models.

## 5 Results and Analysis

### 5.1 DS-DP and DS-FLOW Performance

Table 2 shows the performance of the DS-DP and DS-FLOW algorithms on the STAC and Molweni-clean corpora. Precisely, the results for each model include the vanilla version, as well as versions incorporating the penalization factor ( $P_{dist}(d)$ ) and the speech turn limitation (STL). Generally, algorithms utilizing vanilla models perform worse compared to those with constraints; however, they show potential in predicting distant links, as discussed in the following Section 5.3. Applying the STL constraint consistently enhances performance across all metrics. For instance, DS-FLOW on STAC shows an increase in the micro-F<sub>1</sub> score for Vicuna (from 47.2% to 47.7%) and Mistral (from 46.2% to 46.7%). Similarly, DS-DP on STAC improves for Vicuna (from 54.3% to 54.4%) and Mistral (from 53.8% to 54.8%). Comparable improvements are observed on Molweni-clean. These findings suggest that while the STL constraint yields marginal improvements, it reduces complexity by limiting the analysis to fewer sub-dialogues, facilitating a cohesive sub-dialogue examination. Despite predicting complex links with vanilla LLMs and the STL constraint, temporal disparity lowers precision scores (see Section 3.4). When applying the penalization factor  $P_{dist}(d)$ , significant improvements are noted, as shown in the third row of each group in Table 2. The factor  $P_{dist}(d)$  improves results by discouraging longer-distance links and favoring shorter ones, which are more prevalent, as discussed in Section 5.3. Consequently, the best

<sup>2</sup><https://github.com/EleutherAI/lm-evaluation-harness>

performance on STAC is achieved with the DS-FLOW algorithm using STL and  $P_{dist}(d)$ . Similarly, the optimal performance on Molweni-clean is obtained with the DS-DP algorithm using STL and  $P_{dist}(d)$ .

Leveraging the dynamic programming strategy, DS-DP analyzes a larger number of sub-dialogues compared to the greedy approach employed by DS-FLOW, tending to select more short links. This is highlighted by the best performance on Molweni-clean, which involves fewer long-distance links compared to STAC. Conversely, DS-FLOW better predicts longer-distance links, achieving the best performance on the STAC corpus. Overall, when comparing average micro-F<sub>1</sub> scores of DS-DP and DS-FLOW under optimal settings across both corpora, DS-DP slightly outperforms DS-FLOW with scores of 66% versus 65.5%, respectively<sup>3</sup>.

Regarding backbone LLMs, Vicuna consistently outperforms Mistral across all settings, highlighting the advantage of models fine-tuned on conversational data for dialogue analysis tasks. However, Mistral demonstrates satisfactory performance, validating the efficacy of the proposed algorithms.

### 5.2 Unsupervised Method Comparison

We compare our top-performing DS-DP and DS-FLOW settings with other unsupervised methods. Precisely, we consider the following benchmarks<sup>4</sup>: (i) LAST baseline predicts local attachments between adjacent EDUs. Despite its high performance on STAC and Molweni (Muller et al., 2012), it can only extract a single sub-dialogue and cannot detect sub-dialogue structures like our method. (ii) The unsupervised method by Li et al. (2023), which extracts dependency trees from BART model attention matrices (Lewis et al., 2020), fine-tuned through the Sentence Ordering (SO) task. (iii) ChatGPT in a zero-shot setting with a novel prompt for multi-party dialogue discourse parsing, achieving a micro-F<sub>1</sub> score of 52% on STAC, significantly improving over the 20.5% reported by Chan et al. (2023). The prompt is detailed in Appendix H. (iv) Vicuna-13b and Mistral-7b models, prompted identically to ChatGPT in a zero-shot setting.

Table 3 shows the comparison results. Using Vicuna-13b, the DS-DP and DS-FLOW algorithms excel on the STAC corpus, achieving micro-F<sub>1</sub>

<sup>3</sup>Qualitative analysis of generated structures is presented in Appendix G.

<sup>4</sup>See Appendix D for additional results pertaining to a smaller Vicuna model.

Model	Algorithm	STAC			Molweni-clean		
		F <sub>1</sub>	P	R	F <sub>1</sub>	P	R
Vicuna-13b	DS-DP	54.3	52.9	55.8	71.5	68.0	75.3
	DS-DP + STL	54.4	53.3	55.7	72.0	68.9	75.3
	DS-DP + STL + $P_{dist}(d)$	57.3	55.1	59.7	<b>74.7</b>	<b>70.1</b>	<b>79.9</b>
	DS-FLOW	47.2	40.0	57.4	58.1	48.1	73.2
	DS-FLOW + STL	47.7	40.0	59.0	59.3	49.3	74.5
	DS-FLOW + STL + $P_{dist}(d)$	<b>58.1</b>	<b>57.1</b>	59.2	72.9	69.1	77.2
Mistral-7b	DS-DP	53.8	52.2	55.5	71.1	68.8	73.7
	DS-DP + STL	54.8	53.0	56.6	71.5	68.6	74.5
	DS-DP + STL + $P_{dist}(d)$	56.7	53.4	<b>60.4</b>	74.1	69.1	<b>79.9</b>
	DS-FLOW	46.2	39.3	55.9	57.0	49.0	68.1
	DS-FLOW + STL	46.7	39.5	57.3	57.3	49.5	68.1
	DS-FLOW + STL + $P_{dist}(d)$	57.0	56.5	57.4	71.0	66.5	76.1

Table 2: Experiment results of proposed approaches on STAC and Molweni-clean corpora. STL: Speech turn limitation.  $P_{dist}(d)$ : Penalization factor. F<sub>1</sub>: Micro-F<sub>1</sub>. P: Precision. R: Recall.

Corpus	Baseline	PLM	ChatGPT	Vicuna-13b			Mistral-7b		
	LAST	BART-SO	ZS	ZS	DS-DP	DS-FLOW	ZS	DS-DP	DS-FLOW
STAC	56.8	57.2	52.0	22.8	57.3	<b>58.1</b>	30.2	56.7	57.0
Molweni-clean	76.9	-	65.6	35.2	<b>74.7</b>	72.9	36.7	74.1	71.0

Table 3: Micro-F<sub>1</sub> scores on STAC and Molweni-clean for the LAST baseline, unsupervised PLM, LLMs within a zero-shot (ZS) setting, and proposed approaches.

scores of 57.3% and 58.1%, respectively, surpassing LAST baseline and BART-SO model. It is noteworthy that the BART-SO model is previously fine-tuned with the SO task on STAC. When employing a vanilla BART model, the performance decreases to 56.6%, representing a 2.6% lower result compared to our method. In comparison, our solution does not require domain-specific data or a fine-tuning process, rendering it easily adaptable to any scenarios. Using Mistral-7b, DS-FLOW outperforms LAST but not BART-SO. On the Molweni-clean corpus, DS-DP and DS-FLOW algorithms lag behind LAST, which achieves 76.9% due to a larger amount of adjacent links in the corpus. Even the strategy proposed by Li et al. (2024a), involving cross-domain training on STAC, only attains a micro-F<sub>1</sub> score of 75.6% on the Molweni-clean corpus, thus trailing behind the LAST baseline. Consequently, a micro-F<sub>1</sub> score of 74.7% (for the DS-DP algorithm incorporating the Vicuna model) may be deemed satisfactory in a fully unsupervised setting. Owing to reproducibility challenges encountered with the BART-SO model on the Molweni-clean corpus, a comparative analysis with our algorithms is not feasible. Finally, in zero-shot settings, Vicuna and Mistral perform abysmally (from -47% to -61% compared to DS-DP and DS-FLOW). ChatGPT outperforms both open-source models, while still falling behind our

proposed unsupervised algorithms. Mistral excels at generating structured responses, while Vicuna struggles with lengthy dialogues but outperforms Mistral in our algorithms on both corpora.

### 5.3 Link Length Analysis

The LAST baseline’s limitation is its inability to predict indirect links. To assess the accuracy of our algorithms in predicting distant links, we investigate the performance concerning different link lengths. Figure 3 shows recall scores for different link lengths for DS-FLOW and DS-DP using Vicuna on STAC and Molweni-clean, respectively. We test different settings, including vanilla Vicuna and STL individually for both algorithms. For DS-FLOW on STAC, using vanilla Vicuna accurately predicts long-distance links up to distances of 12 and 13 but increases false positives, as discussed in Section 5.1, affecting precision. Adding STL (DS-FLOW+STL) improves performance for shorter links (distances 1, 2, and 3) and predicts long-distance links up to distance 10. Incorporating  $P_{dist}(d)$  with STL (DS-FLOW+STL+PF) achieves over 90% recall for direct links and maintains some ability to predict long-distance links, though performance drops for links over distance 6. Long-distance links ( $\geq 6$ ) are rare in STAC, under 5% of all links. For DS-DP on Molweni-clean, like DS-FLOW on STAC, both vanilla Vi-

cuna and STL (DS-DP+STL) predict indirect links but not those longer than 4. Including  $P_{dist}(d)$  (DS-DP+STL+PF) achieves nearly perfect recall for direct links, with slight performance drops for links at distances 2 and 3 compared to DS-DP+STL. Long-distance links ( $\geq 4$ ) are rare in Molwени-clean, under 3% of all links. The LAST baseline achieves perfect recall for direct links but fails on long-distance links. In terms of precision and  $F_1$  scores, the STL+PF setting demonstrates higher precision for short-distance links but somewhat lower precision for direct links. All settings experience a decline in  $F_1$  scores as link length increases. An exception is observed for DS-FLOW using the vanilla Vicuna on STAC, which maintains relatively high  $F_1$  scores for links at distances of 12 and 13. Further evaluation results are in Appendix A.

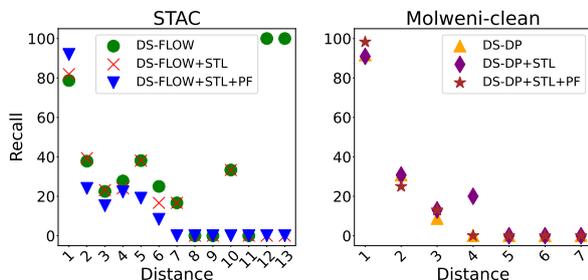


Figure 3: Recall scores for different link lengths. The left plot shows three settings with DS-FLOW on STAC; the right plot depicts the same settings for DS-DP on Molwени-clean. Both algorithms use Vicuna as LLM. STL: Speech turn limitation. STL+PF: Speech turn limitation in combination with penalization factor.

## 6 Discussion

Despite its significantly lower parameter count compared to the Vicuna-13b model, we used the Mistral-7b model for our algorithms’ assessment owing to its superior performance relative to larger models like LLaMa 2-13b (Touvron et al., 2023b) and LLaMa 1-34b (Touvron et al., 2023a) across multiple benchmarks.

Although our algorithms exhibit polynomial complexity, employing exceedingly large models increases the computational time required for calculating perplexity in extended dialogues<sup>5</sup>. Presently, the state-of-the-art lacks alternative unsupervised metrics with efficient time complexity for evaluating dialogue quality. Metrics such as FED (Mehri and Eskénazi, 2020) and FULL (Bruyn et al., 2022) entail computing multiple log-likelihood values for

<sup>5</sup>Detailed insights into the algorithm execution times are provided in Appendix F.

dialogue assessment, contrasting with perplexity, which mandates the computation of a singular log-likelihood value<sup>6</sup>. We leave the comparison among these metrics for future investigations.

## 7 Conclusion and Future Work

In this paper, we introduce an innovative, fully unsupervised method for extracting discourse structures in multi-party dialogues. To this end, we leverage open-source LLMs and introduce two algorithms, DS-DP and DS-FLOW, to detect coherent sub-dialogues within a dialogue. On the STAC and Molwени-clean corpora, we achieve micro- $F_1$  scores of 58.1% and 74.7%, respectively, demonstrating the efficacy of our solution in constructing dialogue structures without the need for labeled data. In the future, we intend to enhance the coherence evaluation metric, particularly addressing cognitive aspects as in Li et al. (2024b), and explore applying LLMs for unsupervised prediction of rhetorical relation types to deduce full discourse structures. Furthermore, we aim to improve algorithm link selection by incorporating linguistically motivated constraints as in Perret et al. (2016). Lastly, we plan to evaluate our architectural choices across diverse corpora and discourse parsing tasks to further validate their efficacy in assessing dialogues in real-world scenarios.

## Acknowledgments

The authors thank the anonymous reviewers for their insightful comments and suggestions. This research was supported by Mitacs as part of the Globalink Research Award (GRA) program. We acknowledge their assistance and funding, which made this study possible.

## References

- Stergos Afantenos, Nicholas Asher, Farah Benamara, Anais Cadilhac, Cedric Dégremont, Pascal Denis, Markus Guhe, Simon Keizer, Alex Lascarides, Oliver Lemon, et al. 2012. Modelling strategic conversation: model, annotation design and corpus. In *Proceedings of SemDial 2012 (SeineDial): The 16th Workshop on the Semantics and Pragmatics of Dialogue*, pages 167–168.
- Nicholas Asher, Julie Hunter, Mathieu Morey, Farah Benamara, and Stergos D. Afantenos. 2016. *Discourse structure and dialogue acts in multiparty dialogue*:

<sup>6</sup>An examination of the constraints associated with perplexity can be found in Appendix G.

- the STAC corpus. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation LREC 2016, Portorož, Slovenia, May 23-28, 2016*. European Language Resources Association (ELRA).
- Nicholas Asher and Alex Lascarides. 2003. *Logics of conversation*. Cambridge University Press.
- Sonia Badene, Kate Thompson, Jean-Pierre Lorré, and Nicholas Asher. 2019a. [Data programming for learning discourse structure](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 640–645. Association for Computational Linguistics.
- Sonia Badene, Kate Thompson, Jean-Pierre Lorré, and Nicholas Asher. 2019b. [Weak supervision for learning discourse structure](#). 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 2296–2305. Association for Computational Linguistics.
- Zineb Bennis, Julie Hunter, and Nicholas Asher. 2023. [A simple but effective model for attachment in discourse parsing with multi-task learning for relation labeling](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 3404–3409. Association for Computational Linguistics.
- Maxime De Bruyn, Ehsan Lotfi, Jeska Buhmann, and Walter Daelemans. 2022. [Open-domain dialog evaluation using follow-ups likelihood](#). In *Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022*, pages 496–504. International Committee on Computational Linguistics.
- Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song. 2023. [Chatgpt evaluation on sentence level relations: A focus on temporal, causal, and discourse relations](#). *CoRR*, abs/2304.14827.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%\\* chatgpt quality](#).
- Davide Colla, Matteo Delsanto, Marco Agosto, Benedetto Vitiello, and Daniele Paolo Radicioni. 2022. [Semantic coherence markers: The contribution of perplexity metrics](#). *Artif. Intell. Medicine*, 134:102393.
- Xiachong Feng, Xiaocheng Feng, and Bing Qin. 2022. [A survey on dialogue summarization: Recent advances and new frontiers](#). In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, pages 5453–5460. ijcai.org.
- Raquel Fernández, Matthew Frampton, Patrick Ehlen, Matthew Purver, and Stanley Peters. 2008. [Modelling and detecting decisions in multi-party dialogue](#). In *Proceedings of the SIGDIAL 2008 Workshop, The 9th Annual Meeting of the Special Interest Group on Discourse and Dialogue, 19-20 June 2008, Ohio State University, Columbus, Ohio, USA*, pages 156–163. The Association for Computer Linguistics.
- Matthew Frampton, Jia Huang, Trung H. Bui, and Stanley Peters. 2009. [Real-time decision detection in multi-party dialogue](#). In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, EMNLP 2009, 6-7 August 2009, Singapore, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1133–1141. ACL.
- Prakhar Gupta, Cathy Jiao, Yi-Ting Yeh, Shikib Mehri, Maxine Eskenazi, and Jeffrey Bigham. 2022. [InstructDial: Improving zero and few-shot generalization in dialogue through instruction tuning](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 505–525, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yuchen He, Zhuosheng Zhang, and Hai Zhao. 2021. [Multi-tasking dialogue comprehension with discourse parsing](#). In *Proceedings of the 35th Pacific Asia Conference on Language, Information and Computation, PACLIC 2021, Shanghai International Studies University, Shanghai, China, 5-7 November 2021*, pages 551–561. 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, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *CoRR*, abs/2310.06825.
- Shafiq R. Joty, Giuseppe Carenini, Raymond T. Ng, and Gabriel Murray. 2019. [Discourse analysis and its applications](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics: Tutorial Abstracts, ACL 2019, Florence, Italy, July 28, 2019, Volume 4: Tutorial Abstracts*, pages 12–17. Association for Computational Linguistics.
- Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2021. [Dialogue state tracking with a language model using schema-driven prompting](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, November 7-11, 2021*, pages 4937–4949. 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, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880. Association for Computational Linguistics.
- Chuyuan Li, Chloé Braud, Maxime Amblard, and Giuseppe Carenini. 2024a. [Discourse relation prediction and discourse parsing in dialogues with minimal supervision](#). In *Proceedings of the 5th Workshop on Computational Approaches to Discourse (CODI 2024)*, pages 161–176.
- Chuyuan Li, Patrick Huber, Wen Xiao, Maxime Amblard, Chloé Braud, and Giuseppe Carenini. 2023. [Discourse structure extraction from pre-trained and fine-tuned language models in dialogues](#). In *Findings of the Association for Computational Linguistics: EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 2517–2534. Association for Computational Linguistics.
- Jiaqi Li, Ming Liu, Min-Yen Kan, Zihao Zheng, Zekun Wang, Wenqiang Lei, Ting Liu, and Bing Qin. 2020. [Molweni: A challenge multiparty dialogues-based machine reading comprehension dataset with discourse structure](#). In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 2642–2652. International Committee on Computational Linguistics.
- Jiaqi Li, Ming Liu, Bing Qin, and Ting Liu. 2022. [A survey of discourse parsing](#). *Frontiers Comput. Sci.*, 16(5):165329.
- Manling Li, Lingyu Zhang, Heng Ji, and Richard J. Radke. 2019. [Keep meeting summaries on topic: Abstractive multi-modal meeting summarization](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 2190–2196. Association for Computational Linguistics.
- Xue Li, Jia Su, Yang Yang, Zipeng Gao, Xinyu Duan, and Yi Guan. 2024b. [Dialogues are not just text: Modeling cognition for dialogue coherence evaluation](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 18573–18581. AAAI Press.
- Annie Louis, Aravind Joshi, and Ani Nenkova. 2010. [Discourse indicators for content selection in summarization](#). In *Proceedings of the SIGDIAL 2010 Conference*, pages 147–156.
- Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. [The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems](#). In *Proceedings of the SIGDIAL 2015 Conference, The 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue, 2-4 September 2015, Prague, Czech Republic*, pages 285–294. The Association for Computer Linguistics.
- Tiezheng Mao, Tianyong Hao, Jialing Fu, and Osamu Yoshie. 2024. [Hierarchical graph fusion network and a new argumentative dataset for multiparty dialogue discourse parsing](#). *Inf. Process. Manag.*, 61(2):103613.
- Shikib Mehri and Maxine Eskénazi. 2020. [Unsupervised evaluation of interactive dialog with dialogpt](#). In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue, SIGdial 2020, 1st virtual meeting, July 1-3, 2020*, pages 225–235. Association for Computational Linguistics.
- Philippe Muller, Stergos D. Afantenos, Pascal Denis, and Nicholas Asher. 2012. [Constrained decoding for text-level discourse parsing](#). In *Proceedings of the 24th International Conference on Computational Linguistics (COLING 2012), December 8-15, 2012, Mumbai, India*, pages 1883–1900. Indian Institute of Technology Bombay.
- Noriki Nishida and Yuji Matsumoto. 2022. [Out-of-domain discourse dependency parsing via bootstrapping: An empirical analysis on its effectiveness and limitation](#). *Trans. Assoc. Comput. Linguistics*, 10:127–144.
- Jérémy Perret, Stergos D. Afantenos, Nicholas Asher, and Mathieu Morey. 2016. [Integer linear programming for discourse parsing](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL 2016): Human Language Technologies (HLT), San Diego California, USA, June 12-17, 2016*, pages 99–109. Association for Computational Linguistics.
- Vipula Rawte, Amit P. Sheth, and Amitava Das. 2023. [A survey of hallucination in large foundation models](#). *CoRR*, abs/2309.05922.
- Emanuel A Schegloff. 2007. *Sequence organization in interaction: A primer in conversation analysis I*, volume 1. Cambridge university press.
- Kai Sun, Su Zhu, Lu Chen, Siqiu Yao, Xueyang Wu, and Kai Yu. 2016. [Hybrid dialogue state tracking for real world human-to-human dialogues](#). In *Proceedings of the 17th Annual Conference of the International Speech Communication Association (Interspeech 2016), San Francisco, CA, USA, September 8-12, 2016*, pages 2060–2064. ISCA.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal

- Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. [Llama: Open and efficient foundation language models](#). *CoRR*, abs/2302.13971.
- 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, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.
- Ante Wang, Linfeng Song, Hui Jiang, Shaopeng Lai, Junfeng Yao, Min Zhang, and Jinsong Su. 2021. [A structure self-aware model for discourse parsing on multi-party dialogues](#). In *Proceedings of the 30th International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pages 3943–3949. ijcai.org.
- 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, EMNLP 2020 - Demos, Online, November 16-20, 2020*, pages 38–45. Association for Computational Linguistics.
- Wen Xiao, Patrick Huber, and Giuseppe Carenini. 2020. [Do we really need that many parameters in transformer for extractive summarization? discourse can help !](#) *CoRR*, abs/2012.02144.
- Jiacheng Xu, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. [Discourse-aware neural extractive text summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 5021–5031. Association for Computational Linguistics.
- Chen Zhang, Luis Fernando D’Haro, Yiming Chen, Malu Zhang, and Haizhou Li. 2024. [A comprehensive analysis of the effectiveness of large language models as automatic dialogue evaluators](#). In *Proceedings of the 38th AAAI Conference on Artificial Intelligence (AAAI 2024), 36th Conference on Innovative Applications of Artificial Intelligence (IAAI 2024), 14th Symposium on Educational Advances in Artificial Intelligence (EAAI 2014), February 20-27, 2024, Vancouver, Canada*, pages 19515–19524. AAAI Press.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2023. [Instruction tuning for large language models: A survey](#). *CoRR*, abs/2308.10792.

## A Further Evaluation Results for Link Length Analysis

To further evaluate our approaches, we analyze the precision and  $F_1$  scores for link lengths using DS-FLOW and DS-DP with Vicuna on the STAC and Molweni-clean corpora. As in Section 5.3, we explore different settings. Figure 4 shows that the STL+PF setting, which includes both the STL constraint and the penalization factor, provides the best precision scores for links with distances ranging from 2 to 5 in STAC and from 2 to 3 in Molweni-clean. This setting improves the evaluation of short-distance links, resulting in fewer false positives, but slightly lower precision ( $\sim 3\%$ ) for direct links due to the penalization factor. Additionally, although STL and vanilla settings predict long-distance links, they introduce several false positives. For instance, DS-FLOW with vanilla Vicuna on STAC predicted incorrect links with distances ranging from 14 to 31. Figure 5 highlights that in both corpora, all settings show decreasing  $F_1$  scores as link lengths increase, except for DS-FLOW with vanilla Vicuna on STAC, which achieves an  $F_1$  score of 40% for links of length 13 and 13% for links of length 12.

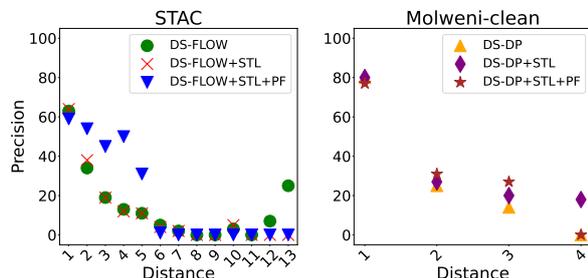


Figure 4: Precision scores for different link lengths. The left plot shows three settings with DS-FLOW on STAC; the right plot depicts the same settings for DS-DP on Molweni-clean. Both algorithms use Vicuna as LLM. STL: Speech turn limitation. STL+PF: Speech turn limitation in combination with penalization factor.

## B Filtering Mechanism Prompt Template

Drawing inspiration from the prompt utilized by Zhang et al. (2024) for evaluating dialogue qualities, we devise a new prompt specifically tailored to the task of predicting potential dialogue continuations, as depicted in Figure 6. We adapt the instruction template to align with the format used by Vicuna and Mistral in their instruction-tuning process.

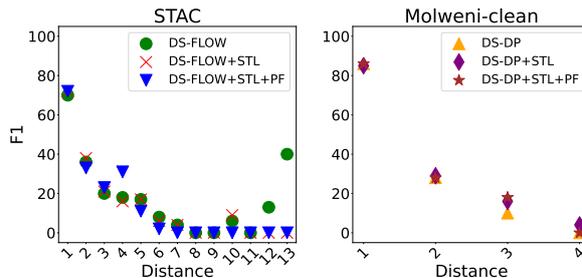


Figure 5:  $F_1$  scores for different link lengths. The left plot shows three settings with DS-FLOW on STAC; the right plot depicts the same settings for DS-DP on Molweni-clean. Both algorithms use Vicuna as LLM. STL: Speech turn limitation. STL+PF: Speech turn limitation in combination with penalization factor.

```

### Context:
[Here is a dialogue utterance]

### Response:
[Here is the potential continuation]

### Instruction:
Please evaluate whether the response is a plausible continuation of the given utterance within a dialogue context and provide a definitive answer Yes or No.

### Your Answer:
[Here is LLM's output in terms of "Yes" or "No"]

```

Figure 6: An example of how open-source LLMs can be prompted to determine if an utterance could potentially follow a preceding one.

## C Filtering Mechanism Evaluation

We assessed the performance of the filtering mechanism under the optimal setting for the DS-FLOW algorithm, specifically leveraging the STL constraint, the penalization factor, and the Vicuna-13b model as the backbone on the STAC corpus. The following scenarios were considered:

- **True Positives:** Links that should be filtered and are correctly identified as such by the LLM (112 instances).
- **False Positives:** Links that should not be filtered but are incorrectly identified as filtered by the LLM (52 instances).
- **True Negatives:** Links that should not be filtered and are correctly identified as such by the LLM (698 instances).
- **False Negatives:** Links that should be filtered but are incorrectly identified as not filtered by the LLM (604 instances).

Model	Algorithm	STAC			Molweni-clean		
		F <sub>1</sub>	P	R	F <sub>1</sub>	P	R
Vicuna-7b	DS-DP	53.8	51.7	56.0	69.6	66.4	73.2
	DS-DP + STL	54.4	52.3	56.6	69.9	66.7	73.5
	DS-DP + STL + $P_{dist}(d)$	56.9	53.2	61.3	73.2	67.6	79.9
	DS-FLOW	46.0	38.6	56.9	57.6	47.7	72.7
	DS-FLOW + STL	46.4	39.3	56.6	58.5	48.7	73.2
	DS-FLOW + STL + $P_{dist}(d)$	56.3	55.1	57.5	72.8	68.4	77.7

Table 4: Experiment results of proposed approaches on STAC and Molweni-clean corpora for Vicuna-7b. STL: Speech turn limitation.  $P_{dist}(d)$ : Penalization factor. F<sub>1</sub>: Micro-F<sub>1</sub>. P: Precision. R: Recall.

Based on these outcomes, we calculated the Accuracy, Precision, Recall, and F<sub>1</sub> scores, as detailed in Table 5. The filtering mechanism exhibits good reliability, demonstrated by a Precision score of 68.2%. However, it only filters out a small number of incorrect potential continuations, resulting in a Recall score of 15.6%, which in turn affects the F<sub>1</sub> score. The overall Accuracy score of 55.3% is consistent with the algorithm’s performance on the STAC corpus. Enhancing the filtering mechanism is expected to improve the algorithm’s performance, a subject we plan to address in future work.

Metric	Value (%)
Precision	68.2
Recall	15.6
Accuracy	55.3
F <sub>1</sub>	25.1

Table 5: Performance metrics of the filtering mechanism under the optimal setting for the DS-FLOW algorithm, leveraging the STL constraint, the penalization factor, and the Vicuna-13b model on the STAC corpus.

## D Experimental Analysis with Smaller Vicuna Model

To analyze how performance changes with LLM model size, we conduct supplementary analyses using a smaller version of Vicuna, comprising 7b parameters. As shown in Table 4, akin to Vicuna-13b (see Table 2), both algorithms exhibit optimal performance when incorporating both the STL constraint and the penalization factor, with a slight improvement observed when integrating the STL constraint compared to the vanilla version. Regarding the best settings, the results indicate that with the downsized LLM, the micro-F<sub>1</sub> scores are slightly lower, with DS-FLOW achieving 56.3% and DS-DP achieving 73.2%, compared to Vicuna-13b, which achieved 58.1% on STAC and 74.7%

on Molweni-clean, respectively. This suggests that employing a larger LLM could potentially yield superior outcomes.

## E An Example of DS-DP Algorithm Execution

Figure 7 illustrates the application of DS-DP to the dialogue depicted in Figure 1. The algorithm begins by calculating perplexity scores for sub-dialogues of length 2 during the initialization phase. It then progresses to compute the perplexity scores for sub-dialogues of length 3. Specifically, when  $k = 1$ , the algorithm analyzes all pairs of EDUs ( $e_i, e_j$ ) with  $i > 1$  and  $j > i$ . Given the constraint that sub-dialogues must start from the initial EDU (see Section 3.1), every pair of sequential EDUs has the initial EDU as the preceding one.

For  $k = 2$ , the algorithm computes sub-dialogues of length 4. For the cells  $M_{co}[3][4][2]$ ,  $M_{co}[3][5][2]$ , and  $M_{co}[3][6][2]$ , there is only one possible sub-dialogue, and the algorithm computes their perplexity scores. When evaluating the sub-dialogue passing through  $e_4$  and ending in  $e_5$ , the algorithm analyzes the incoming links in  $e_4$ . According to the input graph,  $e_4$  has incoming links from  $e_1$ ,  $e_2$ , and  $e_3$ . Since there is no sub-dialogue passing through  $e_1$ , ending in  $e_4$ , and involving  $k - 1 = 1$  EDU,  $e_1$  is disregarded. Instead, the algorithm considers  $e_2$  and  $e_3$  as intermediary EDUs to conclude in  $e_5$  through  $e_4$ . Thus, it analyzes two sub-dialogues (A): ( $e_1, e_2, e_4, e_5$ ) and ( $e_1, e_3, e_4, e_5$ ). The algorithm computes the perplexity scores for both sub-dialogues and retains the one with the lowest perplexity, e.g., ( $e_1, e_3, e_4, e_5$ ). The same method applies to sub-dialogues passing through  $e_4$  and ending in  $e_6$  (B), and those passing through  $e_5$  and ending in  $e_6$  (C). Consider the selection of sub-dialogues ( $e_1, e_3, e_4, e_6$ ) and ( $e_1, e_2, e_5, e_6$ ).

For  $k = 3$ , the algorithm examines pairs of EDUs that can involve three preceding EDUs, such as ( $e_4, e_5$ ), ( $e_4, e_6$ ), and ( $e_5, e_6$ ). The first two pairs

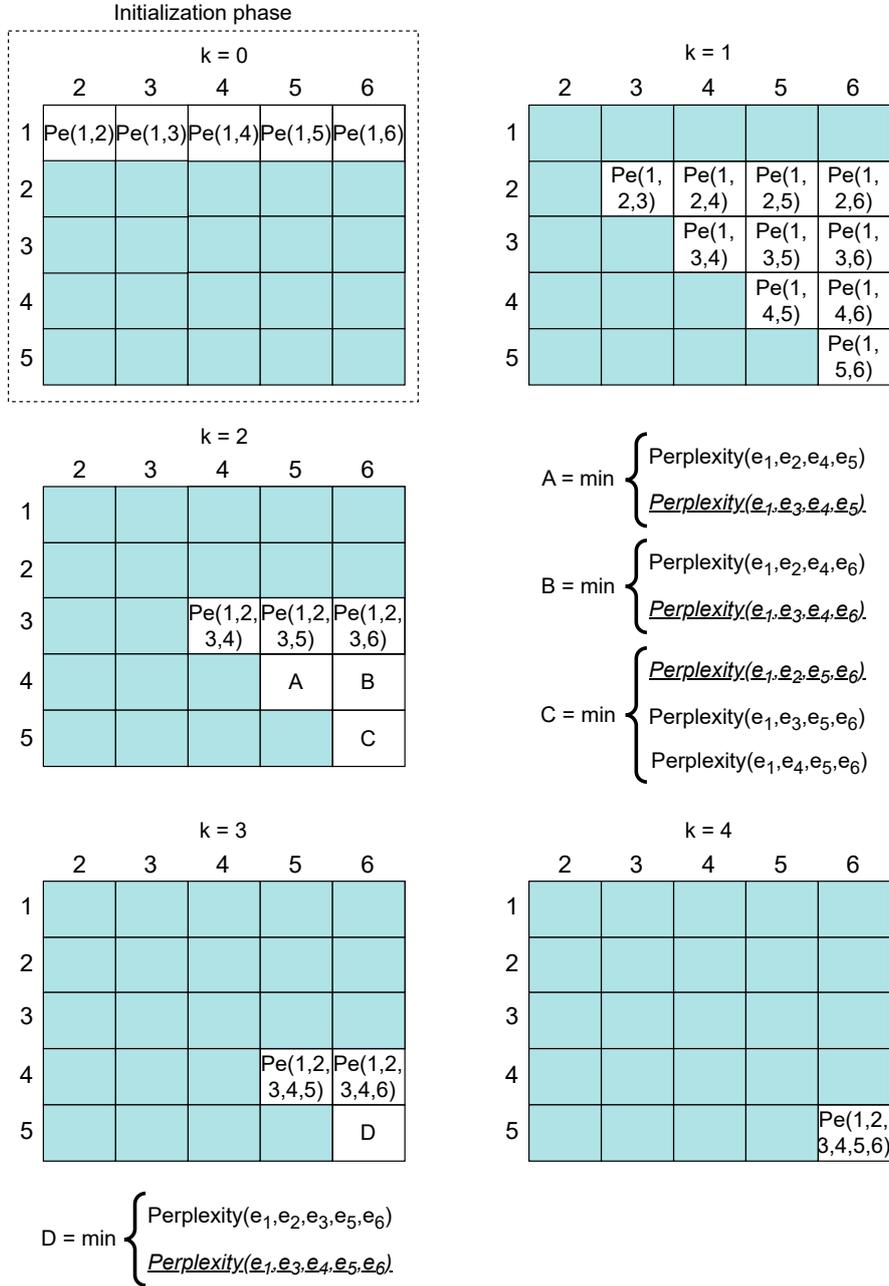


Figure 7: An example of DS-DP execution. For brevity, we use the notation  $Pe(\text{list of indexes})$  instead of  $Perplexity(\text{list of EDUs})$  within the cells. Underlined texts denote selected sub-dialogues during the algorithm's execution.

are constrained by preceding EDUs  $e_1$ ,  $e_2$ , and  $e_3$ . For the pair  $(e_5, e_6)$ , multiple triples can serve as the preceding three EDUs. Here, the algorithm considers EDUs with outgoing links towards  $e_5$ , namely  $e_1$ ,  $e_2$ ,  $e_3$ , and  $e_4$ . Only two EDUs,  $e_3$  and  $e_4$ , can reach  $e_5$  and involve two preceding EDUs. Since the algorithm selected the sub-dialogue  $(e_1, e_3, e_4, e_5)$  as the best option for passing through  $e_4$  and ending in  $e_5$  with two preceding EDUs, it does not analyze the sub-dialogue  $(e_1, e_2, e_4, e_5, e_6)$  and just considers  $(e_1, e_3, e_4, e_5, e_6)$ . Conversely, the only sub-dialogue passing through  $e_3$  and ending in  $e_5$  with two preceding EDUs is  $(e_1, e_2, e_3, e_5)$ , leading the algorithm to analyze  $(e_1, e_2, e_3, e_5, e_6)$ . Consider the selection of  $(e_1, e_3, e_4, e_5, e_6)$ .

For  $k = 4$ , the only pair of EDUs that can have four preceding EDUs is  $(e_5, e_6)$ , resulting in the sub-dialogue  $(e_1, e_2, e_3, e_4, e_5, e_6)$ . With this, the algorithm completes the computation of the most coherent sub-dialogues with lengths ranging from 2 to 6. Then, it iterates over the  $k$  value for each EDU  $e_i$  and selects the sub-dialogue ending in  $e_i$  with the minimum perplexity, examining only the column  $i$  for each  $k$ . For example, to find the most coherent sub-dialogue ending in  $e_4$ , it evaluates the perplexity scores of the following sub-dialogues:  $(e_1, e_4)$ ,  $(e_1, e_2, e_4)$ ,  $(e_1, e_3, e_4)$ , and  $(e_1, e_2, e_3, e_4)$ . In the context of the dialogue in Figure 1, the algorithm selects  $(e_1, e_2)$  for  $e_2$ ,  $(e_1, e_3)$  for  $e_3$ ,  $(e_1, e_3, e_4)$  for  $e_4$ ,  $(e_1, e_3, e_4, e_5)$  for  $e_5$ , and  $(e_1, e_3, e_4, e_6)$  for  $e_6$ , resulting in the final DAG:  $\{(e_1, e_2), (e_1, e_3), (e_3, e_4), (e_4, e_5), (e_4, e_6)\}$ .

From the matrices, it is evident that there is a significant number of empty cells (indicated in light blue). For each  $k$  value, the algorithm only needs to examine rows with indices greater than  $k$ , and for each row  $i$ , only columns with indices greater than  $i$ . This is justified by the assumption of not having backward links within the final DAG to be computed.

## F Algorithm Execution Time Analysis

We assessed the execution times of the proposed algorithms using the STAC and Molweni-clean corpora. Although the DS-FLOW algorithm exhibits a time complexity of  $O(|V|^2)$ , which is more favorable compared to the  $O(|V|^4)$  complexity of the DS-DP algorithm, our empirical analysis revealed that the DS-DP algorithm computes discourse structures more efficiently, with execution times sometimes reduced by up to half. This discrepancy oc-

curs because, even though the DS-DP algorithm has a  $O(|V|^4)$  time complexity, it processes fewer values than expected in the worst-case scenario (as detailed in Appendix E). Furthermore, the DS-FLOW algorithm requires LLM computation for both filtering and perplexity calculations, while the DS-DP algorithm uses an LLM solely for evaluating sub-dialogue coherence.

## G Qualitative Analysis in STAC and Molweni-clean

In Figures 8-19, we present several concrete examples generated by the optimal approaches for STAC (utilizing DS-FLOW with the STL constraint and penalization factor, and Vicuna-13b as the backbone) and Molweni-clean (employing DS-DP with identical settings). Specifically, we show three well-predicted examples (depicted in Figures 8, 9, and 10 for STAC, and Figures 14, 15, and 16 for Molweni-clean) and three badly predicted examples (depicted in Figures 11, 12, and 13 for STAC, and Figures 17, 18, and 19 for Molweni-clean). Some patterns observed in predicted structures include: (i) the algorithms struggle to predict very long-distance links, favoring shorter links with distances of 2, 3, and 4; (ii) direct links are often predicted even when the appropriate indirect incoming links for EDUs are accurately identified.

Our qualitative analysis has identified multiple instances wherein the application of perplexity for dialogue evaluation presents limitations. To exemplify this issue, consider the following dialogue excerpt from the STAC corpus (id *s1-league3-game3\_16*):

- ( $e_1$ ) A: can anyone trade ore? I have more wood to trade
- ( $e_2$ ) B: do you have clay, by any chance?
- ( $e_3$ ) A: sorry, no
- ( $e_4$ ) C: i can do that kieran
- ( $e_5$ ) A: how many can you trade?
- ( $e_6$ ) A: 2 for 2?
- ( $e_7$ ) C: just got one, sorry
- ( $e_8$ ) A: ok cool

In this instance, our algorithms evaluated the links  $(e_5, e_7)$  as more likely than  $(e_5, e_6)$ , despite both links being valid:  $(e_5, e_7)$  with a *QAP* relation and  $(e_5, e_6)$  with a *Continuation* relation. This discrepancy is likely because a direct answer like  $e_7$  seems more contextually relevant as a response to  $e_5$ , thus overshadowing the *Continuation* link to  $e_6$ .

The perplexity metric tends to favor more immediate and clear connections, which can sometimes misrepresent the actual flow of dialogue. This limitation indicates that relying solely on perplexity for dialogue evaluation may overlook nuanced conversational dynamics, underscoring the need for supplementary metrics to fully capture dialogue coherence and relevance.

## **H Dialogue Parsing Task Prompt Template**

To enhance the competitive performance of ChatGPT in the multi-party dialogue discourse parsing task, we undertake manual design efforts to refine the prompt. This refinement, illustrated in Figure 20, builds upon the prompt proposed by Chan et al. (2023). Specifically, we provide a more explicit delineation of the task requirements by specifying the extraction of a DAG, in contrast to the broader objective pursued by Chan et al. (2023), which involved predicting all potential discourse relations between utterances. Furthermore, drawing on insights from the findings of Chan et al. (2023), who demonstrated improved performance with the inclusion of descriptions for discourse relations, we develop more comprehensive descriptions within the prompt. This refined prompt has been consistently used in zero-shot experiments conducted with the Vicuna and Mistral models.

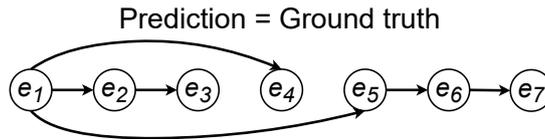


Figure 8: STAC - DS-FLOW - Well predicted example: *pilot02\_12*. #EDUs : 7.

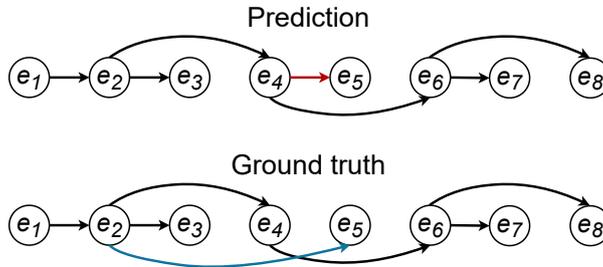


Figure 9: STAC - DS-FLOW - Well predicted example: *pilot02\_21*. #EDUs : 8. In red: False positive edges; in light blue: False negative edges.

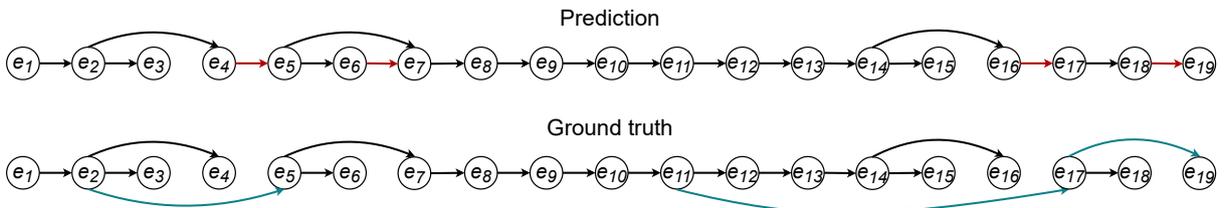


Figure 10: STAC - DS-FLOW - Well predicted example: *pilot02\_13*. #EDUs : 19. In red: False positive edges; in light blue: False negative edges.

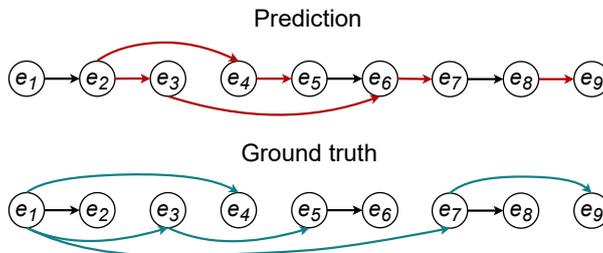


Figure 11: STAC - DS-FLOW - Badly predicted example: *s2-league4-game2\_6*. #EDUs : 9. In red: False positive edges; in light blue: False negative edges.

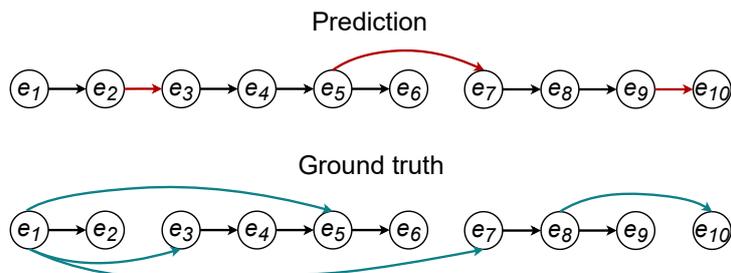


Figure 12: STAC - DS-FLOW - Badly predicted example: *pilot02\_6*. #EDUs : 10. In red: False positive edges; in light blue: False negative edges.

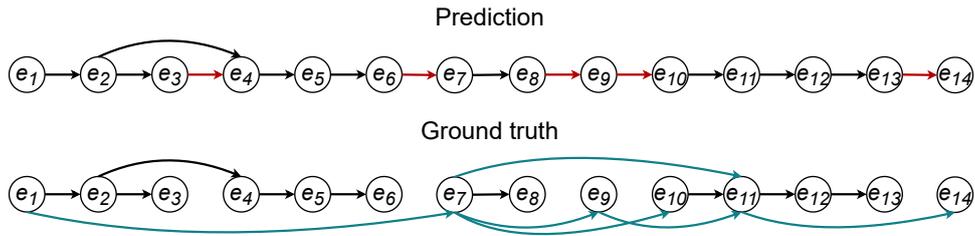


Figure 13: STAC - DS-FLOW - Badly predicted example: *s2-league4-game2\_31*. #EDUs : 14. In red: False positive edges; in light blue: False negative edges.

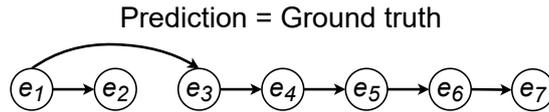


Figure 14: Molweni-clean - DS-DP - Well predicted example: *8031*. #EDUs : 7.

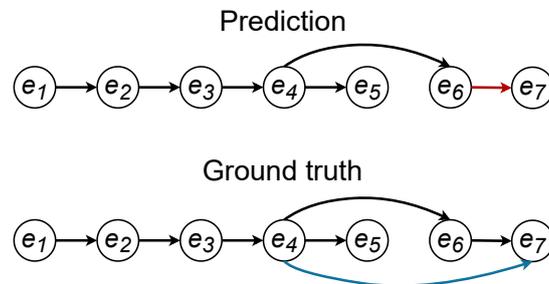


Figure 15: Molweni-clean - DS-DP - Well predicted example: *6037*. #EDUs : 7. In red: False positive edges; in light blue: False negative edges.

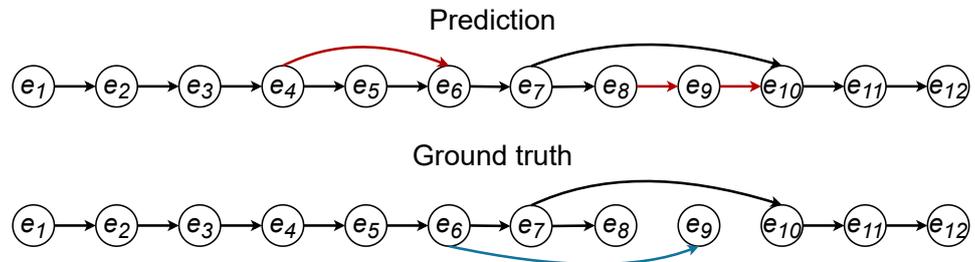


Figure 16: Molweni-clean - DS-DP - Well predicted example: *8026*. #EDUs : 12. In red: False positive edges; in light blue: False negative edges.

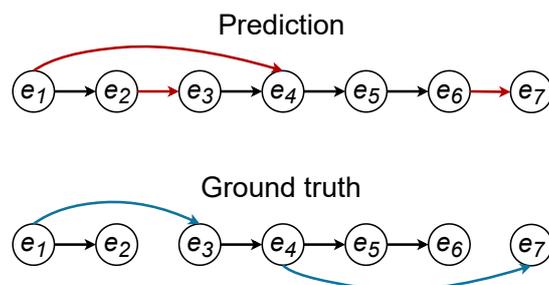


Figure 17: Molweni-clean - DS-DP - Badly predicted example: *5033*. #EDUs : 7. In red: False positive edges; in light blue: False negative edges.

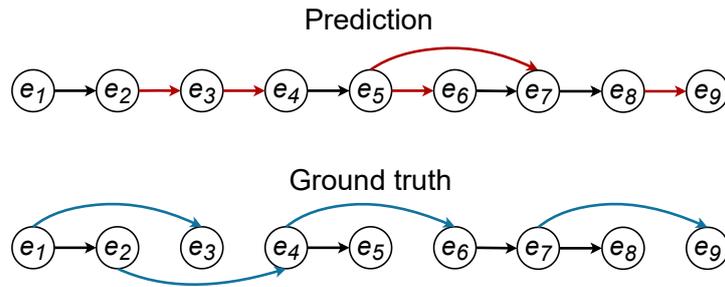


Figure 18: Molweni-clean - DS-DP - Badly predicted example: 8039. #EDUs : 9. In red: False positive edges; in light blue: False negative edges.

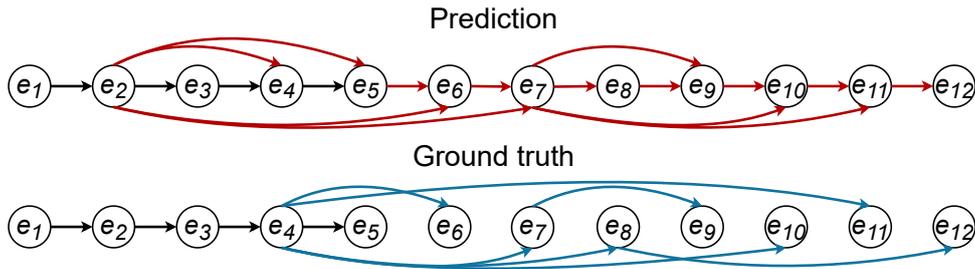


Figure 19: Molweni-clean - DS-DP - Badly predicted example: 8018. #EDUs : 12. In red: False positive edges; in light blue: False negative edges.

Here is a multi-party dialogue:

**[Multi-party dialogue]**

Assume that each utterance represents a node within a graph. Your task is to predict the relations between these utterances based on a provided list of relations. The resulting graph should adhere to the structure of a Directed Acyclic Graph (DAG), wherein edges have a direction, meaning they go from one node to another. A key characteristic of a DAG is that it does not contain cycles, i.e., there are no sequences of edges that form a closed loop. This implies that it is not possible to start from a node, follow the edges, and return to the starting node. It is crucial to emphasize that each node representing an utterance must have at least one incoming edge to ensure that the resulting graph maintains coherence and fosters a connected discourse.

*Relations:*

- 1) *Comment*: This relation type typically indicates that one utterance provides a comment or opinion on the content of another utterance. It shows a speaker's perspective or evaluation of the preceding statement.
- 2) *Clarification Question*: In this relation, one utterance poses a question seeking clarification or additional information about the content of another utterance. It implies a request for further explanation.
- 3) *Elaboration*: Elaboration signifies that one utterance expands upon or provides more details about the content of another utterance. It is used to enhance understanding by offering additional information or context.
- 4) *Continuation*: This relation suggests that one utterance continues the topic or discussion from a previous utterance. It signifies a logical progression in the conversation.
- 5) *Explanation*: Explanation pertains to one utterance offering an explanation or clarification in response to a question or confusion expressed in another utterance. It aids in providing clarity.
- 6) *Conditional*: A conditional relation implies that one utterance presents a condition or hypothetical scenario related to the content of another utterance. It often involves "if-then" statements.
- 7) *Question-Answer Pair*: This relation indicates that one utterance contains a question, and another utterance follows with an answer to that question. It demonstrates a direct question-and-answer interaction.
- 8) *Alternation*: Alternation shows that two utterances present alternative options or choices. It is used when discussing multiple possibilities or courses of action.
- 9) *Q-Elab*: Q-Elab signifies that one utterance asks a question, and another utterance follows with an elaboration or further explanation of the question or its context.
- 10) *Result*: Result indicates that one utterance discusses the outcome or consequence of the content presented in another utterance. It shows a cause-and-effect relation.
- 11) *Background*: In this relation, one utterance provides background information or context that is relevant to the content of another utterance. It helps set the stage for the discussion.
- 12) *Narration*: Narration signifies that one utterance presents a narrative or storytelling element, often in response to a question or to share an anecdote.
- 13) *Correction*: Correction shows that one utterance corrects or revises the content of another utterance. It is used to rectify errors or inaccuracies.
- 14) *Parallel*: Parallel relations occur when two or more utterances share similar or related content, often in a parallel or analogous manner. It emphasizes similarities or comparisons.
- 15) *Contrast*: Contrast signifies that one utterance presents content that is in contrast or opposition to the content of another utterance. It highlights differences or contradictions in the conversation.

Figure 20: Prompt template employed for LLMs in a zero-shot setting for the multi-party dialogue discourse parsing task on the STAC and Molweni-clean corpora.

# Transforming Slot Schema Induction with Generative Dialogue State Inference

James D. Finch and Boxin Zhao and Jinho D. Choi

Department of Computer Science

Emory University

Atlanta, GA, USA

{jdfinch, zinc.zhao, jinho.choi}@emory.edu

## Abstract

The challenge of defining a slot schema to represent the state of a task-oriented dialogue system is addressed by Slot Schema Induction (SSI), which aims to automatically induce slots from unlabeled dialogue data. Whereas previous approaches induce slots by clustering value spans extracted directly from the dialogue text, we demonstrate the power of discovering slots using a generative approach. By training a model to generate slot names and values that summarize key dialogue information with no prior task knowledge, our SSI method discovers high-quality candidate information for representing dialogue state. These discovered slot-value candidates can be easily clustered into unified slot schemas that align well with human-authored schemas. Experimental comparisons on the MultiWOZ and SGD datasets demonstrate that Generative Dialogue State Inference (GenDSI) outperforms the previous state-of-the-art on multiple aspects of the SSI task.

## 1 Introduction

Developing Task-Oriented Dialogue (TOD) systems presents the significant challenge of creating and maintaining a *slot schema*, where each slot defines a type of information that is critical for successfully completing the dialogue task (Budzianowski et al., 2018). Traditionally slot schemas are handcrafted, but manually defining each slot is time-consuming, especially when task domains are complicated or the functionality of the dialogue system is frequently updated. To address this, Slot Schema Induction (SSI) has been proposed to automatically generate slot schemas from unlabeled dialogue data (Chen et al., 2013; Min et al., 2020). This task facilitates the automatic analysis of dialogue structure (Qiu et al., 2022) and identifies key types of information that should be included in dialogue state representations (Min et al., 2020). By reducing the need for manual

schema creation, SSI expedites developing TOD systems for new application domains, and enables continual discovery of new slot types to improve the coverage of existing slot schemata.

The core challenge of SSI is identifying which information presented in unlabeled dialogue data is important for the task domain and should be included in the dialogue state. Once the important information values are identified, a second challenge is defining a minimal set of slots that captures the different types of information the values represent. All previous work on SSI tackles these challenges in an explicit two-step process involving (1) candidate value identification and (2) inducing a slot schema by clustering candidate values into a set of slot clusters. Identifying value candidates has been explored using tagging models trained on other tasks such as NER or SRL (Min et al., 2020; Hudeček et al., 2021; Qiu et al., 2022; Wu et al., 2022), or using token attention distributions produced by a PLM to extract syntactic constituents (Yu et al., 2022). Inducing slots from value candidates has been explored using out-of-the-box clustering algorithms (Qiu et al., 2022), multi-stage clustering pipelines specific to SSI (Hudeček et al., 2021; Wu et al., 2022; Yu et al., 2022), or a neural latent variable model (Min et al., 2020).

Unlike all previous approaches to SSI, we are the first to take a generative approach to value candidate identification.<sup>1</sup> Candidates are identified using a dialogue state generator model, which is trained to summarize the key task-related information in a given dialogue context as a set of state values. Crucially, this state generator also creates a slot name for each value, which serves as a candidate prediction of the name of the slot the value fills. Value candidates are then clustered in conjunction with these predicted slot names to induce a uni-

<sup>1</sup>The code, models, and data for our approach is publicly available at <https://github.com/emorynlp/GenDSI>.

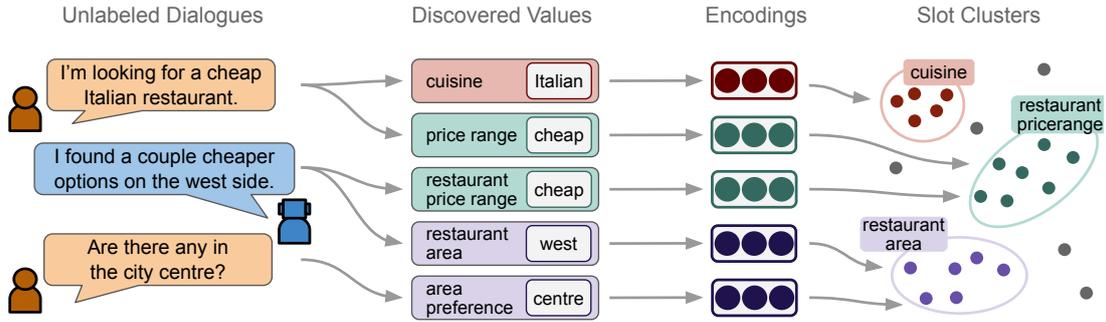


Figure 1: Overview of the GenDSI approach.

fied set of slots. The advantage of this approach is that the type semantics of each value candidate are concretely represented using slot name predictions, whereas previous approaches rely on the similarity of value encodings alone to cluster values into slot types. Predicting slot names has the additional benefit of enabling automatic naming of each slot cluster in the induced schema. We demonstrate the impact of these benefits by evaluating our approach on the MultiWOZ (Eric et al., 2020) and SGD (Rastogi et al., 2020) datasets. Our SSI approach produces slot schemas that better match gold reference schemas when compared to the previous state-of-the-art (SoTA) approaches.

## 2 Approach

Our SSI approach, Generative Dialogue State Inference (GenDSI), induces a slot schema consisting of a set of slot clusters from an unlabeled set of dialogues. The induction procedure is performed in three stages (Fig. 1). First, a dialogue state generator discovers value candidates for each turn in the dialogue data and jointly predicts a slot name with each value. Second, an encoding model produces a dense vector representation for each slot-value candidate. Finally, a clustering algorithm uses the encodings to filter and group candidates into a unified set of slot clusters.

**Dialogue State Generator** Our approach formulates the discovery of value candidates from unlabeled dialogue data as a sequence-to-sequence generation task. The input is a dialogue context  $D_{*..t}$ , and the output is a list of slot-value candidates  $[(s_1, v_1), \dots, (s_k, v_k)]$  represented by the sequence format  $s_1:v_1; s_2:v_2; \dots; s_k:v_k[EOS]$ . Each candidate includes a value  $v_i$  that is inferred to belong to the dialogue state and a slot name prediction  $s_i$  to represent the type of  $v_i$ . To enable the model to generate slot-value pairs that discover important dialogue

state information without any prior knowledge of the task domain, we fine-tune a pretrained encoder-decoder transformer on TOD data that covers a large variety of domains. Section 3 presents an evaluation of the dialogue state generator in which two different training datasets are compared.

**Value Candidate Encoding** Each slot-value candidate  $(s_i, v_i)$  produced by the dialogue state generator is encoded into a single dense vector representation  $e_i$ . To do this, we concatenate the slot name and value candidate with a separator to form a single token sequence  $s_i:v_i$ . We then use the SBERT encoder (Reimers and Gurevych, 2019) to independently encode each candidate sequence.

**Slot Clustering** Given a complete list of all slot-value candidates  $[(s_1, v_1), \dots, (s_n, v_n)]$  produced by the dialogue state generator across all turns of the dialogue dataset, slot-value candidates are jointly filtered and grouped by applying the HDBSCAN algorithm (McInnes et al., 2017) to the candidate encodings  $[e_1, e_2, \dots, e_n]$ . As demonstrated in previous work (Yu et al., 2022), HDBSCAN is a suitable clustering algorithm because (1) it automatically discovers an appropriate number of slot clusters and (2) it filters out examples in low-density regions of the encoding space, which are likely to represent noisy candidates. The result is a set of slot clusters  $[S_1, S_2, \dots, S_k]$  where each cluster  $S_i$  is a list of values that fill the slot it represents.

## 3 State Generator Evaluation

Since our SSI approach relies mainly on the dialogue state generator component to infer high-quality value candidates with appropriate slot names, we first conduct an evaluation of the performance of this component when discovering slot-values from dialogues in unseen task domains.

**Metrics** The dialogue state generator is evaluated by human judges, since discovered slot-value candidates are generated and many surface forms can be equally correct. We recruit three university students as volunteers to evaluate two key aspects of slot value candidate inferences. (1) *Completeness* measures the proportion of dialogue turns for which all key information has been captured as slot-value candidates. (2) *Correctness* measures the proportion of slot-value candidates that accurately represent specific information in their corresponding turns. Details of metrics are presented in Appx. A. This evaluation is performed using a custom annotation software, which was developed to optimize the efficiency of the annotation work. The interface is shown in Appx. C.

To validate our human evaluation metrics, inter-annotator agreement was calculated between the three human judges at 0.43 and 0.27 Krippendorff’s Alpha for Correctness and Completeness respectively. Based on a manual review of the annotation disagreements for Completeness items, we believe the lower agreement occurs because judges are required to consider more information across an entire state update compared to judging the correctness of a single slot-value pair, leading to higher annotation difficulty and thus some noisy judgments.

**Data** Since our goal is to train a dialogue state generator to discover slot-value candidates for unseen domains, we experiment with two domain-diverse datasets for training: SGD (Rastogi et al., 2020) and D0T (Finch and Choi, 2024). SGD is a popular TOD dataset that contains 20 domains and 16, 142 dialogues, with gold dialogue state labels. D0T is a recent dataset that was created using a fully automatic data generation method based on GPT-3.5 and GPT-4. It covers a large 1, 003 domains across 5, 015 dialogues, but it contains some noisy labels from automatic annotation.

We adapt these datasets for slot-value candidate discovery by simply training on dialogue state *updates* instead of full dialogue states, which represent only the slots that are filled by new values. This avoids training the dialogue state generator to predict empty slots, which are not useful for candidate discovery. Additionally, the special request value "?" is removed from D0T state updates. We also replace each slot name in the SGD training split with a random synonymous name from SGD-X (Lee et al., 2022), as we found this augmentation

to improve performance.

Both SGD and D0T are also used as evaluation data by randomly sampling 60 turns from their test splits, each from a unique dialogue. Crucially, we *only* sample turns from domains *not* included in the training split. Since the D0T dataset has no native splits for training and testing, we randomly sample 100 domains out of the total 1, 003 to be held-out for evaluation. The D0T training split thus includes only the remaining 4, 515 dialogues.

**Models** We train two models using T5-3B (Raffel et al., 2020) as a base model: T5-D0T trained on D0T and T5-SGD trained on SGD. We also compare against the GPT-based automatic annotator used to create silver D0T state update labels (GPT-D0T). Implementation details presented in Appx. D.

Model	D0T			SGD		
	CP	CR	HM	CP	CR	HM
T5-SGD	32.3	72.6	44.7	69.3	<b>90.8*</b>	78.6
GPT-D0T	93.3*	<b>82.0*</b>	87.3	90.0*	84.7	87.3
T5-D0T	<b>95.7*</b>	81.2*	<b>87.9</b>	<b>94.7†</b>	81.7	<b>87.7</b>

Table 1: Human evaluation of completeness (CP), correctness (CR), and their harmonic mean (HM) for each dialogue state generator. \*/† denote statistical significance against unstarred/all results in the same column (Agresti-Caffo,  $p < 0.05$ ).

**Results** As shown in Table 1, T5-D0T exhibits the best overall performance, achieving approximately 81% correct slot-value inferences and completely covering all key information in 95% of turns. The fact that there was nearly zero performance drop-off on the out-of-distribution SGD evaluation demonstrates its robustness in discovering useful slot-values for new domains. As expected, GPT-D0T exhibits similar performance, as it generated the labels used to train T5-D0T; however, GPT-D0T is much costlier due to multiple API calls to GPT3.5 and GPT4. T5-SGD achieves the highest correctness score when evaluated on held-out SGD domains, but its completeness score of only 70% demonstrates it is incapable of fully adapting to unseen domains. On the out-of-distribution D0T evaluation, the performance of T5-SGD heavily suffers, achieving only 32% completeness and 73% correctness. This result demonstrates the difficulty of discovering state information in unseen domains, and indicates that SGD is insufficiently diverse as a training resource for this purpose.

Model	MultiWOZ							SGD						
	C	Slot			Value			C	Slot			Value		
		P	R	F1	P	R	F1		P	R	F1	P	R	F1
DSI	522	96.2	80.7	87.7	41.5	57.4	37.2	11992	-	-	<b>92.2</b>	-	-	46.2
USI	290	<b>100.0</b>	93.6	<b>96.7</b>	61.3	67.3	58.7	806	-	-	77.0	-	-	47.5
GenDSI	180	85.6	<b>96.8</b>	90.9	81.4	<b>70.2</b>	70.5	746	<b>92.4</b>	77.9	84.5	65.4	<b>50.0</b>	<b>48.8</b>
- slot names	<b>157</b>	73.9	90.3	81.3	85.2	47.7	55.3	<b>467</b>	76.4	75.6	76.0	<b>70.6</b>	36.3	37.9
+ all domains	161	85.1	<b>96.8</b>	90.6	<b>87.9</b>	68.1	<b>71.0</b>	737	90.8	<b>79.1</b>	84.5	68.0	47.2	47.7

Table 2: Schema induction results showing Precision/Recall/F1 (P/R/F1) for both induced slots and discovered values, as well as the induced Slot Count (C). Note that the optimal Slot Count would equal the gold slot counts of 31 and 82 for MultiWOZ and SGD respectively. DSI and USI results taken from Yu et al. (2022).

## 4 Schema Induction Evaluation

To evaluate our SSI approach, we use the benchmark defined by Yu et al. (2022) on the validation splits of MultiWOZ 2.1 (Eric et al., 2020) and SGD (Rastogi et al., 2020) datasets. This evaluation method measures the quality of an induced set of slot clusters by matching it against a gold reference slot schema.

Matching is performed automatically by computing the centroid of each induced and gold reference slot cluster using BERT encodings (Devlin et al., 2019) of their values. Each induced cluster is mapped to the gold slot whose cluster centroid is nearest by cosine similarity, or to no cluster if there is no match of 80% similarity or higher. Similarly, in order to evaluate the purity and coverage of clustered values, discovered values are matched against the gold value labels that fill each slot. This value matching is performed between the values that fill each gold slot and the discovered values of all induced clusters mapped to that slot using fuzzy string matching.

**Metrics** Given the mapping of induced clusters to gold slots, *Slot Precision* measures the proportion of induced clusters that were able to be matched to a gold slot, *Slot Recall* is the proportion of gold slots that were matched with at least one induced cluster, and *Slot F1* is their harmonic mean. Since multiple induced slots are allowed to map to a single gold slot, the induced *Slot Count* is also reported to measure redundancy. *Value Precision* is the average proportion of discovered values that matched to gold values, averaged across all gold slots. Similarly, *Value Recall* is the average proportion of gold values that were matched to a discovered value, and *Value F1* is the average F1 score across all gold slots. Equations defining these metrics are presented in Appx. B.

**Models** Our SSI approach, GenDSI, uses a T5-3B model trained on the D0T dataset. Since D0T contains some task domains that are related to travel domains appearing in MultiWOZ and SGD, we manually review and filter out 34 domains with overlap and train our dialogue state generator on all D0T dialogues in remaining domains. We also evaluate the performance when using a model trained with all D0T domains (GenDSI +all domains), which simulates extending our approach using the D0T data generation method to create synthetic training resources for target domains. Additionally, we evaluate a version of our approach where value candidates are encoded without their predicted slot names (GenDSI -slot names) to measure the benefit of concretely representing value type information. Implementation details provided are in Appx. D. Finally, we compare to two strong baselines from previous work:

- DSI (Min et al., 2020), which leverages a Part-of-Speech (POS) tagger, Named Entity Recognition (NER) tagger, and coreference resolution model to extract value candidate spans using a set of heuristic rules. Slot clusters are then assigned to value candidates using a neural latent variable model.
- USI (Yu et al., 2022), which is the SoTA SSI approach. It is a fully unsupervised SSI approach that leverages attention scores between token spans estimated using a pretrained language model to extract value candidates. A three-step hierarchical clustering procedure is then used that aims to cluster value types, then domains, then slots, using HDBSCAN.

**Results** As shown in Table 2, GenDSI outperforms the previous SoTA USI on almost every aspect of the SSI task. It contains fewer redundant slot clusters, superior recall of gold slots, higher

cluster purity as measured by value precision, and better coverage of gold slot values. The only metric on which GenDSI did not outperform USI is slot precision on the MultiWOZ evaluation. This is because the state generator model learned to predict boolean slot values from the DØT dataset that represent intent types, such as greeting and requesting information, which are considered as precision errors under this evaluation since gold slots do not encode intent classes. The performance of GenDSI -slot names dropped considerably on all metrics other than slot count, indicating the utility of inferring concrete slot names when discovering value candidates. Surprisingly, GenDSI +all domains did not afford any meaningful benefit, which may indicate that our approach generalizes to new domains without the need to generate in-domain resources.

**Slot Name Evaluation** Our SSI approach is the first to enable automatic naming of slot clusters. Simply labeling each cluster with the most frequent candidate slot name achieves 93.5% correctly named clusters by human evaluation.

## 5 Conclusion

This work presents a new SoTA for SSI, demonstrating the power of a generative approach to value candidate discovery. Our dialogue state generator model shows excellent performance for discovering key dialogue state information from unlabeled dialogues without any prior knowledge of the task domain. Its ability to label discovered values with appropriate slot names provides rich type information, allowing a simple clustering method to induce a quality slot schema for unseen domains. Despite this advancement, there is still room to improve SSI. In particular, current SSI methods produce a far greater number of induced slots compared to human-defined schemas. Although our approach reduces the number of induced slots somewhat, future work should aim for SSI with minimal redundancies in induced slots to further improve the utility of SSI in practical settings.

## Acknowledgments

We gratefully acknowledge the support of the Amazon Alexa AI grant. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the Alexa AI.

## References

- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - A Large-Scale Multi-Domain Wizard-of-Oz Dataset for Task-Oriented Dialogue Modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Yun-Nung Chen, William Yang Wang, and Alexander I Rudnicky. 2013. Unsupervised induction and filling of semantic slots for spoken dialogue systems using frame-semantic parsing. In *2013 IEEE Workshop on Automatic Speech Recognition and Understanding*, pages 120–125. IEEE.
- 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.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. [MultiWOZ 2.1: A Consolidated Multi-Domain Dialogue Dataset with State Corrections and State Tracking Baselines](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 422–428, Marseille, France. European Language Resources Association.
- James D. Finch and Jinho D. Choi. 2024. [Diverse and effective synthetic data generation for adaptable zero-shot dialogue state tracking](#). *Preprint*, arXiv:2405.12468.
- Vojtěch Hudeček, Ondřej Dušek, and Zhou Yu. 2021. [Discovering Dialogue Slots with Weak Supervision](#). 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 2430–2442, Online. Association for Computational Linguistics.
- Harrison Lee, Raghav Gupta, Abhinav Rastogi, Yuan Cao, Bin Zhang, and Yonghui Wu. 2022. [SGD-X: A Benchmark for Robust Generalization in Schema-Guided Dialogue Systems](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):10938–10946. Number: 10.
- Leland McInnes, John Healy, and Steve Astels. 2017. [hdbscan: Hierarchical density based clustering](#). *Journal of Open Source Software*, 2(11):205.
- Qingkai Min, Libo Qin, Zhiyang Teng, Xiao Liu, and Yue Zhang. 2020. Dialogue state induction using

neural latent variable models. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI’20*, pages 3845–3852, Yokohama, Yokohama, Japan.

Liang Qiu, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. [Structure Extraction in Task-Oriented Dialogues with Slot Clustering](#). *arXiv preprint*. ArXiv:2203.00073 [cs].

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.

Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. [Towards Scalable Multi-Domain Conversational Agents: The Schema-Guided Dialogue Dataset](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8689–8696. Number: 05.

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 and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.

Yuxia Wu, Lizi Liao, Xueming Qian, and Tat-Seng Chua. 2022. [Semi-supervised New Slot Discovery with Incremental Clustering](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6207–6218, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Dian Yu, Mingqiu Wang, Yuan Cao, Izhak Shafran, Laurent Shafey, and Hagen Soltau. 2022. [Unsupervised Slot Schema Induction for Task-oriented Dialog](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1174–1193, Seattle, United States. Association for Computational Linguistics.

## A State Generator Evaluation Details

To facilitate a thorough evaluation of dialogue state generators, a human evaluation measures the following two key aspects:

**State Update Completeness** measures the proportion of predicted state updates that humans have judged to fully capture the key information in their associated turns. Human judges are asked to read each turn within its context and make a binary decision on whether or not *any* essential information is missing in the state update such that:

$$CP = \frac{1}{|\mathcal{U}|} \sum_{\forall U \in \mathcal{U}} \mathbb{I}(\mathbf{complete}(U))$$

$\mathcal{U}$  is a list of all state updates across dialogues to be evaluated and  $\mathbb{I}(x)$  is 1 if  $x$  is true; otherwise, 0. Note that the judges are not responsible for finding *all* missing information but identifying at least one to assess completeness for efficient evaluation.

**Slot Value Correctness** measures the proportion of slot-value pairs that humans have judged to accurately represent *specific* information in their corresponding turns. Judges are asked to mark each slot-value pair as correct if it makes sense and is entirely faithful to the content of the associated turn s.t.:

$$CR = \frac{1}{\sum_{\forall U \in \mathcal{U}} |U|} \sum_{\forall U \in \mathcal{U}} \sum_{\forall (s,v) \in U} \mathbb{I}(\mathbf{correct}(s,v))$$

Note that both the slot name  $s$  and value  $v$  must be accurate for  $\mathbb{I}(\mathbf{correct}(s,v))$  to be 1.

## B SSI Evaluation Metrics

An SSI model produces a list of slot clusters  $\hat{S} = [\hat{s}_1, \hat{s}_2, \dots, \hat{s}_n]$  where each slot  $\hat{s}_i$  is a cluster of values  $\hat{s}_i = [\hat{v}_1, \hat{v}_2, \dots, \hat{v}_{|\hat{s}_i|}]$ . The quality of these slot clusters is measured by matching them against a list of gold reference slots  $S = [s_1, s_2, \dots, s_m]$ , each of which can be similarly represented as a list of the gold value labels that fill each slot such that  $s_i = [v_1, v_2, \dots, v_{|s_i|}]$ .

Matching is performed by assigning each induced slot  $\hat{s}_i$  to one or zero gold slots, creating a mapping  $M : \hat{S} \rightarrow S \oplus [\text{none}]$ . This matching is performed automatically. First, a centroid representation  $c_i$  is computed for each induced slot cluster and each gold slot cluster using the average BERT (Devlin et al., 2019) encoding of each value:

$$c_i = \frac{\sum_{v_j \in s_i} \text{BERT}(v_j)}{|s_i|}$$

Each induced cluster is mapped to the gold cluster whose centroid is closest by cosine distance, or to none if no gold centroid is within  $\geq 0.8$  cosine similarity.

Given the mapping  $M$  from predicted to gold slots, the evaluation metrics are calculated follows:

**Slot Precision** is the proportion of predicted slots that were able to be mapped to a gold slot:

$$\text{SP} = \frac{\sum_{\hat{s}_i \in \hat{S}} 1_S(M(\hat{s}_i))}{|\hat{S}|}$$

**Slot Recall** is the proportion of gold slots for which there is at least one corresponding predicted slot:

$$\text{SR} = \frac{|\{M(\hat{s}_i) : \hat{s}_i \in \hat{S}\} - \{\text{none}\}|}{|S|}$$

**Slot F1** is calculated normally as the harmonic mean of precision and recall:

$$\text{S-F1} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

**Slot Count** In the above Slot Precision calculation, multiple predicted clusters are allowed to be mapped to a single gold slot. This choice of formulation was made by previous work to avoid punishing the schema induction approach for inducing a finer-grained schema than what the gold schema provides, but fails to reflect the number of redundant clusters that are induced. To mitigate this, the number of induced slots is reported as an additional evaluation metric, where a lower number of induced slots is considered preferable.

**Value Precision** is meant to measure the purity of predicted slot clusters. It is calculated only between matched predicted clusters  $\hat{S}_{matched}$  and matched gold clusters  $S_{matched}$ . For each gold slot with at least one match  $s_i \in S_{matched}$ , the proportion of predicted values in the mapped predicted slots that have a fuzzy match to some gold slot value is measured using fuzzy match boolean function  $f$ :

$$\text{VP}_{s_i} = \frac{|\{\hat{v}_{kl} : \hat{v}_{kl} \in \hat{v}_k, M(\hat{s}_k) = s_i, v_{ij} \in s_i, f(v_{ij}, \hat{v}_{kl})\}|}{|\{\hat{v}_{kl} : \hat{v}_{kl} \in \hat{v}_k, M(\hat{s}_k) = s_i\}|}$$

The final Value Precision score is an average across matched gold slots calculated in this way:

$$\text{VP} = \frac{\sum_{s_i \in S_{matched}} \text{VP}_{s_i}}{|S_{matched}|}$$

**Value Recall** is calculated similarly to Value Precision. For each gold slot with a mapping to one or more predicted clusters, recall is measured as the proportion of gold values that have a fuzzy match to some value in the corresponding predicted clusters:

$$\text{VR}_{s_i} = \frac{|\{v_{ij} : \hat{v}_{kl} \in \hat{v}_k, M(\hat{s}_k) = s_i, v_{ij} \in s_i, f(v_{ij}, \hat{v}_{kl})\}|}{|s_i|}$$

The final Value Recall is also averaged across matched gold slots:

$$\text{VR} = \frac{\sum_{s_i \in S_{matched}} \text{VR}_{s_i}}{|S_{matched}|}$$

## C State Generator Evaluation Interface

Figure 2 shows a screenshot of the interface when performing completeness annotations, and Figure 3 shows a screenshot of the interface when performing correctness annotations. Note that the application interface relies on custom keybindings (e.g. pressing the y or n keys to indicate “yes” or “no”) for annotators to record their evaluation judgements.

## D Implementation Details

**Dialogue State Generator** All dialogue state generator models were trained using the original version of T5-3B using the huggingface transformers library<sup>2</sup>. All training was performed using a learning rate of  $1e - 4$ , weight decay of  $5e - 3$ , batch size 128, and for exactly 1 epoch, using the Adam optimizer.

**Slot Schema Induction** All SSI models used a T5-3B dialogue state generator model trained with the configuration presented above. The all-MiniLM-L6-v2 model from SentenceTransformers<sup>3</sup> was used for slot-value encoding. All HDBSCAN runs used the CUMML<sup>4</sup> library with a min. samples of 5, minimum cluster size 25, and cluster merge epsilon 0.3.

<sup>2</sup><https://huggingface.co/docs/transformers>

<sup>3</sup><https://www.sbert.net>

<sup>4</sup><https://docs.rapids.ai/api/cuml/stable/>

Instructions	Dialogue History
<p>Your objective is to evaluate the quality of the dialogue state update for the last dialogue turn.</p>	<p>I would like to make a private transaction with Jerry for 169 bucks.</p>
<p>The dialogue state update is organized into slots (types of information) and values (specific instances of these types of information), like so:</p>	<p>Will the transfer come from the app balance or credit card?</p>
<pre>flight destination:   New York City flight date:   January 4 flight available:   yes</pre>	<p>Please send the money from my mastercard.</p>
<p>To evaluate dialogue state updates, direct your attention to the dialogue history in the middle pane, focusing on the most recent dialogue turn and corresponding dialogue state update.</p>	<p>Please confirm you want me to make a private transfer to Jerry from your credit card in the amount of \$169.</p>
<p>The Task pane on the bottom left will prompt you with a yes/no question about the most recent turn: answer this prompt with a YES/ACCEPT answer by pressing [a] on your keyboard, or answer with NO/REJECT by pressing [r] on your keyboard</p>	<pre>recipient:   Jerry transfer type:   Private money source:   Credit card transfer amount:   \$169 confirm transfer:   ? ----- Complete?</pre>
<p>Given a dialogue turn, a complete dialogue state update is one that covers all key information shared in the turn. To decide whether the state update is complete:</p>	
<ol style="list-style-type: none"> <li>1. Identify all key information shared in the turn-- key information represents information that is necessary for the listener to understand in order for the dialogue to be successful.</li> <li>2. For each piece of key information, check whether the dialogue state update contains a slot-value pair that covers the key information.</li> <li>3. All key information must be covered by the slot-value pairs for the dialogue state update to be complete.</li> </ol>	
<p>Note that if the dialogue state update contains ADDITIONAL information that is irrelevant or incorrect, it does not affect the completeness of the dialogue state update: in other words, a dialogue state update can be complete even if it contains redundant or incorrect information, as long as the key information</p>	

Figure 2: Annotation interface with instructions for human evaluation of Completeness of predicted state updates.

Instructions	Dialogue History
<p>Your objective is to evaluate the quality of the dialogue state update for the last dialogue turn.</p>	<p>I would like to make a private transaction with Jerry for 169 bucks.</p>
<p>The dialogue state update is organized into slots (types of information) and values (specific instances of these types of information), like so:</p>	<p>Will the transfer come from the app balance or credit card?</p>
<pre>flight destination:   New York City flight date:   January 4 flight available:   yes</pre>	<p>Please send the money from my mastercard.</p>
<p>To evaluate dialogue state updates, direct your attention to the dialogue history in the middle pane, focusing on the most recent dialogue turn and corresponding dialogue state update.</p>	<p>Please confirm you want me to make a private transfer to Jerry from your credit card in the amount of \$169.</p>
<p>The Task pane on the bottom left will prompt you with a yes/no question about the most recent turn: answer this prompt with a</p>	<pre>recipient:   Jerry ----- Correct?</pre>
<p>Given a dialogue turn, a correct slot-value pair in the dialogue state update is one where:</p>	
<ul style="list-style-type: none"> <li>* the slot name represents a type of information that is relevant to the current turn</li> <li>* the value represents information mentioned or strongly implied by the turn (values representing information that is ONLY mentioned or implied in PREVIOUS turns are not correct)</li> <li>* if the value is "?", the slot information type was requested (implicitly or explicitly) by the turn's speaker</li> </ul>	
<p>Is the slot and value shown below an accurate extraction of information shared in the last dialogue turn?</p>	
<pre>recipient:   Jerry</pre>	

Figure 3: Annotation interface with instructions for human evaluation of Correctness of predicted slot-value pairs.

# Using Respiration for Enhancing Human-Robot Dialogue

Takao Obi and Kotaro Funakoshi

Tokyo Institute of Technology, Tokyo, Japan  
{smalltail, fukanoshi}@lr.pi.titech.ac.jp

## Abstract

This paper presents the development and capabilities of a spoken dialogue robot that uses respiration to enhance human-robot dialogue. By employing a respiratory estimation technique that uses video input, the dialogue robot captures user respiratory information during dialogue. This information is then used to prevent speech collisions between the user and the robot and to present synchronized pseudo-respiration with the user, thereby enhancing the smoothness and engagement of human-robot dialogue.

## 1 Introduction

For spoken dialogue robots to be effectively used in various scenarios, it is crucial for human-robot dialogue to be as natural as human-human dialogue. In human-human dialogue, communication occurs not only through verbal language but also through non-verbal cues. Thus, incorporating non-verbal information is essential for enhancing the naturalness of human-robot dialogue. Previous research has shown that integrating non-verbal cues such as nodding and body movements into robots can improve dialogue fluency (Watanabe et al., 2002), confirming the benefits of including non-verbal information in dialogue robots.

Our focus is on a specific type of non-verbal information: respiration. We believe that integrating respiratory information can significantly enhance human-robot dialogue, as respiration is intimately connected to speech. Our research has demonstrated that respiratory information is effective in predicting user speech onset (Obi and Funakoshi, 2023). Based on this finding, we developed a spoken dialogue system that uses user respiratory information to predict user speech onset, helping to prevent speech collisions in human-robot dialogue (Obi and Funakoshi, 2024). Furthermore, to ensure accurate capture of user respiratory informa-

tion in real-world dialogue settings, we have implemented a respiratory estimation method that uses video input. This method employs the first-ever deep learning model to provide a robust estimation of the respiratory waveform, even in the presence of speech movements, marking an improvement over an existing method (Obi and Funakoshi, 2023).

Building on these developments, we created a spoken dialogue robot that uses respiration. The dialogue robot estimates user respiratory waveform values and uses them for enhancing human-robot dialogue. Initially, the dialogue robot predicts user speech onset and initiates dialogue responses only when user utterances are not predicted, thereby preventing speech collisions. This approach is designed to facilitate smoother human-robot dialogue by ensuring that the conversation flows without interruptions. A previous study in human-robot dialogue has confirmed that avoiding speech collisions contributes to smoother turn-taking (Funakoshi et al., 2008). Secondly, the dialogue robot presents synchronized pseudo-respiration with user respiration. This approach is designed to enhance the dialogue robot's impression. Research in robots has shown that the presentation of pseudo-respiration can enhance the robot's impression (Terzioğlu et al., 2020). Furthermore, a study in human-human dialogue has demonstrated that synchronized respiration during turn-taking leads to smoother transitions between speakers (Rochet-Capellan and Fuchs, 2014).

To confirm the impact of these approaches on human-robot dialogue, we conducted a dialogue experiment in which 50 participants each interacted individually with the dialogue robot. We used actual respiratory waveform values obtained from a respiratory measurement device as user respiration (Obi and Funakoshi, 2024). Preliminary analysis, conducted after the initial report, indicates that adjusting the timing of robot speech onset using user speech prediction effectively reduces speech

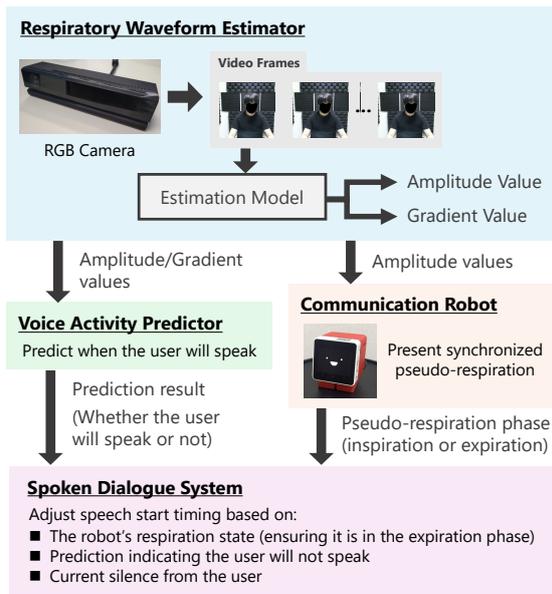


Figure 1: Overview of spoken dialogue robot using respiration

collisions. Additionally, initial user impression evaluations suggest that both the robot speech adjustment and synchronized pseudo-respiration presentation make users feel that their speech is less likely to overlap with the robot’s responses. These findings suggest that when the respiratory waveform estimation operates ideally, it can facilitate smooth human-robot dialogue.

## 2 System Overview

Our dialogue robot comprises a respiratory waveform estimator, a voice activity predictor, a communication robot, and a spoken dialogue system. Figure 1 provides an overview of these components and their arrangement.

**Respiratory Waveform Estimator:** We developed a respiratory waveform estimator using a deep-learning model comprising 3DCNN-ConvLSTM, which is robust against speaker motion (Obi and Funakoshi, 2023). This estimator uses RGB video frames of a user as input to estimate the user’s respiratory waveform amplitude and gradient at the time of the final frame. The model was trained using VRWiDataset<sup>1</sup>. The estimated values are then transmitted to both the voice activity predictor and the communication robot.

**Voice Activity Predictor:** We developed a voice activity predictor using a single-layer Long Short-Term Memory (LSTM) network that processes es-

timated respiratory waveform values as input. This model predicts whether user speech will occur within the next 200 ms during non-speaking periods. It was trained on a dataset created using the VRWiDataset, which extracts data from user non-speaking intervals. This dataset pairs the respiratory waveform values over a specific period with the user’s voice activity occurring 200 ms later. The prediction results are then transmitted to the spoken dialogue system.

**Communication Robot:** We use an open-source robot named stack-chan<sup>2</sup> for the communication robot. The communication robot performs a pseudo-respiratory movement, represented by its vertical motion. The movement is based on the user’s respiratory waveform amplitude values obtained from the respiratory waveform estimator. The communication robot uses these values to synchronize the timing of its inspiration and expiration with the user. Additionally, the communication robot’s inspiration/expiration phase information is transmitted to the spoken dialogue system to determine the speech timing.

**Spoken Dialogue System:** We developed a spoken dialogue system facilitating dialogues on arbitrary topics. This system uses GPT-4 Turbo<sup>3</sup> for the generation of dialogue responses. For speech processing, it employs both Google Cloud speech-to-text<sup>4</sup> and say command in macOS. The system initiates responses when the voice activity predictor confirms no imminent user speech onset, and only during the communication robot’s expiration phases. Additionally, it is designed not to respond during user speaking, ensuring that there are no response overlaps between the user and the dialogue robot in dialogue. If no user speech is detected, it will autonomously initiate responses continuously to maintain dialogue. The intervals between the continuous responses are set randomly between 0.5 and 3.5 seconds to simulate a realistic dialogue pace.

## 3 Use Case

Our dialogue robot is designed to demonstrate the effectiveness of integrating respiratory information into human-robot dialogue.

<sup>1</sup><https://github.com/fnkslab/VRWiDataset>

<sup>2</sup><https://github.com/meganetaaan/stack-chan>

<sup>3</sup><https://openai.com/gpt-4>

<sup>4</sup><https://cloud.google.com/speech-to-text>

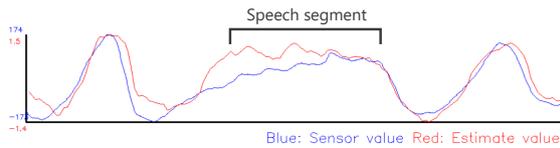


Figure 2: Example of real-time respiratory waveform comparison’s window with added speech segment

### 3.1 Scenario

Initially, a user sits in a chair positioned so that the upper body is visible to the camera, which captures the movements of the chest and abdomen associated with respiration. The respiratory waveform estimator continuously uses the captured video frames to estimate the user’s respiratory waveform values. As each estimation completes, the estimated values are immediately and continuously transmitted in real-time to both the voice activity predictor and the communication robot. During this process, the estimated waveform can be plotted and visually verified, enabling real-time confirmation of the accuracy of the estimations. Once the estimator begins transmitting the estimated values, both the voice activity predictor’s prediction and the communication robot’s pseudo-respiration presentation are initiated. These components start sending data to the spoken dialogue system simultaneously with their activation. Upon receiving these values, the system initiates a greeting, beginning the dialogue with the user. The system engages in dialogues on a variety of topics, capturing user speech through a microphone and considering it to generate contextually relevant responses.

### 3.2 Advanced Validation Features

The dialogue robot is equipped with various features to explore the effectiveness of using respiratory information.

**Real-time Waveform Comparison:** A user can attach a respiratory measurement device to their upper body, enabling real-time comparisons of actual waveform amplitudes with the estimated ones (Figure 2). Since the estimated waveforms and the actual waveforms have different ranges, they are displayed overlaid in a manner that aligns them to the same scale for comparison. This feature enables the user to directly observe the accuracy with which the respiratory waveform estimator is able to capture the user’s respiratory waveform. This real-time feedback is crucial for validating the performance of the respiratory waveform estimator.

**Using Actual Respiratory Waveform:** The respiratory waveform estimator can also transmit actual respiratory waveform values obtained from a respiratory measurement device instead of the estimated ones. When the actual waveform values are transmitted, they are normalized to align with the scale of the estimated waveforms before transmission. Using these actual waveform values, we can verify the effectiveness of the dialogue robot in using the user respiratory information for human-robot dialogue, assuming that the waveform estimation is accurate.

**Customizing Pseudo-Respiration Modes:** The communication robot’s pseudo-respiration presentation features three distinct modes: synchronized, steady, and no-presentation. In the steady mode, the communication robot follows a consistent, internally generated waveform, presenting pseudo-respiration independent of the user respiration. In the no-presentation mode, the communication robot does not move, and the spoken dialogue system responds based solely on the input from the voice activity predictor and the current absence of user speech, without considering the communication robot’s respiratory phase. These options allow for a comprehensive evaluation of how respiratory synchronization and pseudo-respiration presentation affect human-robot dialogue.

**Options for Voice Activity Predictor:** The voice activity predictor offers a choice between using amplitude or gradient values as inputs. This feature enables to verify which input is more effective in real-world dialogue settings. In our experimental environment, using the estimated gradient values as inputs resulted in higher prediction accuracy than using the estimated amplitude ones (Obi and Funakoshi, 2023). Additionally, the predictor can be turned off, allowing one to observe the impact of its presence or absence on human-robot dialogue.

## 4 Conclusion and Future Work

In pursuit of facilitating smooth and engaging human-robot dialogue, we have developed a spoken dialogue robot that uses respiration. This dialogue robot employs a respiratory estimation method using video input to capture user respiratory information, which serves two primary purposes: predicting user speech onset to prevent speech collisions in dialogues, and presenting pseudo-respiration synchronized with the user’s respiration. These approaches are expected to enhance the smoothness

and engagement of human-robot dialogue.

While adjusting speech timing contributes to smoother human-robot dialogue, reducing the number of robot utterances could detract the naturalness of the dialogue. To address this concern, future work will focus on incorporating non-verbal cues such as gaze into the voice activity prediction model, aiming to enhance its accuracy and ensure the dialogue robot does not unnecessarily remain silent. Additionally, accurate capture of user respiration is essential for the prediction in natural dialogue, so we will also work on developing a more robust respiratory waveform estimation method. Furthermore, we aim to develop a pseudo-respiration presentation that considers robot utterances, preventing a decrease in robot utterances while maintaining pseudo-respiration. Through these enhancements, we aim to use respiratory information more effectively, achieving more natural human-robot dialogue.

## Acknowledgement

This work was supported by JSPS KAKENHI Grant Number JP22H04859. We thank Dr. Ludovico Minati, formerly with Tokyo Tech, for his help with the respiratory measurement device.

## References

- Kotaro Funakoshi, Kazuki Kobayashi, Mikio Nakano, Seiji Yamada, Yasuhiko Kitamura, and Hiroshi Tsujino. 2008. [Smoothing human-robot speech interactions by using a blinking-light as subtle expression](#). In *Proceedings of the 10th International Conference on Multimodal Interfaces, ICMI '08*, page 293–296. Association for Computing Machinery.
- Takao Obi and Kotaro Funakoshi. 2023. [Video-based respiratory waveform estimation in dialogue: A novel task and dataset for human-machine interaction](#). In *Proceedings of the 25th International Conference on Multimodal Interaction, ICMI '23*, page 649–660. Association for Computing Machinery.
- Takao Obi and Kotaro Funakoshi. 2024. [Respiration-enhanced human-robot communication](#). In *Companion of the 2024 ACM/IEEE International Conference on Human-Robot Interaction, HRI '24*, page 813–816, New York, NY, USA. Association for Computing Machinery.
- Amélie Rochet-Capellan and Susanne Fuchs. 2014. [Take a breath and take the turn: how breathing meets turns in spontaneous dialogue](#). *Philosophical Transactions of the Royal Society B: Biological Sciences*, 369(1658):20130399.
- Yunus Terzioğlu, Bilge Mutlu, and Erol Şahin. 2020. [Designing social cues for collaborative robots: The role of gaze and breathing in human-robot collaboration](#). In *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction, HRI '20*, page 343–357, New York, NY, USA. Association for Computing Machinery.
- T. Watanabe, R. Danbara, and M. Okubo. 2002. [Inter-actor: Speech-driven embodied interactive actor](#). In *Proceedings. 11th IEEE International Workshop on Robot and Human Interactive Communication*, pages 430–435.

# Interactive Dialogue Interface for Personalized News Article Comprehension

**Tomoya Higuchi**  
The University of  
Electro-Communications  
1-5-1, Chofugaoka, Chofu,  
Tokyo, Japan  
h2430109@edu.cc.uec.ac.jp

**Michimasa Inaba**  
The University of  
Electro-Communications  
1-5-1, Chofugaoka, Chofu,  
Tokyo, Japan  
m-inaba@uec.ac.jp

## Abstract

We developed an interface to explain news articles through dialogue by considering the user's comprehension level. The interface generates several pertinent questions based on the ongoing dialogue and news article, and users advance the conversation by selecting a question. Based on the user's selected questions, the interface estimates their comprehension level of the news article and adjusts the difficulty of the generated questions accordingly. This enables a personalized dialogue tailored to each user's comprehension needs. The results of the baseline comparison experiments confirmed the usefulness of the interface.

## 1 Introduction

Organization for Economic Co-operation and Development reported that 21% of individuals in the United States aged 16 years and over have low English literacy. This limitation results in them having difficulty reading longer texts such as news articles (Desjardins et al., 2013). In this paper, we construct an interactive news commentary dialogue interface that considers the user's comprehension level to increase the understandability of news articles.

We constructed an interface, shown in Figure 1, that interactively explains news articles. The interface generates several pertinent questions based on the ongoing dialogue and news article, and users advance the conversation by selecting a question. This interactive approach is expected to provide users the opportunity to organize information and deepen their understanding. When users select a question, the system predicts their comprehension level of the news article based on the selected question and adjusts the difficulty of the next question accordingly. This feature allows the dialogue to be tailored to each user's knowledge level.

The screenshot displays a dialogue interface with the following content:

- Title:** Brightest-ever cosmic explosion solved but new mysteries sparked
- System:** Researchers have identified the source of the brightest burst of light ever recorded, which originated from an exploding star observed in 2022.
- User:** What does it mean when a star explodes?
- System:** When a star explodes, it undergoes a supernova, a powerful and luminous stellar explosion. This event is significant because it is thought to be a key source of production for heavy elements in the universe, such as gold and platinum, although the recent findings challenge this understanding as no heavy elements were detected in the observed explosion.
- Candidate questions:**
  - Level 1 question: What are heavy elements, and why are they important?
  - Level 2 question: How does the recent finding of no heavy elements in the supernova challenge the existing theories about the origin of elements like gold?
  - Level 2 question: What are the implications of this discovery for our understanding of how precious metals are formed in the universe?

Figure 1: Proposed Dialogue Interface. The dialogue history, which includes the user-selected questions and the system's answers, is displayed at the top. The user selects a question from the candidate questions at the bottom. The difficulty level of the displayed question candidates varies from turn to turn based on the estimated user comprehension level.

## 2 Related Works

Media often use question-and-answer series about news stories because of their clarity, but reporters manually create them; thus, the process is inefficient due to high costs and time requirements. Laban et al. (2020) proposed a method to automatically build chatbots from news articles. In these chatbots, the system presents the user with candidate questions. Podcasts with automatic construction containing QA sessions on news were also analyzed, confirming the effectiveness of a conversational format (Laban et al., 2022). However, these questions are not customized to individual user characteristics, and the tailored dialogues are not generated for individual users. User characteristics, including their social group, influence the types of questions asked. Stewart and Mihalcea (2022) introduced a method for producing ques-

tions that consider these characteristics, training a text generation model on social media data that accounts for social groups and their expertise in specific domains. An et al. (2021) developed a conversational agent prototype that incorporates a strategy to account for the user’s knowledge and ignorance regarding speech, validating the importance of considering user knowledge. Drawing inspiration from these methodologies, we generate questions that incorporate the user’s comprehension level of the news story.

### 3 System Overview

We constructed an interactive interface to explain the contents of news articles, illustrated in Figure 1. The interface first presents the user with a brief introduction to the news article and three candidate questions. The user then selects one of these questions, and the system explains the article’s content by answering the selected question. This approach reduces the user’s burden by eliminating the need to think about questions independently. The purpose of this interface is to facilitate understanding by providing explanations tailored to individual users. Candidate questions should reflect those the user might want to ask, considering the user’s comprehension level. Therefore, we generate and present candidate questions, considering the user’s comprehension level. The system comprises three modules, as shown in Figure 2: an introduction generation module, a candidate question generation module, and an answer generation module. Each module uses GPT-4 (OpenAI et al., 2024) as the large language model.

#### 3.1 Introduction Generation

The introduction generation module create an introductory summary of the news article intended for the user. This process involves utilizing the news article’s content and providing guidelines within the prompt to craft a concise introduction in one sentence. GPT-4 generates this introduction in a zero-shot.

#### 3.2 Question Candidate Generation

The question candidate generation module creates a set of questions to present to the user. A dialogue tailored to individual users must pose questions that match the user’s comprehension level. Additionally, the questions should follow the dialogue naturally and be answerable based on the news. We use GPT-4 to overgenerate candidate questions that

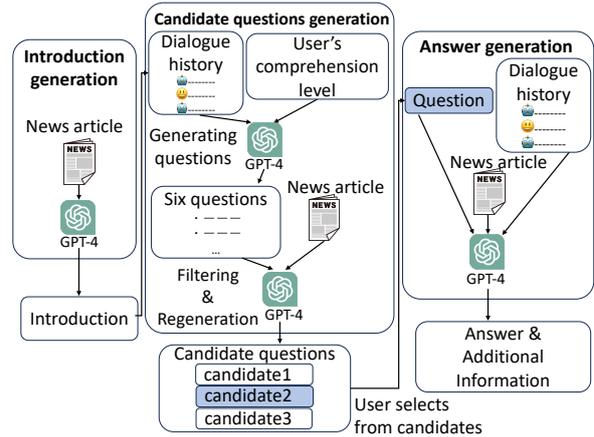


Figure 2: Architecture of Proposed system. The system comprises three modules: an introduction generation module that creates an introduction to the news article, a candidate question generation module that generates three candidate questions to present to the user, and an answer generation module that generates answers to the user’s chosen questions.

Table 1: The user’s comprehension level  $x$  and the number of candidate questions to present for each level

$x$	level 1	level 2	level 3
$x < 1.5$	1	2	0
$1.5 \leq x < 2$	1	1	1
$2 \leq x < 2.5$	0	2	1
$2.5 \leq x$	0	1	2

consider the user’s comprehension level from the dialogue history, filter them based on their answerability from the articles, and regenerate them if the required number is not fulfilled, and present three to the user.

#### 3.2.1 Candidate Question Generation

We categorize the levels of difficulty for the questions as 1, 2, and 3. Level 1 pertains to queries regarding the interpretation of terms, tailored for novices in the field; level 2 contains general inquiries about the article’s content; and level 3 comprises more intricate questions designed for field experts. The quantity of questions presented to the user for each difficulty level is modified based on Table 1, according to the user’s comprehension level estimated for each turn.

The average difficulty level of questions selected by the user up to the current turn is used to estimate the user’s comprehension level. Because the user’s comprehension level cannot be estimated in the first turn, candidate questions are presented to the user at difficulty levels 1, 2, and 3.

The difficulty level of the candidate questions is determined by referring to Table 1. The dialogue history and the difficulty condition of the candidate question are input into GPT-4, which generates a total of six candidate questions at difficulty levels 1-3, two for each level, that naturally follow the dialogue. Generating questions at levels 1 to 3 in every turn increases the likelihood that questions appropriate to each level are generated.

### 3.2.2 Filtering & Regeneration

When GPT-4 generated questions with a news article as input, the questions often anticipate information from the article not included in the dialogue history. This can prevent users from organizing information. Therefore, candidate questions are first overgenerated solely from the dialogue history, and only those questions that can be answered using the content of the news article are selected through filtering. If the number of candidate questions after filtering is less than required, the candidate questions are regenerated.

Of the six candidate questions generated, level 2 and level 3 candidate questions and the text of the news article are input into GPT-4, which determines whether they can be answered from the article’s content and filters them accordingly. No filtering for level 1 questions because they are designed to help users obtain prerequisite knowledge not in the content of the news article.

If the number of questions that can be answered from the content of a news article at each difficulty level is less than the number specified based on Table 1, GPT-4 generates the missing candidate questions by using the text of the news article and the dialogue history as input.

### 3.3 Answer Generation

The answer generation module generates responses to user-selected questions from candidates produced by the question candidate generation module. Inputs include the news article text, dialogue history, and questions, with GPT-4 generating answers and supplementary information to transition to the next topic in one shot. This approach helps prevent subsequent question candidates from focusing solely on one topic by providing context beyond just the answers.

## 4 Experiment

To evaluate the usefulness of this interface in reading news articles, we conducted a subject exper-

iment to compare it with three different baseline settings. Crowdworker read a randomly assigned article in a randomly assigned setting and completed a comprehension test and questionnaire.

### 4.1 Comparative Methods

As comparison methods, we conducted experiments using the following three settings:

#### 1. Reading news article

We adopted this baseline to compare interactive and non-interactive formats. The user reads a news article using a web browser.

#### 2. Microsoft Copilot<sup>1</sup>

We adopted this baseline to examine the effectiveness of presenting users with candidate questions. The user opens the news article in Microsoft Edge<sup>2</sup> and asks the Copilot a question without reading the body of the news article.

#### 3. W/o comprehension level

This baseline is compared with the proposed method to investigate the effectiveness of presenting comprehension-aware candidate questions. In this baseline interface, the system presents candidate questions at a single difficulty level to the user without considering their comprehension level. The difficulty level of the candidate questions corresponds to Level 2 in the proposed method.

### 4.2 Experimental Setup

We recruited 80 crowdworkers using Amazon Mechanical Turk<sup>3</sup> and conducted a subject experiment. Four news articles were selected from BBC news, one each from the fields of natural science<sup>4</sup>, technology<sup>5</sup>, international<sup>6</sup>, and finance/economics<sup>7</sup>. Each article was chosen based on the criteria that it was at least 500 words long and contained specialized content. Each crowdworker was randomly assigned one of four settings and one of the four news articles. After reading the article, the crowdworkers took a comprehension test on the content of the news article and completed a questionnaire.

<sup>1</sup><https://www.microsoft.com/ja-jp/microsoft-copilot>

<sup>2</sup><https://www.microsoft.com/ja-jp/edge?form=MA13FJ>

<sup>3</sup><https://www.mturk.com/>

<sup>4</sup><https://www.bbc.com/news/science-environment-68787534>

<sup>5</sup><https://www.bbc.com/news/business-68225115>

<sup>6</sup><https://www.bbc.com/news/world-us-canada-68883659>

<sup>7</sup><https://www.bbc.com/news/world-europe-68761491>

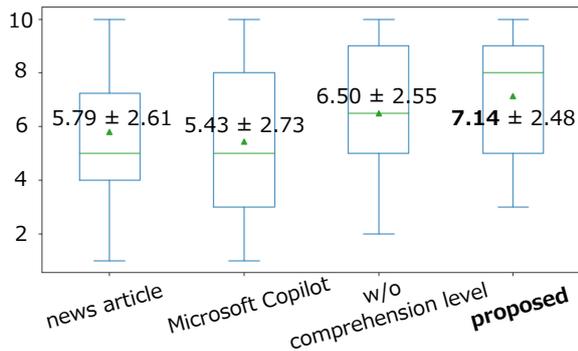


Figure 3: Boxplot of comprehension test scores.  $\triangle$  is the average score.

The comprehension test consisted of 10 four-choice questions, with one point per question. We collected a total of 80 data, five for each combination of setting and article.

### 4.3 Result

Figure 3 shows the aggregate results of the comprehension test. The average score was highest for the proposed interface. We can conclude that the proposed method, which presents candidate questions considering the comprehension level, is useful for reading comprehension. There are statistically significant differences in average comprehension test scores between the proposed interface and Microsoft Copilot ( $p < 0.05$ ) and no significant differences in the other combinations. Figure 3 shows that in the proposed interface, few people scored low on the comprehension test. This finding indicates that candidate questions considering the user’s comprehension level are effective for users with low comprehension.

According to the open-ended questionnaire, participants found w/o comprehension level interface and proposed interface to be fun and easy (e.g., "Fun to read, easy to use"). In the proposed interface, participants expressed favorable opinions about the candidate questions that considered comprehension level (e.g., "Great range of follow-up questions," "This will be sufficient for first-time readers"). However, some participants complained about the response time (e.g., "Was a bit slow to generate responses.")

## 5 Conclusion

We developed a news commentary dialogue interface that considers the user’s comprehension level. The interface alleviates the user’s burden by pre-

senting multiple automatically generated question candidates and having the system answer the selected questions. It aims to tailor the dialogue to each user by predicting their comprehension level of the news content on the basis of the selected questions and presenting question candidates that account for this level of understanding. In our experiment, we implemented the proposed interface and quantitatively evaluated its effectiveness by using crowd-sourcing. Comparison experiments between the proposed interface and baselines confirmed that the proposed interface enhances users’ reading comprehension.

In the open-ended responses to the questionnaire during the experiment, one participant said, "Conveying the content of news solely through text is difficult." Based on this feedback, our further research will consider incorporating images along with text.

## References

- Sungeun An, Robert Moore, Eric Young Liu, and Guang-Jie Ren. 2021. [Recipient design for conversational agents: Tailoring agent’s utterance to user’s knowledge](#). In *Proceedings of the 3rd Conference on Conversational User Interfaces, CUI ’21*.
- Richard Desjardins, William Thorn, Andreas Schleicher, Glenda Quintini, Michele Pellizzari, Viktoria Kis, and Ji Eun Chung. 2013. *Oecd skills outlook 2013: First results from the survey of adult skills*.
- Philippe Laban, John Canny, and Marti A. Hearst. 2020. [What’s the latest? a question-driven news chatbot](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 380–387.
- Philippe Laban, Elicia Ye, Srujay Korlakunta, John Canny, and Marti Hearst. 2022. [Newspod: Automatic and interactive news podcasts](#). In *Proceedings of the 27th International Conference on Intelligent User Interfaces, IUI ’22*, page 691–706.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, and Shyamal Anadkat. et al. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Ian Stewart and Rada Mihalcea. 2022. [How well do you know your audience? toward socially-aware question generation](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 255–269.

# Enhancing Dialogue Speech Recognition with Robust Contextual Awareness via Noise Representation Learning

Wonjun Lee <sup>\*1</sup>, San Kim <sup>\*2</sup> and Gary Geunbae Lee <sup>1,2</sup>

<sup>1</sup> Department of Computer Science and Engineering, POSTECH, Republic of Korea

<sup>2</sup> Graduate School of Artificial Intelligence, POSTECH, Republic of Korea

{lee1jun, sankm, gblee}@postech.ac.kr

## Abstract

Recent dialogue systems rely on turn-based spoken interactions, requiring accurate Automatic Speech Recognition (ASR). Errors in ASR can significantly impact downstream dialogue tasks. To address this, using dialogue context from user and agent interactions for transcribing subsequent utterances has been proposed. This method incorporates the transcription of the user's speech and the agent's response as model input, using the accumulated context generated by each turn. However, this context is susceptible to ASR errors because it is generated by the ASR model in an auto-regressive fashion. Such noisy context can further degrade the benefits of context input, resulting in suboptimal ASR performance. In this paper, we introduce Context Noise Representation Learning (CNRL) to enhance robustness against noisy context, ultimately improving dialogue speech recognition accuracy. To maximize the advantage of context awareness, our approach includes decoder pre-training using text-based dialogue data and noise representation learning for a context encoder. Based on the evaluation of speech dialogues, our method shows superior results compared to baselines. Furthermore, the strength of our approach is highlighted in noisy environments where user speech is barely audible due to real-world noise, relying on contextual information to transcribe the input accurately.

## 1 Introduction

Automatic Speech Recognition (ASR) is central in accurately interpreting human speech, serving as a fundamental resource for numerous subsequent downstream tasks. The advent of robust ASR modules, such as wav2vec2.0 (Baevski et al., 2020) and Whisper (Radford et al., 2023), has significantly enhanced the capabilities of ASR systems, facilitating their integration into a wide array of

research and application domains. The integration of ASR modules into various works highlights the pivotal role of ASR in enhancing human-computer interaction, signifying a notable development in interactive technologies.

Despite the successful advancement of the ASR system, its inaccuracy poses significant risks to the efficacy of downstream tasks, such as speech-to-text translation (Liu et al., 2020; Le et al., 2024; Tang et al., 2021) and spoken language understanding (Serdyuk et al., 2018; Arora et al., 2022; Huang and Chen, 2020). These tasks predominantly rely on the textual output generated by ASR systems, highlighting the importance of accuracy in the initial speech recognition process. Especially for the dialogue system, the quality of the ASR system is paramount to ensure seamless interaction between user and dialogue agent, as models trained on written conversations perform poorly on spoken data (Kim et al., 2021). To minimize the impact of ASR error on the dialogue model, various endeavors have been made. Jiang et al. (2023) used an ASR correction module which employs multiple ASR models, while others focused on augmenting data with plausible ASR errors (Park et al., 2023; Wang et al., 2020; Tian et al., 2021). However the limitation is evident as they primarily focus on the robustness of dialogue models, which may not address the core issue compared to directly rectifying ASR models.

Conversely, incorporating a context encoder for dialogue history to improve the ASR model has been proposed, resulting in notable performance enhancements (Ortiz and Burud, 2021; Shenoy et al., 2021; Hou et al., 2022; Hori et al., 2020). Nevertheless, since the context is transcribed at each turn by the ASR model, it may contain errors, potentially disrupting the use of contextual information.

In this work, we present a novel Context Noise Representation Learning (CNRL) method to encode accurate contextual information, even from

\*Equally contributed

noisy ASR transcriptions. This approach aims to improve the performance of speech recognition in Task Oriented Dialogue (TOD) by minimizing the impact of ASR errors in dialogue history as context. Furthermore, we explore the advantages of decoder pre-training in context-aware ASR systems, emphasizing their improved robustness in noisy environments. The overall training pipeline can be decomposed by three steps: 1) Decoder pre-training on text-based dialogue data between user and agent. 2) ASR fine-tuning with speech encoder and context encoder jointly. 3) CNRL on context encoder to minimize the impact of ASR-noise context. Our contributions are as follows:

- We propose a novel training pipeline for dialogue speech recognition that leverages the dialogue history between user and agent.
- We demonstrate the effectiveness of CNRL by comparing it to various baseline models, showing a relative 13% reduction in Word Error Rate (WER) compared to the current state-of-the-art ASR model (Radford et al., 2023).
- In evaluations conducted in highly noisy environments, our model exhibits robust transcription accuracy, achieving up to a 31.4% reduction in WER compared to the baseline.

## 2 Related work

### 2.1 Context-aware speech recognition

Several studies have shown that leveraging contextual information in dialogue scenarios can enhance ASR performance. Shenoy et al. (2021) used a context carry-over mechanism to enhance the recurrent model’s accuracy. Hou et al. (2022) proposed utilizing a context encoder in RNN-T architecture, adopting the semantic embedding of dialogue context from BERT (Devlin et al., 2019). Hori et al. (2020) targeted considering long-context by sliding-window fashion. Wang et al. (2023) and Wang et al. (2024) proposed an audio-augmented retriever to directly transcribe and track the dialogue state. These Context-Aware ASR(CA-ASR) models have a potential drawback: the context generated for each turn is based on ASR transcriptions, which inevitably contain errors, potentially degrading context-awareness. In this paper, we introduce the CNRL method, which trains only the context encoder independently. The goal is to enable the

context encoder to produce similar encoding for noisy (ASR output) contexts to match clean context.

### 2.2 Decoder pre-training

Compared to pre-training encoder layers (Baevski et al., 2020; Hsu et al., 2021; Chen et al., 2022), pre-training the decoder for ASR has received comparatively less attention. Notably, in scenarios where input speech is flawed or incomplete, the decoder can still play a crucial role in transcribing user utterances by leveraging contextual language modeling. To harness the decoder’s capabilities, the use of external datasets like phoneme-to-grapheme paired data (Masumura et al., 2020) or text data (Gao et al., 2021) has been suggested. This approach enables the model to benefit from numerous external, non-paired data sources. Tsunoo et al. (2023) trained decoder for both ASR task and language modeling task, enabling improved linguistic understanding and leading to better ASR performance. Following these works, we pre-trained the decoder for a context-aware ASR model using voluminous text-only data. Specifically, we focus on turn-based dialogue data between user and agent, where each utterance is highly related to each other.

### 2.3 Noise Representation Learning

Noise in input data is inevitable in various forms across many datasets. Training models with such data negatively impacts their generalization performance. To address this challenge, numerous studies have adopted contrastive learning to enhance model robustness. Ma et al. (2023) improved named entity recognition performance by employing a token-level dynamic loss function and contrastive learning, leveraging noisy data and accounting for noise-distribution changes during training. Xu et al. (2023) enhanced contrastive learning through a dimension-wise method to mitigate feature corruption in sentence embeddings. Sun et al. (2023) used a K-NN graph to identify confident samples and applied mixup supervised contrastive learning to create robust representations, leading to improved relation extraction performance. Zheng et al. (2023) utilized both class-wise and instance-wise contrastive learning in their novel representation learning module. In this work, we adopt representation learning to enhance context awareness when noisy ASR transcriptions are used for context. The proposed CNRL is integrated solely with the context encoder in the CA-ASR model to

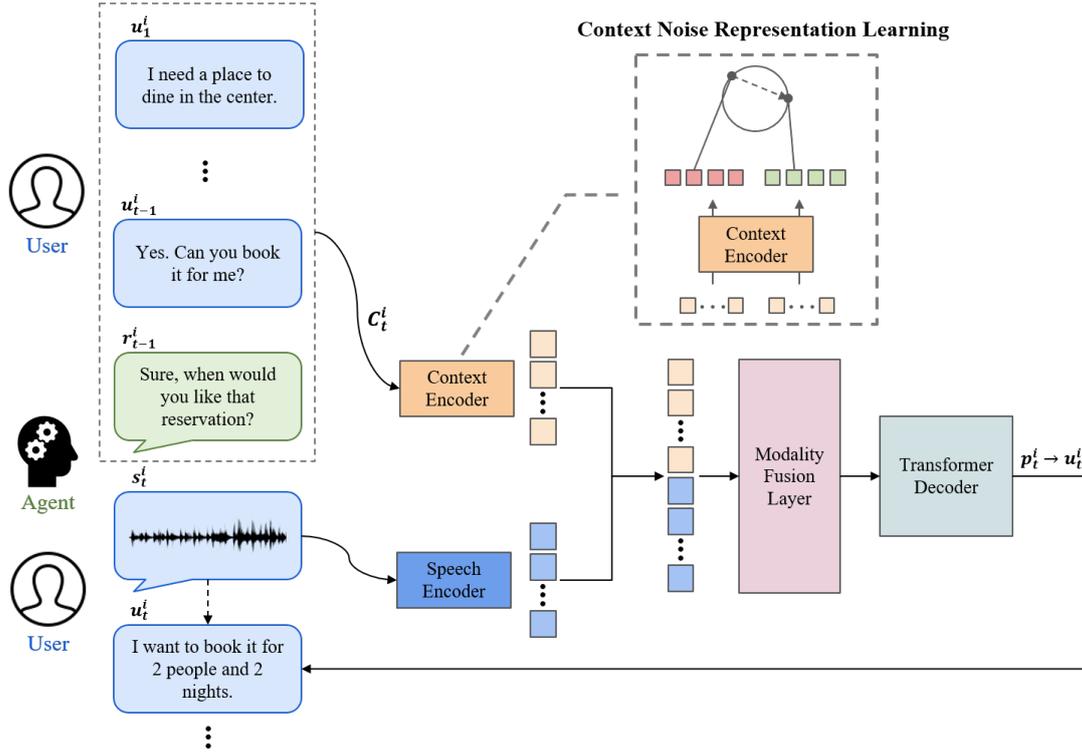


Figure 1: The architecture of a Context-Aware ASR(CA-ASR), featuring separate speech and context encoders to process the user’s current speech  $s_t^i$  and dialogue history  $C_t^i$ , respectively. These representations are concatenated and fused using a modality fusion layer and transcribed to the predicted user utterance  $p_t^i$  by the transformer decoder. The predicted user utterance will be added to context ( $p_t^i \rightarrow u_t^i$ ) for the next turn ( $t + 1$ ). After the training, the context encoder can improve itself by our CNRL method, detailed in Figure 2 and Section 3.3.

minimize training costs.

### 3 Methodology

#### 3.1 Preliminary

We define  $D_t^i$  as the turn-based dialogue dataset for turn  $t$  in the  $i$ -th dialogue, which includes the speech input  $s_t^i$ , the corresponding text labels  $u_t^i$  (transcriptions) of user utterances, and the turn-based dialogue history  $C_t^i = (u_1^i, r_1^i, \dots, u_{t-1}^i, r_{t-1}^i)$ , accumulating up to turn  $t - 1$ , where  $r_t^i$  represents the agent’s response at turn  $t$ . Each dialogue instance at the  $k$ -th turn, denoted as  $(u_k^i, r_k^i)$ , comprises a single-turn conversation consisting of both a user utterance and an agent response. During inference, the predicted utterance (transcription) from model  $p_t^i$  is used instead of  $u_t^i$  for user utterance to form context  $C_t^i$ .

The CA-ASR model integrates the user’s speech and dialogue history. For each turn  $t$ , the model predicts the current user utterance  $u_t^i$  from the speech input  $s_t^i$  and the context  $C_t^i$ . The dialogue history comprises text logs from both the user and the agent, where the user’s speech is transcribed in real-time, while the agent’s responses are given in

text format. To transcribe the user’s speech at turn  $t$ , the model draws upon past conversations from turn 1 to  $t - 1$ . Utilizing an encoder-decoder architecture for the CA-ASR model, dedicated encoders initially process each input type—speech and text. These encodings are then concatenated and fused through a modality fusion layer, yielding a fused representation. Subsequently, the fused representation is passed through a decoder layer to transcribe the user utterance. Figure 1 illustrates the CA-ASR architecture, highlighting the interaction between user utterances and agent responses.

#### 3.2 Decoder pre-training for Dialogue

We adopt a pre-training method specifically targeting decoders in the CA-ASR model. This method employs an encoder-decoder architecture, where the model takes the text-form dialogue history  $C_t^i$  as input. For the output, since the decoder is eventually used for transcribing user utterances, it aims to predict the next user utterance  $u_t^i$ . Additionally, the utilization of text data as input enables the training process to use external text datasets, further enhancing the decoder’s performance. We demonstrate

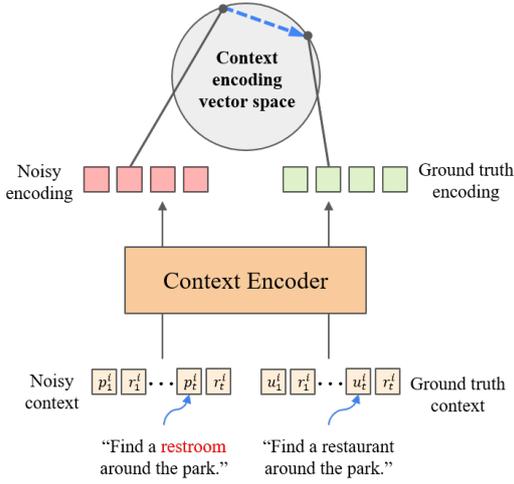


Figure 2: **Context Noise Representation Learning:** The noisy context including user utterances generated by the CA-ASR model during inference ( $p_t^i$ ), and the ground truth context with clean user utterances ( $u_t^i$ ), are encoded by the context encoder. The noisy encoding is adjusted to closely match the ground truth encoding in the context encoding vector space.

this efficacy in Section 5.2 . This approach enables the decoder to anticipate the subsequent user utterance based on contextual information derived from the dialogue history. This training method is particularly effective because dialogues in TOD are more predictable from the dialogue history than other types of conversations. In typical user-agent interactions, the agent often asks specific questions, and the user responds with relevant answers, making the dialogue structure more consistent and easier to predict.

When integrated into the CA-ASR model and fine-tuned for ASR tasks, the pre-trained decoder can significantly enhance transcription performance. By leveraging its ability to anticipate user responses from the agent’s response (or the entire dialogue history), the decoder contributes to more accurate and robust transcription results, even with imperfect input speech, such as noisy audio signals.

### 3.3 Context Noise Representation Learning

During inference, the CA-ASR model uses context from previous transcriptions of user utterances and agent responses. However, inaccuracies in the ASR-generated transcriptions can degrade the advantage of using context, as training typically uses only ground truth context for each turn. To address this, we introduce CNRL. This method involves an additional training step where the model tran-

scribes and utilizes noisy transcriptions to train the context encoder in a representation learning manner, as illustrated in Figure 2. The context encoder is fine-tuned to generate similar encoding for noisy input context as it does for the ground truth context. This method focuses solely on enhancing the context encoder, maintaining training efficiency.

To create the training set for CNRL, we first generate noisy transcriptions using the CA-ASR model with the ASR training set (See Section 4.1) divided into 10 folds. In each fold, 90% of the training set is used to train the CA-ASR model, and the remaining 10% is used to generate noisy ASR transcriptions. By iterating through all 10 folds, we obtain a complete noisy context training set. The dataset for CNRL comprises pairs of noisy and ground truth contexts, each containing multiple conversation turns. Each turn pairs a user utterance with an agent response, except for the initial turn, which consists only of the user’s utterance.

We trained context encoder with cosine embedding loss:

$$\text{loss}(x, y) = \begin{cases} 1 - \cos(x_1, x_2), & \text{if } y = 1 \\ \max(0, \cos(x_1, x_2) - \text{margin}), & \text{if } y = -1 \end{cases} \quad (1)$$

Where  $x_1$  is the encoding vector from the context encoder within the ASR-generated context and  $x_2$  is the encoding from ground truth context.  $y$  is the label that indicates these two ( $x_1$  and  $x_2$ ) are of the same class ( $y = 1$ ) or not ( $y = -1$ ). Since we trained the context encoder to generate a similar output encoding for the noisy input ( $x_1$ ) to match the clean ground truth ( $x_2$ ), we set  $y = 1$  for training. During training,  $x_1$  gets close to  $x_2$  on context encoding vector space, ensuring the context encoder produces similar encoding for a given noisy context. By using CNRL, the context encoder can maintain accurate context information, leading to improved speech recognition accuracy.

## 4 Experimental setup

### 4.1 Datasets

**The DSTC11 Challenge Dataset** The DSTC11 (Soltau et al., 2022) dataset is derived from the MultiWoZ 2.1 (Eric et al., 2020) by adding speech recordings and synthesized voices generated by a TTS model. The training set is built using the TTS model, while the evaluation sets are recorded by human volunteers. Each dialogue consists of audio files of user utterances paired with corre-

sponding agent responses. In every dialogue, the user initiates the conversation, making the first user utterance has no preceding context.

Since the official transcription for the DSTC11 test split (test-dstc11.human-verbatim) is not publicly available, we evaluate our experiments on the DSTC11 development split with human recording (dev-dstc11.human-verbatim)<sup>1</sup> as test set. Additionally we randomly sampled 3000 audios from the training set and used them as our development set during training.

The DSTC11 training set consists of 8,434 dialogues comprising 56,750 user utterances synthesized by four TTS voices, generating a total of 227,000 audio files. Our development set, randomly sampled from the training set, contains 3000 user utterances and is excluded from the training data. The test set includes human recordings of 7,374 user utterances from 1,000 dialogues. The average audio duration is 3.31 seconds for the training and development sets and 5.35 seconds for the test set.

**Evaluation in Noisy Environments** Environmental noise is a significant challenge for ASR systems in real-world scenarios. However, contextual information can mitigate this issue. To test our ASR system’s resilience to real-world noises, we use the ESC-50 dataset (Piczak, 2015), which includes 50 classes of common urban noises, such as drilling and sirens. Noise samples are randomly selected from 2000 audio files and injected into our test set at Signal-to-Noise Ratios (SNR) of 20dB and 0dB, representing soft and hard noise conditions, respectively. This evaluation replicates challenging acoustic environments to test the ASR system’s robustness rigorously. Note that the noisy audio is used exclusively for evaluation, not training. Our goal is to show that contextual information can be helpful in noisy environments where the audio signal is significantly degraded.

**Decoder pre-training** To facilitate the use of context information, we first trained CA-ASR’s decoder using exclusively text-based data before ASR fine-tuning. For this purpose, we employ large datasets of turn-based dialogue text, combining the Schema-Guided Dialogue (SGD) (Rastogi et al., 2020) dataset with the DSTC11 text dataset to pre-train the decoder. SGD consists of over 20,000 task-oriented conversations between human and

virtual assistant. From 8434 English dialogues from DSTC11 and approximately 16,000 English dialogues from the SGD training dataset, we use about 260,000 turn conversations. To evaluate the effect of decoder pre-training, we varied the volume of text data used for this process. The effects of these variations are detailed in Table 2.

## 4.2 Model configuration

**Baselines** We compare our CA-ASR model against several baselines, including those reported in DSTC11 (Soltau et al., 2022) and the current state-of-the-art ASR model Whisper (Radford et al., 2023). Additionally, we present a model that uses wav2vec2.0 (Baevski et al., 2020) as the encoder and BART (Lewis et al., 2019) as the decoder. This model shares the same architecture as the CA-ASR model, except for removing the context encoder and modality fusion. For transcription post-processing, we normalize common English patterns (e.g., "I've" to "I have"), remove punctuation, and normalize digits to ensure a fair comparison between models.

**Context-Aware ASR** Compared to the baselines, the CA-ASR model leverages previous user utterances and agent responses as textual input to enhance transcription accuracy. To encode this contextual information, CA-ASR uses the BART encoder as the context encoder. The speech encoder is wav2vec2.0 with the checkpoint *wav2vec2-large-960h*<sup>2</sup>, and the pretrained BART encoder and decoder with the checkpoint *bart-large*<sup>3</sup> are utilized as the context encoder and the CA-ASR decoder, respectively. Given that the maximum token length for BART-large is limited to 1024, we truncate the context to the last 1024 tokens if necessary.

For modality fusion, the wav2vec2.0 speech encoder and the BART context encoder each produce hidden representations with dimensions of token  $\times$  1024. Since the BART decoder requires an encoder hidden state with a dimension of 1024, we concatenate these hidden representations along the 1024 dimension. This concatenated representation is then passed through a linear layer (1024, 1024) with ReLU activation to create a fused representation. This fused representation is subsequently fed into the BART decoder to transcribe the user utterance.

Total parameter size of our model is 774M, consisting of 315M for the speech encoder, 203M for

<sup>1</sup>[https://storage.googleapis.com/gresearch/dstc11/dstc11\\_20221102a.html](https://storage.googleapis.com/gresearch/dstc11/dstc11_20221102a.html)

<sup>2</sup><https://huggingface.co/facebook/wav2vec2-large-960h>

<sup>3</sup><https://huggingface.co/facebook/bart-large>

the BART context encoder, 254M for the BART decoder, and 1M for the linear fusion layer.

### 4.3 Training configuration

Our training pipeline consists of three sequential steps: decoder pretraining, ASR fine-tuning with audio masking, and CNRL. We evaluate the effect of each step in the subsequent Result & Analysis section.

**Decoder pre-training** We initially adopt the BART encoder-decoder model to pre-train the decoder, which is subsequently used for ASR fine-tuning. The optimization is performed using the AdamW algorithm (Loshchilov and Hutter, 2017) with  $(\beta_1, \beta_2) = (0.9, 0.999)$ , learning rate of  $5e-5$ , weight decay of  $1e-5$ , and a batch size of 32. We select the best model based on the lowest validation loss over 10 epochs of training, spanning 50 hours. The encoder functions as the context encoder, while the decoder serves as the transformer decoder in the CA-ASR model. Utilizing Cross-Entropy loss, we aim to input the dialogue history with the agent’s response, which is the last turn, into the encoder and generate the user’s response as the output from the decoder.

**ASR fine-tuning** In ASR fine-tuning stage, a speech encoder (wav2vec2.0) is attached to the pre-trained BART decoder from decoder pre-training. We adopt a batch size of 64 and an Adam optimizer with a learning rate of  $2e-5$ . Across 10 epochs of training for 20 hours, the model with the lowest WER on development set at the end of each epoch was chosen as the best model for the speech encoder.

**Audio masking** Motivated by other multi-modal ASR study (Shi et al., 2022), a small portion of the speech data is obscured by masking to reduce the model’s reliance on speech input. Specifically, 10% of speech data are randomly chosen for masking, and each selected data is masked for 20% of its total duration. Note that this configuration of masking probability and duration was empirically determined to yield optimal results in our experiments, with the proportion of masked data and masking length varied between 10% to 30% and 10% to 50%, respectively. To implement the masking process, we segment each audio into discrete chunks of 1-second duration. These chunks serve as the minimum unit for the masking, e.g. in an audio input with a duration of 10 seconds, two randomly chosen chunks would be masked. Unless otherwise specified, all results of the CA-ASR

model include audio masking during training.

**CNRL Setup** We utilized the noisy context training set from the 10-folds described in Section 3.3. The average WER for the noisy context was 6.53% across the 10 folds. We filtered out transcriptions with a WER exceeding 20% to prevent interference with CNRL, resulting in the exclusion of 8.2% of the noisy context training set. We evaluated the effect of CNRL noisy context data by modifying the dialogue turns and introducing arbitrary word drops. For arbitrary word drop, we remove words for user utterances from golden context by 10% of change for each word and iterate it until we match the WER for each dialogue up to 6.5%, which is similar to WER with 10-folds. The training data setups are listed below:

- **S1:** Arbitrarily remove words from the golden context (user utterance only) to match an average WER of 6.5%.
- **S2:** Using the 10-fold training set, only the last user utterance contains noisy text.
- **S3:** Using the 10-fold training set, all user utterances may contain noisy text.
- **S4:** Combining S1 with S3. If a user utterance for each turn does not contain noisy text, arbitrary word drops are applied to increase the noise.

Unless otherwise specified, subsequent experimental results with CNRL use the **S4** setup. We use a batch size of 128 and the Adam optimizer with a learning rate of  $5e-4$ . Training is conducted for up to 5 epochs, selecting the epoch with the lowest cosine embedding loss on our development set.

All experiments are conducted using 4 NVIDIA A6000 GPUs.

## 5 Result & Analysis

### 5.1 Context Aware-ASR

Table 1 illustrates the WER across various models and noise levels. The CA-ASR model significantly improves performance on our test set, reducing relative WER by **33.4%** compared to the RNN-T (Soltau et al., 2022) baseline (**7.92% vs. 11.90%**) and by **14.2%** compared to the wav2vec2.0 with BART baseline, even without additional methods like CNRL or decoder pre-training. This highlights the advantage of using multi-modality with a context encoder for dialogue speech recognition.

Configurations			Audio Noise Level		
Model	Modality	Parameter size	No Noise	SNR:20dB	SNR:0dB
DSTC11 RNN-T (Soltau et al., 2022)	Speech	220M	11.90%	-	-
DSTC11 Whisper (Soltau et al., 2022)*	Speech	1550M	8.50%	-	-
Whisper-large-v2 (Radford et al., 2023)**	Speech	1550M	8.10%	8.45%	14.82%
Wav2Vec2.0+BART (baseline)	Speech	569M	9.23%	11.89%	18.45%
CA-ASR (Ours)	Speech+Text	774M	7.92%	8.23%	15.65%
+CNRL	Speech+Text	774M	7.66%	8.10%	15.03%
+Decoder Ptr.	Speech+Text	774M	7.39%	7.51%	13.33%
+Decoder Ptr. & CNRL	Speech+Text	774M	<b>7.04%</b>	<b>7.24%</b>	<b>12.65%</b>

Table 1: WER comparison of various baselines and CA-ASR settings under different noise conditions. Our proposed CA-ASR model is evaluated with and without Context Noise Representation Learning (CNRL) and Decoder Pretraining (Decoder Ptr.) enhancements. \*: reported. \*\*: re-evaluated with our post-processing.

Decoder pre-training further enhances the performance of the CA-ASR model, significantly reducing relative WER by **6.7%**, especially under severe noise conditions (SNR:0dB) where the voice is barely audible. This is expected since the decoder is initially tuned to the dialogue domain, enabling it to predict the user’s subsequent probable response from the context even with incomplete speech input.

The benefits are maximized when CNRL is applied, resulting in a relative WER reduction of **11.1%** in clean conditions and **19.1%** in noisy environments compared to the basic CA-ASR model. Since CNRL is designed to make the context encoder resilient to context errors, it significantly enhances the model’s robustness against strong noise.

Under the noisy audio test set (refer to Section 4.1), each model’s performance declines as the noise level increases (SNR:20dB to SNR:0dB). However, incorporating decoder pre-training and CNRL significantly mitigates this performance drop compared to the basic CA-ASR model (**12.65% vs. 15.65%**).

While the Whisper model shows robust performance under severe noise conditions (SNR:0dB), our CA-ASR model with CNRL and decoder pre-training demonstrates even greater robustness (**12.65% vs. 14.82%**).

## 5.2 Decoder Pre-training for Dialogue

Table 2 demonstrates the effectiveness of pre-training the decoder with varying the number of turns and pre-training dataset sizes. Note that the baseline model is the same as in Table 1, consisting only of a speech encoder (wav2vec2.0) and a BART decoder. As illustrated, pre-training the decoder on

Model	Input Dialogue	Decoder Pre-training	WER
baseline	-	BART(Lewis et al., 2019)	9.23%
baseline	-	+ MultiWoZ 2.1	8.95%
baseline	-	+ SGD	8.88%
CA-ASR	single-turn	BART	8.14%
CA-ASR	single-turn	+ MultiWoZ 2.1	7.98%
CA-ASR	single-turn	+ SGD	7.64%
CA-ASR	multi-turn	BART	7.92%
CA-ASR	multi-turn	+ MultiWoZ 2.1	7.45%
CA-ASR	multi-turn	+ SGD	<b>7.39%</b>

Table 2: WER across various accumulated datasets and a number of turn-takings. Note that CNRL and noise evaluation are not applied in this result to focus on the efficacy of decoder pre-training.

the dialogue domain benefits both the speech-only model (baseline) and the speech-text multimodal model (CA-ASR). Compared to the best result of baseline, the inclusion of the context encoder leads to significant improvements, resulting in a relative WER reduction of approximately **16.7%** at best in CA-ASR with multi-turn (**8.88% vs. 7.39%**). This finding suggests that the efficacy of pre-training the decoder is maximized when the model incorporates information from previous dialogues. Additionally, the WER of CA-ASR with multi-turn improves relatively by up to 6.7% as the dataset size increases (adding SGD), indicating the utility of incorporating external datasets as long as they involve user-agent conversations. Moreover, models considering multiple turns of dialogue exhibit a relatively 3.2% better WER compared to those considering a single turn, as shown in the comparison of best results (**7.64% vs 7.39%**). This highlights the importance of considering a longer context.

Model (Modality)	Audio Masking	No Noise	SNR:20db	SNR:0db
baseline (Speech)	No	8.94%	11.20%	18.02%
baseline (Speech)	Yes	8.88%	10.58%	17.61%
CA-ASR (Speech + Text)	No	7.45%	7.88%	14.28%
CA-ASR (Speech + Text)	Yes	<b>7.04%</b>	<b>7.24%</b>	<b>12.65%</b>

Table 3: WER comparison between modality and audio masking in clean and noisy samples. Each model’s decoder is pre-trained with Multi-WoZ 2.1 and SGD, and CNRL is additionally applied to CA-ASR.

### 5.3 Effect of Audio masking

Since audio masking can serve as data augmentation, we conducted additional experiments to compare the performance improvement between the baseline (speech-only) model and the CA-ASR (multimodal) model. As shown in Table 3, audio masking enhances ASR performance in both the baseline and CA-ASR models. While the baseline models exhibit marginal performance improvements of about 0.6% in clean sample evaluations, CA-ASR benefits from audio masking with a **5.5%** relative WER reduction. The improvement in CA-ASR becomes more pronounced in noisy environments as noise levels increase. Although the WER is highest at SNR:0dB, indicating the strongest noise, the relative WER reduction is **11.4%**, compared to 8.12% at SNR:20dB. These results suggest that while audio masking is beneficial in both clean and noisy environments, its effect is maximized when the model can utilize contextual information.

### 5.4 Context Noise Representation Learning

To investigate the impact of noise data on CNRL, we conducted experiments using different types of noise (S1-S4) as described in the CNRL setup in Section 4.3. In Table 4, compared to the model without CNRL, S1 (which arbitrarily removes words) degraded performance, indicating that using only artificial noise is not beneficial for CNRL. S2 and S3, which use real ASR noise from 10-fold data generation, showed better performance, with multi-turn noise (S3) outperforming single-turn noise (S2).

In our evaluation, we found that S4, which combines S1 with S3, performed the best, with WERs of 7.04%, 7.24%, and 12.65% for No-Noise, SNR:20dB, and SNR:0dB conditions, respectively. For comparison, we evaluated our model with ground truth context during inference without CNRL, serving as the upper bound of our experiment. As expected, using ground truth context showed robust results across noise levels, while

CNRL.	No Noise	SNR:20db	SNR:0db
No	7.39%	7.51%	13.33%
S1	7.53%	7.45%	13.45%
S2	7.30%	7.41%	12.94%
S3	7.22%	7.29%	12.83%
S4	7.04%	<b>7.24%</b>	12.65%
Ground Truth Context*	<b>7.01%</b>	7.25%	<b>12.28%</b>
full fine-tune w/ S4	7.24%	7.63%	13.50%

Table 4: CNRL result on different training data settings (S1, S2, S3 and S4) including evaluation result with ground truth context (\*) and full fine-tuning result.

CNRL with S4 produced similar results with a small margin. This demonstrates that CNRL enables the context encoder to handle noisy contexts effectively, generating representations close to the ground truth.

We also experimented with training the full CA-ASR model, not just the context encoder, using S4 with corresponding audio for ASR fine-tuning. Training the full model showed lower performance gains than CNRL (**7.04% vs. 7.24%**) and required much larger training costs. We believe this is because training all components with noisy data can disrupt optimization. CNRL allows us to maintain ASR performance against noisy contexts while keeping training efficient.

## 6 Conclusion

This study introduced Context Noise Representation Learning (CNRL) to improve context-aware ASR systems, especially in noisy environments. By integrating decoder pre-training with dialogue data, ASR fine-tuning, and CNRL, we significantly reduced transcription errors. Our training pipeline demonstrated significant improvements in dialogue speech recognition, even in noisy environments where speech input is defective. Experiments showed CNRL’s efficacy, reducing Word Error Rate (WER) by up to 11.1% in clean conditions and 19.1% in noisy settings. By making the model more robust against noisy context, our approach consistently outperformed baselines in various settings.

## Limitations

Due to the scarcity of spoken turn-based dialogue datasets, we could only validate our method on a single dataset DSTC11. However validating on the various test datasets would improve its credibility if applicable.

Our primary goal is to enhance ASR performance. However, these enhancements could be even more valuable for downstream Dialogue State Tracking (DST) tasks. Future work could explore optimizing ASR specifically for DST applications to further increase the impact and value of our contributions.

## Acknowledgements

This work was supported by the Technology Innovation Program(20015007, Development of Digital Therapeutics of Cognitive Behavioral Therapy for treating Panic Disorder) funded By the Ministry of Trade, Industry & Energy(MOTIE, Korea).

This research was supported by the MSIT(Ministry of Science and ICT), Korea, under the ITRC(Information Technology Research Center) support program(IITP-2024-2020-0-01789) supervised by the IITP(Institute for Information Communications Technology Planning Evaluation)

## References

- Siddhant Arora, Siddharth Dalmia, Pavel Denisov, Xuankai Chang, Yushi Ueda, Yifan Peng, Yuekai Zhang, Sujay Kumar, Karthik Ganesan, Brian Yan, et al. 2022. Espnet-slu: Advancing spoken language understanding through espnet. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7167–7171. IEEE.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
- Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518.
- 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.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. **MultiWOZ 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking baselines**. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 422–428, Marseille, France. European Language Resources Association.
- Changfeng Gao, Gaofeng Cheng, Runyan Yang, Han Zhu, Pengyuan Zhang, and Yonghong Yan. 2021. Pre-training transformer decoder for end-to-end asr model with unpaired text data. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6543–6547. IEEE.
- Takaaki Hori, Niko Moritz, Chiori Hori, and Jonathan Le Roux. 2020. Transformer-based long-context end-to-end speech recognition. In *Interspeech*, pages 5011–5015.
- Junfeng Hou, Jinkun Chen, Wanyu Li, Yufeng Tang, Jun Zhang, and Zejun Ma. 2022. Bring dialogue-context into rnn-t for streaming asr. In *INTERSPEECH*, pages 2048–2052.
- Wei-Ning Hsu, Yao-Hung Hubert Tsai, Benjamin Bolte, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: How much can a bad teacher benefit asr pre-training? In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6533–6537. IEEE.
- Chao-Wei Huang and Yun-Nung Chen. 2020. Learning asr-robust contextualized embeddings for spoken language understanding. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8009–8013. IEEE.
- Ridong Jiang, Wei Shi, Bin Wang, Chen Zhang, Yan Zhang, Chunlei Pan, Jung Jae Kim, and Haizhou Li. 2023. **Speech-aware multi-domain dialogue state generation with ASR error correction modules**. In *Proceedings of The Eleventh Dialog System Technology Challenge*, pages 105–112, Prague, Czech Republic. Association for Computational Linguistics.
- Seokhwan Kim, Yang Liu, Di Jin, Alexandros Papanagelis, Karthik Gopalakrishnan, Behnam Hedayatnia, and Dilek Hakkani-Tür. 2021. “how robust ru?”: Evaluating task-oriented dialogue systems on spoken conversations. In *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1147–1154. IEEE.
- Chenyang Le, Yao Qian, Long Zhou, Shujie Liu, Yanmin Qian, Michael Zeng, and Xuedong Huang. 2024.

- Comsl: A composite speech-language model for end-to-end speech-to-text translation. *Advances in Neural Information Processing Systems*, 36.
- 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.
- Yuchen Liu, Jiajun Zhang, Hao Xiong, Long Zhou, Zhongjun He, Hua Wu, Haifeng Wang, and Chengqing Zong. 2020. Synchronous speech recognition and speech-to-text translation with interactive decoding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8417–8424.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Zhiyuan Ma, Jintao Du, and Shuheng Zhou. 2023. Noise-robust training with dynamic loss and contrastive learning for distantly-supervised named entity recognition. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10119–10128.
- Ryo Masumura, Naoki Makishima, Mana Ichori, Akihiko Takashima, Tomohiro Tanaka, and Shota Orihashi. 2020. Phoneme-to-grapheme conversion based large-scale pre-training for end-to-end automatic speech recognition. In *INTERSPEECH*, pages 2822–2826.
- Pablo Ortiz and Simen Burud. 2021. Bert attends the conversation: Improving low-resource conversational asr. *arXiv preprint arXiv:2110.02267*.
- Cheonyoung Park, Eunji Ha, Yewon Jeong, Chi-young Kim, Haeun Yu, and Joo-won Sung. 2023. Copyt5: Copy mechanism and post-trained t5 for speech-aware dialogue state tracking system. In *Proceedings of The Eleventh Dialog System Technology Challenge*, pages 89–94.
- Karol J. Piczak. 2015. [Esc: Dataset for environmental sound classification](#). In *Proceedings of the 23rd ACM International Conference on Multimedia*, MM '15, page 1015–1018, New York, NY, USA. Association for Computing Machinery.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8689–8696.
- Dmitriy Serdyuk, Yongqiang Wang, Christian Fuegen, Anuj Kumar, Baiyang Liu, and Yoshua Bengio. 2018. Towards end-to-end spoken language understanding. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5754–5758. IEEE.
- Ashish Shenoy, Sravan Bodapati, Monica Sunkara, Srikanth Ronanki, and Katrin Kirchhoff. 2021. [Adapting Long Context NLM for ASR Rescoring in Conversational Agents](#). In *Proc. Interspeech 2021*, pages 3246–3250.
- Bowen Shi, Wei-Ning Hsu, Kushal Lakhota, and Abdelrahman Mohamed. 2022. Learning audio-visual speech representation by masked multimodal cluster prediction. *arXiv preprint arXiv:2201.02184*.
- Hagen Soltau, Izhak Shafran, Mingqiu Wang, Abhinav Rastogi, Jeffrey Zhao, Ye Jia, Wei Han, Yuan Cao, and Aramys Miranda. 2022. Speech aware dialog system technology challenge (dstc11). *arXiv preprint arXiv:2212.08704*.
- Xin Sun, Qiang Liu, Shu Wu, Zilei Wang, and Liang Wang. 2023. Noise-robust semi-supervised learning for distantly supervised relation extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13145–13157.
- Yun Tang, Juan Pino, Xian Li, Changhan Wang, and Dmitriy Genzel. 2021. Improving speech translation by understanding and learning from the auxiliary text translation task. 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 4252–4261.
- Xin Tian, Xinxian Huang, Dongfeng He, Yingzhan Lin, Siqi Bao, Huang He, Liankai Huang, Qiang Ju, Xiyuan Zhang, Jian Xie, Shuqi Sun, Fan Wang, Hua Wu, and Haifeng Wang. 2021. [Tod-da: Towards boosting the robustness of task-oriented dialogue modeling on spoken conversations](#). *Preprint*, arXiv:2112.12441.
- Emiru Tsunoo, Hayato Futami, Yosuke Kashiwagi, Sidhant Arora, and Shinji Watanabe. 2023. Decoder-only architecture for speech recognition with ctc prompts and text data augmentation. *arXiv preprint arXiv:2309.08876*.
- Longshaokan Wang, Maryam Fazel-Zarandi, Aditya Tiwari, Spyros Matsoukas, and Lazaros Polymenakos. 2020. Data augmentation for training dialog models robust to speech recognition errors. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 63–70.
- Mingqiu Wang, Izhak Shafran, Hagen Soltau, Wei Han, Yuan Cao, Dian Yu, and Laurent El Shafey. 2024. Retrieval augmented end-to-end spoken dialog models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 12056–12060. IEEE.

Mingqiu Wang, Izhak Shafran, Hagen Soltau, Wei Han, Yuan Cao, Dian Yu, and Laurent El Shafey. 2023. Speech-to-text adapter and speech-to-entity retriever augmented llms for speech understanding. *arXiv preprint arXiv:2306.07944*.

Jiahao Xu, Wei Shao, Lihui Chen, and Lemao Liu. 2023. Simcse++: Improving contrastive learning for sentence embeddings from two perspectives. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

Xiaolin Zheng, Mengling Hu, Weiming Liu, Chaochao Chen, and Xinting Liao. 2023. Robust representation learning with reliable pseudo-labels generation via self-adaptive optimal transport for short text clustering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10493–10507.

# Local Topology Measures of Contextual Language Model Latent Spaces With Applications to Dialogue Term Extraction

Benjamin Matthias Ruppik, Michael Heck, Carel van Niekerk, Renato Vukovic,  
Hsien-Chin Lin, Shutong Feng, Marcus Zibrowius, Milica Gašić

Heinrich Heine University Düsseldorf, Germany  
{ruppik, heckmi, niekerk, renato.vukovic,  
linh, shutong.feng, marcus.zibrowius, gasic}@hhu.de

## Abstract

A common approach for sequence tagging tasks based on contextual word representations is to train a machine learning classifier directly on these embedding vectors. This approach has two shortcomings. First, such methods consider single input sequences in isolation and are unable to put an individual embedding vector in relation to vectors outside the current local context of use. Second, the high performance of these models relies on fine-tuning the embedding model in conjunction with the classifier, which may not always be feasible due to the size or inaccessibility of the underlying feature-generation model.

It is thus desirable, given a collection of embedding vectors of a corpus, i.e. a datastore, to find features of each vector that describe its relation to other, similar vectors in the datastore. With this in mind, we introduce complexity measures of the local topology of the latent space of a contextual language model with respect to a given datastore.

The effectiveness of our features is demonstrated through their application to dialogue term extraction. Our work continues a line of research that explores the manifold hypothesis for word embeddings, demonstrating that local structure in the space carved out by word embeddings can be exploited to infer semantic properties.

## 1 Introduction

The prevailing approach to sequence tagging tasks such as named entity recognition or dialogue term extraction involves a two-step process: start with a general contextual vector representation for text sequences, for instance the embedding vectors created by a pretrained language model, then train a separate tagging model on top of the vector representations (Lample et al., 2016; Ramshaw and Marcus, 1995). Optionally, assuming differentiability of the model and target function, one can fine-tune

the representation model such that its embeddings are more suitable for the tagging task (Panchendrarajan and Amaresan, 2018). While highly effective, the representations may be expensive to compute, and fine-tuning a language model is not always feasible, for instance if the underlying model is hidden behind an application programming interface (API). Thus, it is desirable to develop tagging methods which achieve the best performance on the given representations. In fact, the performance of prompting-based approaches with large language models (LLMs) on named entity recognition tasks has lagged behind that of supervised sequence tagging approaches (Wang et al., 2023). Additionally, this leads to problems such as hallucinations and potential dataset contamination, which prevent a fair evaluation.

A more fundamental limitation of the prevailing paradigm is that the relation of a single input sequence to other sequences in the dataset cannot be taken into account. Both the representation module and the tagging module commonly have a limited maximum context length. They cannot process the entire dataset at once, but rather need to be provided with single sentences or paragraphs at a time. The limited context can lead to suboptimal performance (Amalvy et al., 2023). For example, consider named entity recognition for an isolated sentence such as *Prince was prominently featured at the event*. The word *Prince* is ambiguous. In a corpus containing news articles, *Prince* or *Prince Harry* likely appear in many articles related to the British royal family. In a different corpus, the term *Curry Prince* might appear frequently in the context of restaurant reviews. So only with regard to the entire corpus under consideration, an informed choice on how to tag *Prince* in the example sentence could be made.

In this work, we show that the relation between the representation of a single token and its containing corpus can be captured by studying the *latent*

*space* – the collection of the language model’s hidden states – surrounding the corresponding embedding vector. The geometry of these hidden states is known to capture both syntactic and semantic properties of the underlying text. For instance, [Cohen et al. \(2019\)](#) find that distances between the contextual vectors of bidirectional encoder representations from transformers (BERT) ([Devlin et al., 2019](#)) correspond to parse tree embeddings based on the grammatical structure of the input phrases. Here, we study neighborhoods of embedding vectors from a *topological* viewpoint, and introduce descriptors of the shapes of these neighborhoods that are stable under symmetries such as permutations, translations, and rotations. In particular, we define descriptors based on *persistent homology*, a well-established tool of topological data analysis ([Carlsson and Vejdemo-Johansson, 2021](#)).

## 1.1 Contribution

Consider the latent space of a language model in the neighborhood of a given contextual embedding vector. For instance, the neighborhood of an embedding of the word *cheap* in the context *I am looking for options for cheap dinner* contains other occurrences of the word *cheap* in different contexts, but also different words expressing a similar meaning (*inexpensive, good-value*) and words connected to the center word, such as *restaurant*. In this work, we show that:

- (a) this neighborhood contains information that is not present in the language model next-token prediction distribution, and that cannot be ‘distilled’ into the language model via naive fine-tuning,
- (b) this additional information can be used to improve the performance of sequence tagging tasks, and
- (c) this information can be efficiently summarized using low-dimensional topological feature descriptors.

Our topological descriptors are codensity at multiple scales ([Carlsson et al., 2008](#); [Carlsson, 2014](#)), topological singularity measures based on Wasserstein norms ([Cohen-Steiner et al., 2010](#)), and vectorized persistence modules. Towards (a), we show that several of our one-dimensional numerical measures show minimal correlation with language model perplexity, indicating that they contain independent information. Towards (b) and

(c), we empirically demonstrate improvements on the natural language processing task of variants of term extraction. In each case, we build the latent space through a masked language model from a dialogue corpus. As a baseline, we employ a tagging model trained directly on the original language model vectors, and compare with models that take as input a combination of the language model vectors with our topological descriptors of the neighborhood within the latent space of a contextual language model. Furthermore, we compare with models trained on features from [Vukovic et al. \(2022\)](#), which are based on neighborhoods in a *static* word embedding space. We show that utilizing the *contextually* augmented vectors yields statistically significant improvements.

Observation (a) is not completely new. For example, it is present in the idea of  $k$ -nearest neighbor language models ([Khandelwal et al., 2020](#); [Xu et al., 2023](#)), where the current hidden state is augmented by the nearest neighbors from a datastore. Our low-dimensional descriptors, on the other hand, have not been deployed before, and our experiments for (b) provide the first application of contextual topological features to token level sequence tagging tasks. Note with reference to point (c) that summarizing a collection of vectors in a permutation-invariant way is a challenging problem in representation learning ([Zaheer et al., 2017](#)), which we tackle in this work via tools from persistent homology.

Our work is complementary to other recent applications of topological methods to the study of contextual embedding spaces. [Tulchinskii et al. \(2023\)](#) demonstrate that the topology of a point cloud derived from a text paragraph can be utilized in a sequence classification task, namely to differentiate human-written from artificially generated paragraphs. Their approach takes into account solely the given paragraph’s embedding vectors, and does not explore how these reside within the larger latent space. Another approach involves constructing filtered graphs from the attention scores in a transformer model, followed by sequence-level classification based on persistent homology ([Kushnareva et al., 2021](#); [Perez and Reinauer, 2022](#)). However, this approach only applies to supervised sequence classification tasks, and does not yield local features required for tagging. In a more qualitative direction, [Valeriani et al. \(2023\)](#) investigate the intrinsic dimension of the latent space through the different layers of a transformer, and [Ethayarajh](#)

(2019) and Cai et al. (2021) identify isolated clusters and low dimensional manifolds in the latent spaces of various language models. However, they do not apply their quantitative local analysis to a practical task.

## 2 Background and Methods

### 2.1 Latent Spaces of Contextual Language Models

We consider the encoder part of a contextual language model, which can be thought of as a map

$$e: (\mathbb{R}^d)^{\times N} \rightarrow (\mathbb{R}^h)^{\times N}.$$

Here,  $d$  is the dimension of the input layer,  $h$  is the hidden dimension ( $h \ll d$ ), and  $N$  is the maximum sequence length, after which sequences will be truncated. This maximum length is usually fixed and finite. The input of the encoder is a sequence of vectors  $\mathbf{X} \in (\mathbb{R}^d)^{\times N}$  representing a tokenized context. *Tokenization* describes the process in which an input string is decomposed into a sequence of vectors. In our setting, tokenization can be thought of as a lookup layer converting short text segments to vectors (together with position information). Typically, longer words are split into several token vectors in this process.

The output of the encoder is a sequence of so-called hidden states. Commonly, these hidden states are inputs to the “prediction head” of the language model, which produces a probability distribution over the token space for the corresponding token location.

We think of a language corpus  $C$  as a collection of tokenized portions of text. From the point of view of a language model, each instance  $i$  of a particular token appears in a specific context  $\mathbf{X}(i)$ . These contexts are filled with padding tokens so that they always have length  $N$ , permitting construction of the embedding sets  $e(\mathbf{X}(i))$ .

**Definition 2.1** *Given an encoder  $e$  derived from a pretrained language model, the ambient corpus datastore with respect to a corpus  $C$  is the multi-set<sup>1</sup>/point cloud of all the embeddings  $e(\mathbf{X}(i))$  of all instances  $i$  of all tokens in  $C$ .*

Note that we cannot explore the entire latent space of the language model, but only the subspace “carved out” by the datastore under consideration,

<sup>1</sup>We write *multi-set* to allow for repetitions/multiplicities. This is relevant in our setting, because strings might appear multiple times in the corpus.

as in Definition 2.1. In other words, by selecting the task-dependent ambient corpus for sampling the language model hidden states, we are making a choice of how we explore the hidden state space. This choice of ambient corpus may have a big impact on the derived features.

### 2.2 Local Topological Measures

All our topological measures are based on neighborhoods of a given contextual embedding vector  $v$  with respect to a collection of contextual embedding vectors coming from an ambient corpus datastore. Given an integer  $n \geq 1$ , we define the *neighborhood*  $\mathcal{N}_n(v)$  as the multi-set consisting of  $v$  and its  $(n - 1)$  nearest neighbors. To avoid adding another copy of the query center vector  $v$  when building the neighborhood from the datastore, we first check for similarity to existing vectors in the datastore with a Euclidean distance threshold of  $10^{-4}$ , and take a possible match as center vector if applicable. For a schematic illustration of the neighborhood extraction process and the feature computation, see Figure 1. We consider the following local features:

**Persistence Images** For a positive integer  $d$ , *persistent homology of degree  $d$*  associates with a point cloud a multi-set that encodes “ $d$ -dimensional topological features” of the cloud. We refer to Edelsbrunner and Harer (2010) or Otter et al. (2017) for introductions. Various vectorizations of this multi-set have been developed for subsequent use in machine learning. Persistence images are introduced in Adams et al. (2017) as a refined, higher-dimensional vectorization of persistent homology. We define  $\text{PI}^d(v) \in \mathbb{R}^{100}$  as a persistent image vector of the degree  $d$  persistent homology of  $\mathcal{N}_n(v)$ , scaled by a factor of  $\frac{1}{n \cdot 100}$ . The parameter  $n$  is not included in this notation, as it will be fixed to 128 throughout all experiments. For detailed definitions, see Appendix A.1. The factor  $\frac{1}{n \cdot 100}$  appearing in our definition of the persistence image is not important at this point. It is included to avoid instabilities in the training of the BIO-tagger discussed in Section 3, which may otherwise arise from the vastly different scales of the coordinates of the language model embedding vectors and these additional coordinates.

**Wasserstein Measure** A simple one-dimensional vectorization is the Wasserstein norm. We

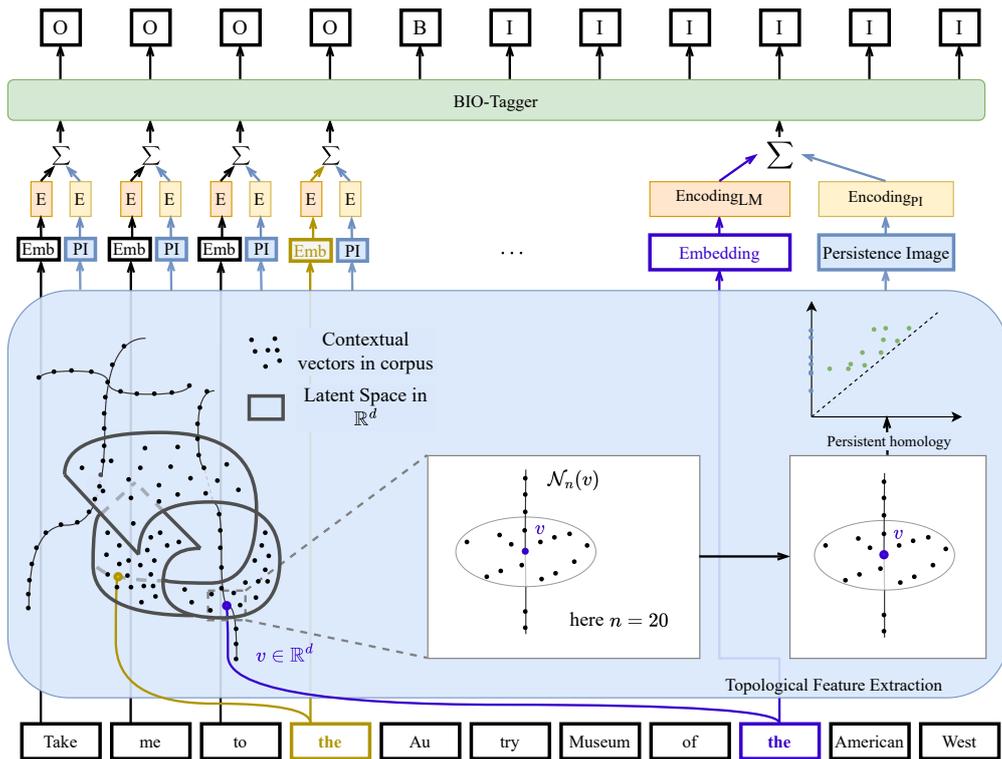


Figure 1: Schematic illustration of the local topological feature extraction and of our topological deep learning pipeline: The blue box illustrates the extraction of neighborhoods  $\mathcal{N}_n(v)$  in the contextualized embedding space, followed by the computation of each neighborhood’s topological features, resulting in a contextualized persistence image vector. Note the color coding of the different occurrences of the token ‘the’; contextuality leads to different language model embedding vectors and persistence images depending on whether it is part of the term ‘Autory Museum of the American West’ or used as a non-content word. For each token, the language model embedding (Emb) and persistence image vectors (PI) are encoded (E), combined ( $\Sigma$ ), and then serve as input to our BIO-tagging transformer (green), which is trained on the token-level term labels (B-TERM (begin), I-TERM (inside), O (outside)).

define  $W_n^d(v) \in \mathbb{R}$  as the Wasserstein norm of the degree  $d$  persistent homology of  $\mathcal{N}_n(v)$ .

**Codensity** We define the  $n$ -th codensity  $\text{coden}_n(v) \in \mathbb{R}$  as the radius of  $\mathcal{N}_{n+1}(v)$ . Higher codensity corresponds to regions where the vectors are farther spread apart.

There are several reasons why we fix the cardinality  $n$  of the neighborhoods rather than, say, their radius. Firstly, fixing the cardinality takes into account sample density of the ambient corpus from the latent space of the language model. If we took a fixed radius, sparse regions of the ambient corpus space would be underrepresented. Secondly, some of the topological features we consider are more readily comparable when computed on fixed cardinalities. Indeed, a reasonable comparison of Wasserstein norms of neighborhoods of different cardinalities seems difficult, and our multiscale definition of (co)density could also not

easily be emulated for neighborhoods of fixed radius. Finally, computational complexity limits the feasibility of approaches that allow for unlimited cardinalities of neighborhoods. For instance, in [Von Rohrscheidt and Rieck \(2023\)](#), where neighborhoods of fixed radii are employed, an additional sampling step is necessary. More on this is discussed in [Appendix A.1](#).

### 3 Application of Local Topology Measures to Token Level Tagging Task

We perform a correlation analysis of local features and conduct a case study to explore the efficacy of local topology measures. Specifically, we apply our proposed topological feature augmentations to the task of dialogue term extraction.

#### 3.1 Set-Up

**Data** For the term extraction case study, we resort to the MultiWOZ2.1 ([Budzianowski et al.](#),

2018; Eric et al., 2020) and schema-guided dialogue (SGD) (Rastogi et al., 2020) task-oriented dialogue datasets. Here, the ambient reference corpus  $C$  is built solely from the language model hidden states for the training corpus of MultiWOZ2.1, so that the local measurements are comparable across both datasets.

**BIO-Tagging** For the sequence tagging tasks, we employ a beginning (B), inside (I), outside (O) labeling schema, as in Qiu et al. (2022). To keep the comparison between our different models fair and to obtain statements about the quality of the underlying features, we choose the architectures so that the trainable BIO-tagging components have a similar number of adjustable parameters. In this way, we can safely attribute any increase in performance to our topological augmentation of the input features rather than a stronger tagging component.

In all cases, the BIO-tagging transformer follows the RoBERTa architecture (Liu et al., 2019) and uses 8 attention heads, 2 hidden layers, and 512 maximum position embeddings. The language model vectors and augmenting feature vectors are fed into the BIO-tagging component through separate two-layer fully connected encoding networks with subsequent individual layer normalization, whose purpose is down- or up-scaling the feature dimension (768 for the language model vectors, 100 for persistence images) to the hidden size (512) of the tagging transformer. For a schematic of our BIO-tagging setup, see Figure 1.

**Features** For creating the language model embeddings, we use the second to last hidden states (at layer 11) of the pretrained RoBERTa base model (Liu et al., 2019), which returns 768-dimensional vectors, with  $L_2$ -normalization. Note that on unit vectors, the cosine distance is proportional to the square of the Euclidean distances, thus for the relative order in which the nearest neighbors occur, it does not matter whether we search with respect to the cosine or the Euclidean distance.

We decide on the *second-to-last* hidden state of the language model, as opposed to another intermediate layer, for two reasons: Cai et al. (2021) show that the local intrinsic dimension tends to increase with the depths in the transformer, thus the resulting neighborhoods should be more expressive. Moreover, Peters et al. (2018) and Tenney et al. (2019) find that deeper layers in language models tend to capture more of the semantic properties, while earlier layers tend to capture the syntax. For

feature based learning, Devlin et al. (2019) report that among single-layer features, the second-to-last layer leads to the highest performance. Note that our setup is not specific to the RoBERTa model or tokenizer. Our contextualized topological features can be computed for any (masked or causal) embedding model, extraction layer, datastore produced by the model, and query dataset.

As a baseline in our term extraction experiments, we use the language model hidden state vectors described in Section 3 as input for the BIO-tagging model. We test these against augmentation of the hidden states with our local persistence image descriptors introduced in Section 2.2.

### 3.2 Correlation Analysis of Local Features

We begin by collecting statistical observations on the correlation between the different local topological feature types, as well as their correlation with *pseudo-perplexity*, on the example of the MultiWOZ2.1 and SGD datasets. The *perplexity* of a causal language model is a model intrinsic measurement of the surprise of seeing the next token, defined as the exponentiation of the cross-entropy between the model prediction and the corpus data. While causal perplexity is not available here, in our the masked language model setting, we apply a pseudo-likelihood score by masking tokens one by one, and computing the prediction loss of the masked out token following Salazar et al. (2020).

In addition to the non fine-tuned version of the language model, here we also include the perplexity of the fine-tuned version for comparison, which uses the MultiWOZ2.1 training portion (fine-tuning for 5 epochs, 0.15 masking proportion, selecting the best model on MultiWOZ2.1 validation loss). All local features are based a non fine-tuned version of the language model.

Note that we cannot directly compute correlations between the vector-valued persistence images and the perplexity measures. For this reason, we are relying on selected codensities  $\text{coden}_n(v)$  with values  $n \in [1; 127; 511]$  and Wasserstein norms as numerical proxy estimates of the neighborhoods' topological complexity. Our term extraction models will take only the persistence images as input, as they provide the most powerful and comprehensive representation of the local topology. The neural network feature extraction model can learn directly from the persistence images to estimate complexity measures approximating the Wasserstein norms, which avoids manual feature engineering.

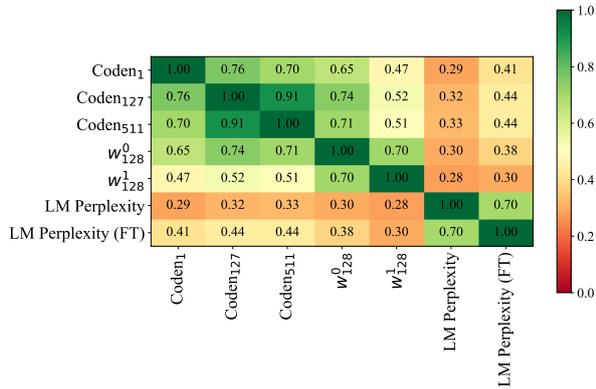


Figure 2: Kendall’s rank correlation coefficients between various local estimates and language model (LM) perplexity for the SGD test dataset. FT stands for LM fine-tuned on the MultiWOZ2.1 training split. All correlations have  $p < 10^{-6}$ .

The results of Kendall’s rank correlation are given in Figure 2. The Wasserstein measures are strongly positively correlated with the codensities. More importantly, the codensities and Wasserstein norms are only weakly correlated with the perplexity, indicating that these topological measures (and thus persistence images) capture information that is not present in the language model masked-token prediction distribution, and that cannot be ‘distilled’ into the language model via naive fine-tuning.

### 3.3 Case Study: Term Extraction in Task-Oriented Dialogue Data

#### 3.3.1 Task Definition

We approach dialogue term extraction as a transfer learning problem. Here, MultiWOZ2.1 serves as our source dataset used for training a term extractor. The trained model is then applied to the SGD dataset, necessitating sufficient transfer learning capabilities to properly handle the distributional shift in the data. We label all phrases in all utterances that match an entry in the respective dataset’s ontology, i.e., that match a value in a non-categorical slot of the current turn’s dialogue state or a value in the current turn’s dialogue act. The dataset ontology entities are normalized and matched to the occurrences in the respective utterances by applying the TripPy-R label map (Heck et al., 2022) and the SGD canonical value mapping. The ontology comprises names of entities, their domains, properties (slots), and values of these slots. We refer to these labelled phrases as *dialogue terms*.

These tagged spans for the dialogue datasets are encoded for the BIO-tagger, resulting in the three

labels: O (outside), B-TERM, I-TERM (begin and inside a term). Since our BIO-tagging model operates on the token-level of the underlying language model, we re-align the tags with the tokenization using the IOB2 schema: for a word with B-tag, the first subtoken is tagged with B, its remaining subtokens with I. For a word with I-tag, all its subtokens are tagged with I; and analogously for the O-tag.

While we employ token-level cross-entropy loss as the differentiable target function in the model training, the objective of term extraction within the context of this work is to retrieve each unique target dialogue term at least once. That is, we do not require the tagger to find all occurrences: We normalize the predicted phrases and ground truth by lower-casing, and deduplicate the resulting collections. A term is considered a true positive if it is identical to exactly one ground truth term. If a term cannot be assigned to any ground truth term, e.g., comprises several ground truth terms or is an incomplete substring of a ground truth term, it is considered a false positive. The left-over ground truth terms without a matched prediction are the false negatives. We call the resulting prediction, recall and F1-scores the *phrasal results*.

We train the BIO-tagger for 10 epochs, using the AdamW optimizer (Loshchilov and Hutter, 2019) with learning rate  $5 \cdot 10^{-5}$ , linear warm-up for 10% of training steps and batch size of 48. The model predictions on the held-out MultiWOZ2.1 validation set are evaluated every 100 batches for the first 3 000 global steps, and the model checkpoint with the best phrasal results on the validation set is selected as the final model.

Our goal is to show that injecting our local topological features into the model yields statistically significant improvements over the original language model embeddings. We run statistical tests over multiple different random seeds for initialization and check for significant changes in evaluation scores on the transfer set, which we take as the full collection of 463 284 utterances from the SGD dataset comprising 20 domains.

**Training Data** In the full data setting, we train on all 113 556 utterances of the MultiWOZ2.1 training split, the results are included in Table 1.

To demonstrate that the contextual topological features are useful in settings with reduced data and might help in mitigating overfitting, we create a variation of the transfer task by restricting training to subsets of the MultiWOZ2.1 dataset. This is

a more realistic transfer setting, since a good model checkpoint needs to demonstrate that it can generalize to the unseen left-out domain which it encounters in the MultiWOZ2.1 validation split for the first time. Given one of the five major domains  $\mathcal{D} \in [\text{attraction}; \text{hotel}; \text{restaurant}; \text{taxi}; \text{train}]$  in the training split, we exclude those utterances contained in dialogues from  $\mathcal{D}$  in the tagger training, which leaves [71 768; 59 222; 58 156; 86 568; 66 736] utterances respectively. Model selection is performed based on phrasal F1-score on the 14 748 validation utterances, which span over all five domains. We report results of this cross-validation setup by macro averaging the phrasal scores on the SGD dataset over 10 seeds for each of the five left out data folds in Table 2.

**Static Topological Features Baseline** We evaluate term extraction performance on the level of phrase predictions. The phrase-level evaluation allows a comparison with Vukovic et al. (2022), who present a method that employs *static* topological descriptors in sequence tagging tasks. The main differences to our approach with contextual topological features are as follows:

- Our local persistent homology descriptors are defined on token level with respect to the tokenization of the language model. This is essential for combining our new features with the language model embeddings to create fusion models which can provide predictions on token-level. The *static* topological features of Vukovic et al. (2022) operate on word level, and they only gained contextuality in the BIO-tagging component of the model. Note that in this and our work, the context of an input sequence is a single dialogue utterance.
- Vukovic et al. (2022)’s features were based on neighborhoods in an ambient static word space composed of the 100 000 most common words in the English language. Thus, their method depends both on having a word-level separation of the input data, and on a given dictionary. Our contextualized features on the other hand can be defined without any additional external data.

For a comparison between static and contextual features, we align the static topological features with the roberta-base tokenization in our BIO-tagging setup, and train BIO-taggers with the

Input features	Precision $\uparrow$	Recall $\uparrow$	F1 $\uparrow$
LM roberta-base	48.89	56.61	52.39
LM roberta-base $\oplus$ static PI <sup>0</sup>	49.33	<b>58.82</b>	53.62
LM roberta-base $\oplus$ contextual PI <sup>0</sup>	<b>50.26*</b>	58.44	<b>53.97*</b>

Table 1: Phrasal-level performance comparison for term extractors trained on the MultiWOZ2.1 training split and evaluated on the full SGD dataset. Results are averages over 15 seeds. \* indicates statistically significant differences (one-sided independent *t*-test) w.r.t. the baseline LM roberta-base with  $p < 0.05$ .

Input features	Precision $\uparrow$	Recall $\uparrow$	F1 $\uparrow$
LM roberta-base	47.97	56.07	51.67
LM roberta-base $\oplus$ static PI <sup>0</sup>	47.92	<b>56.85</b>	51.94
LM roberta-base $\oplus$ contextual PI <sup>0</sup>	<b>48.92*</b> <sup>LM</sup> <sub>stat</sub>	56.23	<b>52.24</b>

Table 2: Cross-validated phrasal-level performance on SGD for term extractors trained on MultiWOZ2.1 training split without selected domain in [attraction; hotel; restaurant; taxi; train] averaged over 10 seeds for each of the five folds. \* indicates statistically significant difference with  $p < 0.05$ , w.r.t. the baseline LM roberta-base (LM) and augmentation with static persistence images (stat).

same architecture and data as in the contextual topological feature setting. To that end, the first constituent subtoken of each word is augmented with the word’s corresponding 100-dimensional  $H_0$ -persistence image feature vector of Vukovic et al. (2022).

### 3.3.2 Results

**Quantitative Analysis** Table 1 lists the term extraction performance for the pure language model baseline, our proposed method of augmenting with contextual topological features, and the alternative approach by augmenting with the static topological features from Vukovic et al. (2022). The main objective of said work was the maximization of recall, and to that end they proposed separate language model and topological feature taggers, with a subsequent union of predictions. In contrast, we show that our unified model augmented with contextual topological features can increase precision, recall, and F1 over the language model baseline.

Table 2 presents averaged results for the models trained on a reduced dataset constructed by omit-

Ground Truth Terms	LM roberta-base (our baseline)	LM roberta-base $\oplus$ static PI <sup>0</sup>	LM roberta-base $\oplus$ contextual PI <sup>0</sup>
“cafe jolie”	“cafe jolie”	“cafe jolie”	“cafe jolie”
“angelina jolie”	–	“angelina jolie”	“angelina jolie”
“yellow chilli” “the yellow chilli by chef sanjeev kapoor”	“yellow chilli” “the yellow chilli” “the yellow chilli by” “sanjeev kapoor”	“yellow chilli” “the yellow chilli” “the yellow chilli by chef” “sanjeev kapoor”	“yellow chilli” “the yellow chilli” “the yellow chilli by chef sanjeev kapoor”
“water seed” “water seed concert”	“water seed” “the water seed” “water seed event”	“water seed”	“water seed” “water seed concert”
“be alright”	–	“alright” “be alright”	“alright economy”

Table 3: Representative examples of predictions where the baseline model fails to retrieve the correct term, while a local topology feature augmented model succeeds. We indicate **true positives** and **false positives** by color.

ting a given domain in the MultiWOZ2.1 training set. Here, on average, the augmentation with the contextual persistence images is again better than the language model vector baseline.

**Qualitative Analysis** To obtain explicit examples, we select a model checkpoint for each feature type after 1 100 global steps, and inspect the differences between predicted normalized phrase sets. In Table 3 we see examples where the topologically augmented model succeeds in finding complete multi-word terms, whereas the baseline model cuts off before the end of a term or misses proper names that should follow a preposition. Such information is highly dependent on the context of a term, and the contextual topological model is able to find long terms more consistently. All models identify the restaurant name “Cafe Jolie” correctly, but only the topological models recognize the actress “Angelina Jolie”. Similarly, the song title “Be alright” containing the frequent word “alright” is not recognized by the language model alone, but can be detected by a topological model.

### 3.4 Relation to the Manifold Hypothesis

At first glance, our results may appear to be at odds with the *manifold hypothesis*, a common assumption underlying many representation learning paradigms. While this hypothesis has been questioned for *static* word embeddings (Jakubowski et al., 2020), it remains uncontested for contextual embeddings. It posits that the latent space of a trained machine learning model is clustered around a disjoint union of lower-dimensional manifolds (Bengio et al., 2013; Brown et al., 2023). This implies that, from a purely topological perspective, the local structure of the latent space is constant,

at least along connected components – every point should have a neighborhood topologically identical to an open ball in some Euclidean space. How then is it possible that we can extract meaningful information from variations of the local topology?

There are at least two answers to this. First, all our measures depend on the way data is *sampled*. There is no reason to assume that the embeddings drawn from a given corpus provide a uniform sample of the latent space. On the contrary, the distribution of these samples will depend heavily on the corpus. And within a given corpus, we might expect the neighborhoods of latent vectors of content words to be “more spread-out” and of higher dimension than those of non-content words, since there are more plausible possibilities for filling in content words in a text than for non-content words. Second, our measures are based on persistent homology, which is known to detect not only topological properties but also differentiable structure such as curvature (Bubenik et al., 2020). Thus, even on a uniformly sampled manifold, these measures are expected to vary.

### 3.5 Computational Complexity

In this section, we address the computational overhead coming from our proposed method of augmenting a sequence tagger with contextual topological information of a given corpus. The one-off computational costs for the datastore, in our study derived from the MultiWOZ2.1 training dataset, and the query datasets (MultiWOZ2.1 training & validation dataset, and SGD dataset) involve a single embedding model forward pass for each input sequence.

For each query dataset relative to the datastore, assuming a constant and small neighborhood

size  $n$ , the asymptotic complexity of the neighborhood search depends on the tokenized cardinality of the query dataset  $|Q|$ , the tokenized cardinality of the datastore  $|C|$ , and the embedding dimension  $d$ . The runtime complexity using the exact search implementation from (Johnson et al., 2021) is  $\mathcal{O}(|Q||C|d)$ , and the storage complexity for neighborhood indices is  $\mathcal{O}(|Q|n)$ .

The persistent homology computation in dimension 0 for each query vector depends on the neighborhood size  $n$  as well. For degree 0, the number of simplices in the Vietoris-Rips complex can be upper-bounded by  $n^2$ . Thus, the persistence diagram for each neighborhood can be computed in  $\mathcal{O}(n^{2\omega})$ , where  $\omega < 2.4$  is the matrix multiplication exponent (Milosavljević et al., 2011). There are at most  $n$  generators in the 0-dimensional persistence diagrams, so the computation of the Wasserstein norms can be achieved in  $\mathcal{O}(n^3)$  (Lacombe et al., 2018). Empirically, the computation of the persistence images is observed to be very quick compared to the computation of the persistence diagrams.

Once computed and cached, these topological features can be reused for different training objectives on the given query dataset. The only overhead in transitioning from the baseline tagger (approximately 35.65 million parameters) to the tagger with input LM roberta-base $\oplus$ contextual PI<sup>0</sup> involves a few additional parameters (roughly 60 000) for the encoding module of the 100-dimensional contextual persistence image. Consequently, once the topological features have been cached, the training and inference of the topologically augmented BIO-tagger are only negligibly slower than the baseline BIO-tagger. Appendix A.1 discusses the software packages used in the implementation and how we handle caching of the precomputed neighborhoods and resulting contextual topological features.

## 4 Conclusion

In this work, we introduce a topological deep learning approach to enrich feature learning-based sequence tagging methods with contextual topological data. Our methods do not depend on access to the underlying feature creation method, nor on external knowledge bases. Once these local topological descriptors are computed, they offer the potential for reuse across different tasks, thereby mitigating the initial computational investment. One limitation lies in our method still requiring labels

on the seed dataset. Though our results in the case study hint at a correlation between dialogue terms and higher Wasserstein norms, we have yet to establish a clear-cut purely feature based criterion for distinguishing terms from non-terms in dialogue data.

Looking ahead, we conjecture that the utility of our approach extends beyond the term extraction task investigated in this study. Given its generic design and challenge, it is plausible that it is applicable to other language models and modalities. Although our empirical evaluations have been confined to masked language models, the difficulty of the term extraction task provides optimism that our method could be advantageous in other scenarios where understanding the relation between individual data points and a datastore is critical.

**Reproducibility Statement** The MultiWOZ2.1 and SGD datasets are publicly available through the ConvLab-3 unified data format (Zhu et al., 2023), and we release our preprocessing, local topological feature computation and tagging model training code.<sup>2</sup>

## 5 Limitations

The experiments have been confined to a small masked language model (RoBERTa base model). Our proposed method can be applied to embedding spaces derived from causal LLMs (BehnamGhader et al., 2024), but current state-of-the-art models typically produce latent spaces with significantly larger embedding dimension. This has great influence on the computational complexity required to generate the contextual topological features. While our BIO-tagger can be trained on a single V100 GPU with 16 GB of VRAM in 2 hours, one should note that for efficiently precomputing the nearest neighbors in the contextual topological feature extraction, the embedded datastore needs to fit into the graphic card memory. This one-off neighborhood computation is not an issue for the MultiWOZ2.1 training set datastore, but might limit applications to larger corpus sizes. One possible remedy could be given by applying embedding space dimension reduction techniques such as (Kusupati et al., 2024) to the datastore before computing our topological features.

<sup>2</sup><https://gitlab.cs.uni-duesseldorf.de/general/dsml/tda4contextualembeddings-public>

## Acknowledgments

BMR and RV are supported by funds from the European Research Council (ERC) provided under the Horizon 2020 research and innovation programme (Grant agreement No. STG2018 804636) as part of the DYMO project. CVN and HL are supported by the Ministry of Culture and Science of North Rhine-Westphalia within the framework of the Lamarr Fellow Network. MH and SF are supported by funding provided by the Alexander von Humboldt Foundation in the framework of the Sofja Kovalevskaja Award endowed by the Federal Ministry of Education and Research. Computational infrastructure and support were provided by Google Cloud. We want to thank the anonymous reviewers whose comments improved the quality of our paper.

## References

- Henry Adams, Tegan Emerson, Michael Kirby, Rachel Neville, Chris Peterson, Patrick Shipman, Sofya Chepushtanova, Eric Hanson, Francis Motta, and Lori Ziegelmeier. 2017. [Persistence Images: A Stable Vector Representation of Persistent Homology](#). *Journal of Machine Learning Research*, 18(8):1–35.
- Arthur Amalvy, Vincent Labatut, and Richard Dufour. 2023. [The Role of Global and Local Context in Named Entity Recognition](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 714–722, Toronto, Canada. Association for Computational Linguistics.
- Ulrich Bauer. 2021. [Ripser: efficient computation of Vietoris-Rips persistence barcodes](#). *J. Appl. Comput. Topol.*, 5(3):391–423.
- Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. 2024. [LLM2Vec: Large Language Models Are Secretly Powerful Text Encoders](#). *Preprint*, arXiv:2404.05961.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2013. [Representation Learning: A Review and New Perspectives](#). *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828.
- Bradley CA Brown, Anthony L. Caterini, Brendan Leigh Ross, Jesse C Cresswell, and Gabriel Loaiza-Ganem. 2023. [Verifying the Union of Manifolds Hypothesis for Image Data](#). In *The Eleventh International Conference on Learning Representations*.
- Peter Bubenik, Michael Hull, Dhruv Patel, and Benjamin Whittle. 2020. [Persistent homology detects curvature](#). *Inverse Problems*, 36(2):025008, 23.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - A Large-Scale Multi-Domain Wizard-of-Oz Dataset for Task-Oriented Dialogue Modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Xingyu Cai, Jiaji Huang, Yuchen Bian, and Kenneth Church. 2021. [Isotropy in the Contextual Embedding Space: Clusters and Manifolds](#). In *International Conference on Learning Representations*.
- Gunnar Carlsson. 2014. [Topological pattern recognition for point cloud data](#). *Acta Numerica*, 23:289–368.
- Gunnar Carlsson, Tigran Ishkhanov, Vin De Silva, and Afra Zomorodian. 2008. [On the Local Behavior of Spaces of Natural Images](#). *International Journal of Computer Vision*, 76:1–12.
- Gunnar Carlsson and Mikael Vejdemo-Johansson. 2021. [Topological Data Analysis With Applications](#). Cambridge University Press.
- Andy Coenen, Emily Reif, Ann Yuan, Been Kim, Adam Pearce, Fernanda Viégas, and Martin Wattenberg. 2019. [Visualizing and measuring the geometry of BERT](#). In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Red Hook, NY, USA. Curran Associates Inc.
- David Cohen-Steiner, Herbert Edelsbrunner, John Harer, and Yuriy Mileyko. 2010. [Lipschitz functions have  \$L\_p\$ -stable persistence](#). *Foundations of Computational Mathematics. The Journal of the Society for the Foundations of Computational Mathematics*, 10(2):127–139.
- 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.
- Paweł Dlotko. 2017. [Persistence representations](#). In *GUDHI User and Reference Manual*. GUDHI Editorial Board.
- Herbert Edelsbrunner and John Harer. 2010. [Computational Topology: An Introduction](#). American Mathematical Society.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. [MultiWOZ 2.1: A Consolidated Multi-Domain Dialogue Dataset with State Corrections and State Tracking Baselines](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages

- 422–428, Marseille, France. European Language Resources Association.
- Kawin Ethayarajh. 2019. [How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings.](#) 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 55–65, Hong Kong, China. Association for Computational Linguistics.
- Michael Heck, Nurul Lubis, Carel van Niekerk, Shutong Feng, Christian Geischauser, Hsien-Chin Lin, and Milica Gašić. 2022. [Robust Dialogue State Tracking with Weak Supervision and Sparse Data.](#) *Transactions of the Association for Computational Linguistics*, 10:1175–1192.
- Alexander Jakubowski, Milica Gašić, and Marcus Zibrowius. 2020. [Topology of Word Embeddings: Singularities Reflect Polysemy.](#) In *Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics*, pages 103–113, Barcelona, Spain (Online). Association for Computational Linguistics.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. [Billion-Scale Similarity Search with GPUs.](#) *IEEE Transactions on Big Data*, 7(3):535–547.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. [Generalization through Memorization: Nearest Neighbor Language Models.](#) In *International Conference on Learning Representations*.
- Laida Kushnareva, Daniil Cherniavskii, Vladislav Mikhailov, Ekaterina Artemova, Serguei Barannikov, Alexander Bernstein, Irina Piontkovskaya, Dmitri Piontkovski, and Evgeny Burnaev. 2021. [Artificial Text Detection via Examining the Topology of Attention Maps.](#) In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 635–649, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham Kakade, Prateek Jain, and Ali Farhadi. 2024. [Matryoshka representation learning.](#) In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Théo Lacombe, Marco Cuturi, and Steve Oudot. 2018. [Large Scale computation of Means and Clusters for Persistence Diagrams using Optimal Transport.](#) In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18*, page 9792–9802, Red Hook, NY, USA. Curran Associates Inc.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. [Neural Architectures for Named Entity Recognition.](#) In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 260–270, San Diego, California. Association for Computational Linguistics.
- 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.](#) *Preprint*, arXiv:1907.11692.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled Weight Decay Regularization.](#) In *International Conference on Learning Representations*.
- Nikola Milosavljević, Dmitriy Morozov, and Primož Škraba. 2011. [Zigzag persistent homology in matrix multiplication time.](#) In *Computational geometry (SCG'11)*, pages 216–225. ACM, New York.
- Nina Otter, Mason A Porter, Ulrike Tillmann, Peter Grindrod, and Heather A Harrington. 2017. [A roadmap for the computation of persistent homology.](#) *EPJ Data Science*, 6(1).
- Rrubaa Panchendrarajan and Aravindh Amaran. 2018. [Bidirectional LSTM-CRF for named entity recognition.](#) In *Proceedings of the 32nd Pacific Asia Conference on Language, Information and Computation*, Hong Kong. Association for Computational Linguistics.
- Ilan Perez and Raphael Reinauer. 2022. [The Topological BERT: Transforming Attention into Topology for Natural Language Processing.](#) *Preprint*, arXiv:2206.15195.
- 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.
- Liang Qiu, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. [Structure Extraction in Task-Oriented Dialogues with Slot Clustering.](#) *Preprint*, arXiv:2203.00073.
- Lance Ramshaw and Mitch Marcus. 1995. [Text Chunking using Transformation-Based Learning.](#) In *Third Workshop on Very Large Corpora*.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. [Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset.](#) In *The Thirty-Fourth AAAI Conference on Artificial Intelligence*,

- AAAI 2020, *The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020*, New York, NY, USA, February 7-12, 2020, pages 8689–8696. AAAI Press.
- Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. 2020. [Masked Language Model Scoring](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2699–2712, Online. Association for Computational Linguistics.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. [BERT rediscovers the classical NLP pipeline](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.
- The GUDHI Project. 2015. *GUDHI User and Reference Manual*. GUDHI Editorial Board.
- Eduard Tulchinskii, Kristian Kuznetsov, Laida Kushnareva, Daniil Cherniavskii, Sergey I. Nikolenko, Evgeny Burnaev, Serguei Barannikov, and Irina Piontkovskaya. 2023. [Intrinsic Dimension Estimation for Robust Detection of AI-Generated Texts](#). 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*.
- Lucrezia Valeriani, Diego Doimo, Francesca Cuturello, Alessandro Laio, Alessio Ansuini, and Alberto Cazaniga. 2023. [The geometry of hidden representations of large transformer models](#). 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*.
- Julius Von Rohrscheidt and Bastian Rieck. 2023. [Topological Singularity Detection at Multiple Scales](#). In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org.
- Renato Vukovic, Michael Heck, Benjamin Ruppik, Carel van Niekerk, Marcus Zibrowius, and Milica Gašić. 2022. [Dialogue Term Extraction using Transfer Learning and Topological Data Analysis](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 564–581, Edinburgh, UK. Association for Computational Linguistics.
- Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. [GPT-NER: Named Entity Recognition via Large Language Models](#). *Preprint*, arXiv:2304.10428.
- Frank F. Xu, Uri Alon, and Graham Neubig. 2023. [Why do Nearest Neighbor Language Models Work?](#) In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org.
- Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabás Póczos, Ruslan Salakhutdinov, and Alexander J Smola. 2017. [Deep sets](#). In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 3394–3404, Red Hook, NY, USA. Curran Associates Inc.
- Qi Zhu, Christian Geischauser, Hsien-chin Lin, Carel van Niekerk, Baolin Peng, Zheng Zhang, Shutong Feng, Michael Heck, Nurul Lubis, Dazhen Wan, Xiaochen Zhu, Jianfeng Gao, Milica Gašić, and Minlie Huang. 2023. [ConvLab-3: A flexible dialogue system toolkit based on a unified data format](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 106–123, Singapore. Association for Computational Linguistics.

## A Appendix

### A.1 Implementation Details

In our term extraction applications, the ambient vector datastore comprises a collection of vectors with cardinality in the millions, making the computation of neighborhoods a major computational bottleneck. To alleviate this issue, we employ the Facebook AI Similarity Search (faiss) module (Johnson et al., 2021) to precompute neighborhood indices using GPU acceleration. These indices can be reused in subsequent computations of our local measurements at varying scales. We obtain the neighborhood indices and distances for 1 024 neighbors using the faiss.IndexFlatL2 build from the MultiWOZ2.1 training datastore. This neighborhood cache allows extraction of the codensity measurements and the vectors required to subsequently compute persistence images for neighborhood size  $n = 128$ . Loading the 2 739 744 many 768-dimensional vectors from the roberta-base MultiWOZ2.1 training datastore into the faiss index requires approximately 8 GB of GPU memory.

The faiss library does not currently offer GPU support for range-based nearest neighbor search. This makes it infeasible to compute range-based neighborhoods at the scale of our dataset for the methods described in Von Rohrscheidt and Rieck (2023). This limitation is especially critical because our BIO-tagger requires topological features for each input token in each context.

Another computational challenge lies in the local persistent homology computations, which become a bottleneck when the goal is training a BIO-tagger based on the resulting features. To address this, we precompute and store the topological features and their vectorizations, including both persistence

images and Wasserstein norms. We use the Ripser library (Bauer, 2021) for computing the persistence modules for  $H_0$  and  $H_1$  with  $\mathbb{F}_2$ -coefficients w.r.t. cosine distance from the precomputed neighborhoods, and GUDHI (The GUDHI Project, 2015; Dlotko, 2017) for the vectorization and persistence representation. The Wasserstein norms, i.e., the order-1 Wasserstein distances between the neighborhood persistence diagrams and the empty diagram with Euclidean ground metric, are computed separately for the  $H_0$  and  $H_1$  persistence diagrams using the GUDHI library.

For the persistence image vectorization of the  $H_0$ -persistence module, we decide on a bandwidth of 0.01, image range on the  $y$ -axis of  $[0.0, 1.0]$ , resolution of  $1 \times 100$  and weight each persistence homology generator by its  $y$ -value. We restrict our computations to 0-dimensional and 1-dimensional persistent homology. This is not only due to Ripser's optimizations, which result in a faster runtime, but also to circumvent the potential exponential increase in the number of simplices in the filtered complex when considering higher dimensional Vietoris-Rips complexes on a point cloud.

# Adaptive Open-Set Active Learning with Distance-Based Out-of-Distribution Detection for Robust Task-Oriented Dialog System

Sai Keerthana Goruganthu  
University of Missouri  
sgmhz@umsystem.edu

Roland Oruche  
University of Missouri  
rro2q2@umsystem.edu

Prasad Calyam  
University of Missouri  
calyamp@missouri.edu

## Abstract

The advancements in time-efficient data collection techniques such as active learning (AL) have become salient for user intent classification performance in task-oriented dialog systems (TODS). In realistic settings, however, traditional AL techniques often fail to efficiently select targeted in-distribution (IND) data when encountering newly acquired out-of-distribution (OOD) user intents in the unlabeled pool. In this paper, we introduce a novel adaptive open-set AL framework viz., “AOSAL” for TODS that combines a distance-based OOD detector using an adaptive false positive rate threshold along with an informativeness measure (e.g., entropy) to strategically select informative IND data points in the unlabeled pool. Specifically, we utilize the adaptive OOD detector to classify and filter out OOD samples from the unlabeled pool, then prioritize the acquisition of classified IND instances based on their informativeness scores. To validate our approach, we conduct experiments that display our framework’s flexibility and performance over multiple distance-based approaches and informativeness measures against deep AL baselines on benchmark text datasets. The results show that our AOSAL consistently outperforms the baselines on IND classification and percentage of acquired IND samples, demonstrating its ability to improve robustness of task-oriented dialog systems.

## 1 Introduction

Recent advances in time-efficient data collection techniques such active learning (AL) (Settles, 2009; Ren et al., 2021) show the promise of significantly improving the performance of task-oriented dialog system (TODS) for tasks related to user intent classification (Zhang and Zhang, 2019; Wu et al., 2024). The time-efficient AL techniques not only improve the model accuracy of the TODS, but also help reduce the annotation budget of human anno-

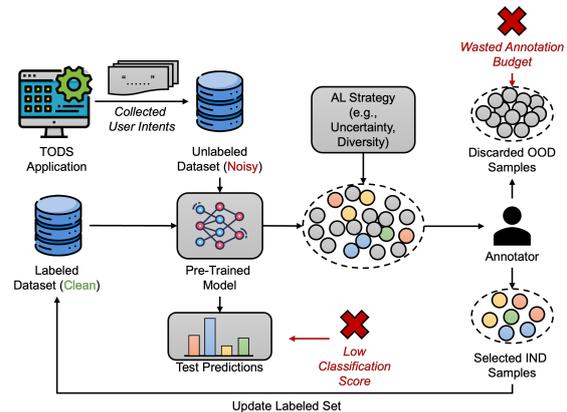


Figure 1: The challenges of traditional AL methods when encountering OOD instances from newly collected user intents in the unlabeled pool which includes low classification score and wasted annotation budget.

tators when querying the most informative samples that accelerate training performance.

In real-world applications, however, existing AL methods often struggle to select in-distribution (IND) data from unlabeled pools containing out-of-distribution (OOD) user intents, leading to inefficiencies in the learning process. Figure 1 illustrates the challenge of employing standard AL frameworks in a TODS application, where an unlabeled dataset of collected user intents are noisy due to instances that are OOD. Typical queries using e.g., uncertainty (Lewis, 1995) and diversity-based methods (Nguyen and Smeulders, 2004) are prone to selecting a high number of OOD instances, which in turn can waste the annotation budget of the human annotator. Consequently, this can lead to low classification performance, and more concretely, incorrect dialog responses if there are insufficient amount of IND samples selected for training, as shown in the example scenario in Figure 2.

Previous works have investigated robust AL frameworks in the context of open-set recognition (Scheirer et al., 2012). The work in (Yang et al., 2024) develops a progressive active learn-

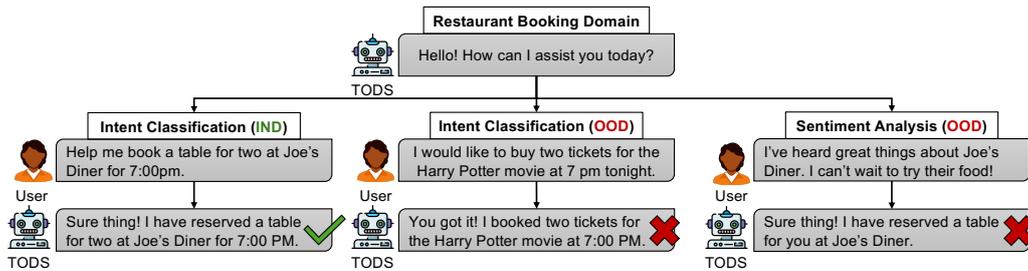


Figure 2: Example scenarios of task-oriented dialog systems (TODS) handling user intents in a restaurant booking domain. TODS can provide incorrect and unwarranted responses when encountering OOD intents.

ing framework that implements an OOD detector for filtering OOD instances in the unlabeled pool. Although other works have proposed similar methods related to open-set active learning (Du et al., 2021; Park et al., 2022; Ning et al., 2022), these works are mainly applied in the computer vision space, and are unsuited for NLP tasks in TODS. In addition, while AL frameworks in the NLP space such as CAL (Margatina et al., 2021) and CounterAL (Deng et al., 2023) address OOD generalization challenges, they are not practical to open-set AL where the unlabeled dataset contains a mixture of IND and OOD samples. Given the emergence and applicability of TODS in various application domains e.g., healthcare, banking, there presents a need to develop robust AL strategies that avoid OOD instances while also acquiring informative IND instances that improve model training.

In this paper, we present a novel adaptive open-set AL framework viz., “AOSAL” for TODS that combines an adaptive distance-based OOD detector with informative sampling measures (e.g., uncertainty, diversity) to effectively acquire IND samples in the unlabeled pool. Our OOD detector features a normalized score function that classifies unlabeled samples based on their distance to each class in the IND labeled space. We enable our OOD detector to be sensitive to distribution shifts by employing an adaptive threshold that is controlled by using a predetermined false positive rate (FPR) over the OOD detection performance. Based on the prioritization of classified pseudo-IND samples, we then leverage sampling measures for selecting the most informative instances for annotation. In addition, we demonstrate the flexibility of our AOSAL approach to multiple distance-based functions (Podolskiy et al., 2021; Frogner et al., 2015) and informative measures (Lewis, 1995).

We perform experiments to validate our AOSAL framework over four NLP benchmark related to intent classification (Larson et al., 2019; Gangal et al.,

2020), and sentiment analysis (Maas et al., 2011; Aslam et al., 2020), comparing its performance against several deep AL baselines that are based on uncertainty, diversity, and hybrid-based approaches. Experimental results suggest that our AOSAL approach consistently outperforms the baselines on metrics such as IND classification accuracy, and percentage of acquired IND/OOD samples.

The remainder of the paper is organized as follows: Section 2 describes related work. In Section 3, we detail the AOSAL methodology. In Section 4, we detail the experimental setup and provide the main results against AL baselines. An analysis to test the robustness of AOSAL is presented in Section 5. Section 6 discusses the limitations of our approach, and finally, Section 7 concludes our work.

## 2 Related Work

### 2.1 Active Learning

Recent advancements in active learning have leveraged pool-based sampling (Settles, 2009), where an agent can select and query a large set of instances to the oracle (i.e., human annotator) from the unlabeled pool. Common methods on the selection process, or *query strategy*, based on how informative a given sample is, include uncertainty (Lewis, 1995; Settles, 2009) and diversity (Nguyen and Smeulders, 2004; Sener and Savarese, 2018) methods. Uncertainty strategies such as Entropy (Settles, 2009) and Least Confidence (Lewis, 1995) aim to select a set of instances from the unlabeled pool in cases where the model is least confident in its prediction. While uncertainty can maintain low computational complexity, diversity-based methods such as Coreset (Sener and Savarese, 2018) and clustering-based methods (Nguyen and Smeulders, 2004) select samples that better represent the distribution of the unlabeled pool.

The advent of deep learning in AL has en-

abled batch-mode active learning (Kirsch et al., 2019), where the sampling of unlabeled instances in batches are sent to the oracle for labeling. Authors in (Kirsch et al., 2019) extend Bayesian Active Learning by Disagreement (BALD) (Houlsby et al., 2011), by presenting BatchBALD, which alleviates the time complexity of calculating the mutual information between an individual sample and the model parameters. Batch-model AL has also engendered recent work in hybrid-based approaches (Yin et al., 2017; Zhdanov, 2019; Ash et al., 2020; Shui et al., 2020). The work in (Ash et al., 2020) develops BADGE, a hybrid query strategy that robustly selects samples by leveraging both the prediction uncertainty and diverse samples from the hallucinated gradient space of the model. Despite such advancements, these methods often fail to improve IND classification performance and limit the oracle’s annotation budget when there is a distribution mismatch between the labeled and the unlabeled set. Thus, traditional active learning methods are not feasible for training agents within TODS systems used in real-world applications.

## 2.2 Open-Set Active Learning

Previous works have aimed to develop AL methods in the context of open-set recognition (Scheirer et al., 2012) that is more suitable for realistic scenarios where there presents a distribution mismatch in the unlabeled pool (Kothawade et al., 2021; Du et al., 2021; Ning et al., 2022; Park et al., 2022; Safaei et al., 2024; Yang et al., 2024). The work in (Du et al., 2021) develops CCAL, which utilizes contrastive learning to extract semantic and distinctive features in the unlabeled pool. The authors propose an AL error when selecting invalid (OOD) samples, which are segmented between valid and invalid query errors. Other works such as in (Kothawade et al., 2021) develop a unified AL framework that addresses OOD samples in the unlabeled pool by utilizing submodular conditional mutual information that jointly models the similarity between the query set and batch of unlabeled samples and their dissimilarity between a conditioning set.

More recent work on open-set AL further addresses distribution mismatches by utilizing OOD samples for training in the unlabeled pool. For instance, progressive active learning (PAL) (Yang et al., 2024) samples both pseudo-IND and pseudo-OOD samples to simultaneously train the ID classifier and a proposed OOD detector using a one-

vs-all classifier. Authors in (Park et al., 2022) demonstrate that balancing between purity (i.e., distinguishing between collected IND and OOD instances), and informativeness (i.e., uncertainty, diversity) consistently improves the classifier accuracy under various noise (OOD) ratio in the unlabeled pool. Similarly, other works such as LfOSA (Ning et al., 2022) and EOAL (Safaei et al., 2024) leverage both known (IND) and unknown (OOD) data instances to effectively informative IND samples while avoid OOD samples during AL rounds. Despite these advancements, the majority of methods from existing work in open-set AL are mainly tailored to the computer vision domain.

Our AOSAL framework for robust TODS is novel because it: (i) detects OOD instances (e.g., user intents) using distance-based approaches coupled with an adaptive threshold to maintain a low false positive rate in text-based datasets, and (ii) utilizes measures over unlabeled instances classified as IND for improving IND accuracy on the labeled set. In addition, we demonstrate that our AL framework can be extended to multiple distance-based approaches and informative measures.

## 3 Methodology

In this section, we describe the problem formulation for open-set AL and then detail the overview and components of our proposed AOSAL framework.

### 3.1 Problem Formulation

We define a TODS problem for identifying user intents as a  $\mathcal{K}$ -class classification task. An IND labeled dataset  $\mathcal{D}_L$  has an input space  $\mathcal{X}$  and a corresponding output label space  $\mathcal{Y} \in \{1, \dots, \mathcal{K}\}$  of  $\mathcal{K}$  IND classes, which are independently and identically distributed (i.i.d.) from  $\mathcal{D}_L$ . The full dataset is defined as  $\mathcal{D}_L = \{(x_i, y_i)\}_{i=1}^{N_L}$ , where  $N_L$  is the length of the initial labeled training set.

We denote an unlabeled dataset as  $\mathcal{D}_U$  for the purposes of re-training the TODS over newly collected user intents. Formally, the unlabeled dataset is defined as  $\mathcal{D}_U = \{(x_j)\}_{j=1}^{N_U}$ , where  $N_U$  is the length of the unlabeled set. We also denote  $N_L \ll N_U$ , highlighting the substantially larger pool of unlabeled dataset  $\mathcal{D}_U$  compared to  $\mathcal{D}_L$ . In real-world AL scenarios, there often presents a distribution mismatch in the unlabeled pool due noisy, OOD class samples. Thus, we define our problem to an open-set AL in a pool-based setting, where

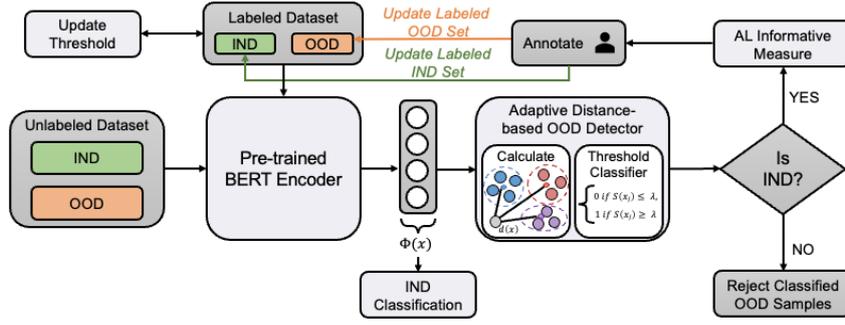


Figure 3: Main Architecture of our AOSAL framework. A pre-trained BERT model encodes samples unlabeled from the unlabeled pool and classifies them using the OOD detector. Classified IND samples are queried using an informative measure and fed to the annotator for updating the labeled set and updating the FPR-controlled threshold.

the unlabeled pool contains a both IND and OOD input samples (i.e.,  $\mathcal{D}_U = \mathcal{X}^{IND} \cup \mathcal{X}^{OOD}$ ) with a label space  $\mathcal{Y}^{IND}$  and  $\mathcal{Y}^{OOD}$ , respectively, and  $\mathcal{Y}^{IND} \cap \mathcal{Y}^{OOD} = \emptyset$ . In simple terms, a sample  $x_j \in \mathcal{D}_U$  may belong to an IND or OOD class in the unlabeled dataset.

Within the AL loop, an AL strategy queries a batch of samples of size  $b$  from  $\mathcal{D}_U$  to form into a query set  $Q$ , which can consist a mixture of IND and OOD samples (i.e.,  $Q = \mathcal{D}_U^{IND} \cup \mathcal{D}_U^{OOD}$ ). This query set is then fed to the human annotator (i.e., oracle) for labeling and updating the initial training set  $\mathcal{D}_L$ .

### 3.2 Adaptive Open-Set Active Learning

We present our novel adaptive open-set AL (AOSAL) framework that couples an adaptive distance-based OOD detector with informativeness measures for efficiently managing OOD instances in the unlabeled pool. We display the main architecture of our AOSAL approach in Figure 3. Unlike previous AL frameworks for NLP (Margatina et al., 2021; Deng et al., 2023), we utilize the unlabeled OOD instances that are queried to the oracle for annotation to improve our distance-based OOD detector with an adaptive threshold controlled by a pre-defined false positive rate (FPR). In the following, we formalize the main components of AOSAL and detail the full sampling procedure in the AOSAL cycle.

#### 3.2.1 Distance-based OOD Detector

To address the challenge of efficiently utilizing annotation resources in AL contexts, we have developed a distance-based OOD detector. This detector classifies an unlabeled sample  $x_j$  as either in-distribution (IND) or out-of-distribution (OOD) based on an adaptive threshold. The classification

decision is made according to the following rule:

$$\text{Classify}(x_j) = \begin{cases} \text{accept}, & \text{if } S(x_j) \leq \lambda, \\ \text{reject}, & \text{if } S(x_j) > \lambda \end{cases} \quad (1)$$

where,  $\lambda$  is the threshold in the range  $[0, 1]$  and  $x_j$  denotes the  $j$ -th unlabeled sample in  $\mathcal{D}_U$ . The threshold separates IND samples, which score at or below the threshold, from OOD samples, which score above it. The scoring function  $S(x_j)$  is designed to measure the proximity of  $x_j$  to the nearest class in the labeled dataset  $\mathcal{D}_L$ .

The scoring function  $S(x_j)$  is conceptualized to enhance the selection of IND samples from the unlabeled dataset  $\mathcal{D}_U$  by calculating the minimal distances between  $x_j$  and each class represented in  $\mathcal{D}_L$ . It is defined as:

$$S(x_j) = \min_{k \in \mathcal{K}} d(x_j, \mu_{x_k}), \quad (2)$$

where,  $\mu_{x_k}$  represents the mean feature vector of class  $k$  from the set of classes  $\mathcal{K}$  in  $\mathcal{D}_L$ . This approach ensures that the scoring function remains adaptable across various distance metrics, each potentially having different mathematical properties and score ranges.

To facilitate uniformity in classification regardless of the absolute scale of distance values, we normalize the scores to a  $[0, 1]$  range:

$$S(x_j) = \frac{S(x_j)}{\max_{x_j \in \mathcal{D}_U} (S(x_j))}. \quad (3)$$

This normalization not only standardizes the score across various distance metrics but also aligns with the thresholding approach to identify between IND and OOD samples. In our experiments, we

utilize the Mahalanobis distance (Podolskiy et al., 2021) and Wasserstein distance (Frogner et al., 2015) to compute  $S(x_j)$ , which are chosen for their robustness in capturing the geometric nuances of data distributions. The specific formulas and their application are detailed further in the Appendix A.1, ensuring a comprehensive exposition of our distance-based OOD detection methodology.

### 3.2.2 Adaptive Threshold

Classifying OOD instances using a constant threshold value creates significant challenges in maintaining high OOD accuracy in real-world settings. This is particularly evident in newly collected unlabeled data in TODS applications, where IND and potential OOD samples can cause a distribution shift. Consequently, this can lead to high false positives (i.e., detecting OOD samples as IND) and ultimately negatively impact the annotation budget as more OOD samples are naturally acquired.

To address this, we implement an adaptive threshold mechanism controlled by a pre-defined false positive rate (FPR), which is essential for maintaining classification integrity under varying data conditions. The FPR is defined as the ratio of IND instances mistakenly classified as OOD to the total number of true negative instances. It is calculated as:

$$FPR = \frac{FP}{FP + TN} \quad (4)$$

where,  $FP$  is the number of false positives, (i.e., IND samples incorrectly classified as OOD), and  $TN$  is the number of true negatives, IND (i.e., samples correctly classified as IND).

To maintain system accuracy and adapt to new data, the adaptive threshold  $\lambda$  is adjusted based on the FPR, which is calculated as:

$$\lambda = FPR(\mathcal{D}_{\text{val}}, \alpha) \quad (5)$$

where,  $\mathcal{D}_{\text{val}}$  the validation set with a mixture of labeled IND and OOD samples (i.e.,  $\mathcal{D}_{\text{val}} = \mathcal{D}_L^{\text{IND}} \cup \mathcal{D}_L^{\text{OOD}}$ ), and  $\alpha$  is the predetermined FPR rate. The benefit of ensuring an adaptive threshold is consistent with the pre-defined FPR that mitigates the risk of the OOD detector from producing high false positives on the unlabeled dataset during AL acquisition. Furthermore,  $\lambda$  is dynamically calibrated to ensure that the proportion of false positives does not exceed  $\alpha$ . This dynamic adjustment is conducted through a meticulous analysis of the

model’s scoring outputs on each validation sample  $x_{\text{val}} \in \mathcal{D}_{\text{val}}$ . The threshold  $\lambda$  is then set such that:

$$\alpha = \frac{|\{x_{\text{val}} \in \mathcal{D}_L : S(x_{\text{val}}) > \lambda \text{ and } y_{\text{val}} = 0\}|}{|\{x_{\text{val}} \in \mathcal{D}_L : y_{\text{val}} = 0\}|} \quad (6)$$

where,  $S(x_{\text{val}})$  is the score function applied to each validation sample, and  $y_{\text{val}}$  indicates the sample’s label, with a label of 0 signifying an IND sample.

The validation set plays a crucial role in accurately updating the adaptive threshold for effective OOD detection in TODS. The informativeness metric derived from calculating uncertainty or diversity on the validation set is utilized to fine-tune the model and threshold. Furthermore, the validation set is continuously updated with newly annotated OOD samples, ensuring that the OOD detector remains up-to-date and capable of handling evolving data patterns. This mechanism enhances the robustness and reliability of TODS, enabling them to maintain high accuracy in OOD detection under varying situations and adapt to dynamic data shifts.

### 3.3 AOSAL Sampling Procedure

The overall AL sampling process for our proposed AOSAL framework is shown in Algorithm 1. We start by training a deep learning model  $M_\theta$  on  $\mathcal{D}_L$  at the initial iteration  $t = 0$  to obtain  $M_{\theta_t}$ . During the validation, we leverage our distance-based OOD detector that computes a score from Equation 3 over samples on the validation set. The normalized scores from the OOD detector are then used to set the initial threshold  $\lambda_t$  based on a pre-defined FPR  $\alpha$  over the validation set  $\mathcal{D}_{\text{val}}$ .

Within our AL loop, we extract the features for each unlabeled sample in  $\mathcal{D}_U$  computed by  $M_{\theta_t}$  as input to our OOD detector using a normalized distance-based function from Equation 3 that computes the distance based scores for classification. After classifying the samples based on Equation 1, we ignore the classified OOD samples and focus on acquiring IND samples using informative measures (e.g., uncertainty, diversity). Following this, human annotators refine these classifications, and the resulting samples consisting of both IND and potential OOD are updated in either the IND train set or the OOD validation set.

The iteration of the model is updated at  $t = t + 1$ , and the threshold is adjusted using the OOD detector with a controlled false positive rate (FPR) at  $\alpha$ . This process is repeated until the annotation budget

---

**Algorithm 1** Adaptive Open-Set Active Learning with Distance-Based OOD Detection

---

**Require:** Labeled IND dataset  $D_L^{IND}$ , labeled OOD dataset  $D_L^{OOD}$ , unlabeled pool  $D_U$ , validation set  $D_{val}$ , model  $M_\theta$ , encoder function  $\Phi$ , acquisition size  $b$ , labeling budget  $B$ , total query set  $Q$ , threshold function  $FPR$ , informativeness measure  $U$ , current iteration  $t$ .

- 1: Train  $M_{\theta_{t=0}}$  on  $D_L^{IND}$  for multi-classification
- 2:  $\lambda_{t=0} \leftarrow FPR(D_{val}, \alpha)$   $\triangleright$  Set initial threshold (Eq. 5)
- 3: **while**  $|Q| < B$  **do**
- 4:   **for each**  $x_j$  in  $D_U$  **do**
- 5:      $\mu_{x_k} \leftarrow \Phi(x_k), k = \{1, \dots, \mathcal{K}\}$
- 6:      $s_{x_j} \leftarrow S(\Phi(x_j), \mu_{x_k})$   $\triangleright$  From Eq. 3
- 7:     **if**  $s_{x_j} \leq \lambda$  **then**  $\triangleright$  IND label
- 8:        $A \leftarrow \{(x_j^{IND}, \hat{y}_j^{IND})\}$
- 9:     **end if**
- 10:   **end for**
- 11:   **for each**  $x_j$  in  $A$  **do**
- 12:      $Q \leftarrow \operatorname{argmax}_{x_j \in U} U(x_j), |Q| = b$   $\triangleright$  Select  $b$  instances with highest informative scores.
- 13:   **end for**
- 14:    $D_L^{IND} \leftarrow D_L^{IND} \cup \{Q^{IND} \setminus D_U\}$   $\triangleright$  Update train set with acquired IND samples
- 15:    $D_{val}^{OOD} \leftarrow D_{val}^{OOD} \cup \{Q^{OOD} \setminus D_U\}$   $\triangleright$  Update validation set with acquired OOD samples
- 16:   Train the model  $M_{\theta_{t+1}}$  on  $D_L^{IND}$
- 17:   Update  $\lambda$  using Eq. 5 on updated  $D_{val}$
- 18:    $t \leftarrow t + 1$
- 19: **end while**
- 20: **return**  $M_{\theta_t}, \lambda_t$   $\triangleright$  Return updated model and threshold

---

$B$  is exhausted, ensuring continuous refinement of the model’s performance and the threshold.

## 4 Experiments and Results

In this section, we provide our experimental setup for open-set AL on benchmark NLP datasets and provide the main results of our AOSAL approach against baseline AL datasets.

### 4.1 Datasets

We validated our AOSAL framework over NLP datasets related to topic classification and sentiment analysis. These datasets are integral for validating the model’s efficacy in handling both IND and OOD samples within varied textual contexts. For topic classification, we test over the CLINC-

Full (Larson et al., 2019) dataset with 150 classes and Real Out-of-distribution Sentences From Task-Oriented Dialog (ROSTD) (Gangal et al., 2020) dataset with 12 classes, which both include OOD samples. For sentiment analysis, we utilize the Stanford Sentiment Treebank (SST)-2 (Aslam et al., 2020) dataset with only 2 classes each for the positive and negative sentiments. In our experiments, we set one dataset to the IND class and the other dataset to the OOD class. For instances where CLINC-Full and ROSTD are assigned to the IND class, we join the remaining OOD samples along with the assigned OOD dataset. We provide the full dataset description and partitions in Appendix A.2.

### 4.2 Baselines and Implementation Details

We compare our approach against five AL baselines that include state-of-the-art approaches for different query strategies. Specifically, we test an uncertainty sampling method, namely Entropy (Joshi et al., 2009), for which samples with the lowest confidence in the model’s predictive probability are selected. For diversity sampling, we test our approach against BERT – KM from the works of (Yuan et al., 2020), where they performs  $k$ -means clustering over the L2-normalized BERT embeddings to select diverse samples in the unlabeled feature space. For hybrid sampling, we compare our approach against BADGE (Ash et al., 2020), which is known to be an AL state-of-the-art method. In addition to state-of-the-art AL methods, we include CAL (Margatina et al., 2021) that selects “contrastive” unlabeled samples based on their feature similarity and divergent predictive probability. Lastly, we include Random sampling as a baseline for randomly acquiring instances in the unlabeled pool.

We implement our approach using a pre-trained BERT model (Devlin et al., 2019) from the HuggingFace library<sup>1</sup> as the backbone model for each approach in our experiments. While we opted to use BERT due to its reliable performance on natural language understanding tasks, our AOSAL framework can be extended to multiple model architectures for intent classification (Liu et al., 2019; Lan et al., 2019; He et al., 2021).

For each dataset, we use 10% of the train set as our initial labeled set  $D_L$  and use 10% of OOD samples and label them in the validation set for OOD detection. In addition, we set the noise ratio (i.e.,

---

<sup>1</sup>HuggingFace BERT model available at: <https://huggingface.co/bert-base-uncased>

percentage of OOD samples in the unlabeled pool) to 30%. This noise ratio presents a realistic consideration of the amount of noise that can be present in the unlabeled pool. During each AL iteration, we fine-tune  $\mathcal{D}_L$  with newly acquired IND samples from the unlabeled pool  $\mathcal{D}_U$ . We set the oracle labeling budget  $B$  to 25% percent of  $\mathcal{D}_U$  for a total number of 5 AL rounds. We pre-trained the base model over 5 epochs on CLINC-Full for training. We run experiments for each AL method 5 times over each dataset and report the average IND accuracy and the percentage of IND/OOD samples in the acquisition size  $|Q|$  for each AL iteration. We provide additional detail on the model implementation using BERT and relevant hyperparameters in Appendix A.3.

### 4.3 Main Results

The main results over the CLINC-Full (IND) and the SST-2 (OOD) dataset are shown in Figure 4. Figures 4a and 4b show the average test accuracy (90.1% ( $\pm 0.01$ )) and the average acquired IND (661.38 ( $\pm 6.66$ )) of our AOSAL approach across all AL iterations compared to the baselines, respectively. When averaging across all AL budgets, our approach shows significant improvement in acquiring IND samples compared to Entropy (57.82 ( $\pm 6.11$ )). The relatively low test accuracy performance for uncertainty methods such as Entropy (87.3% ( $\pm 0.004$ )) may be the result of selecting instances in the unlabeled pool that are least confident in its prediction, which causes Entropy to acquire more OOD instances, thus wasting the annotation budget.

In the scope of diversity- and hybrid-based methods, AOSAL shows comparable test accuracy performance to BADGE (90% ( $\pm 0.00$ )) and BERT-KM (91.1% ( $\pm 0.01$ )). While our AOSAL approach significantly shows higher acquired IND compared to BADGE (388.40 ( $\pm 3.94$ )) and BERT-KM (462.96 ( $\pm 8.62$ )), the high test accuracy results may indicate the benefit of acquiring a diverse set of samples in the unlabeled pool for improving model performance. In addition, CAL shows surprisingly low performance in both IND test accuracy (87.8% ( $\pm 0.01$ )) and average acquired IND (76.56 ( $\pm 10.65$ )) when handling OOD instances from the SST-2 dataset.

Furthermore, AOSAL shows comparable results in acquired IND to Random sampling (676.71 ( $\pm 2.32$ )) and an improvement in IND test accuracy performance (89.8% ( $\pm 0.01$ )). Since Random

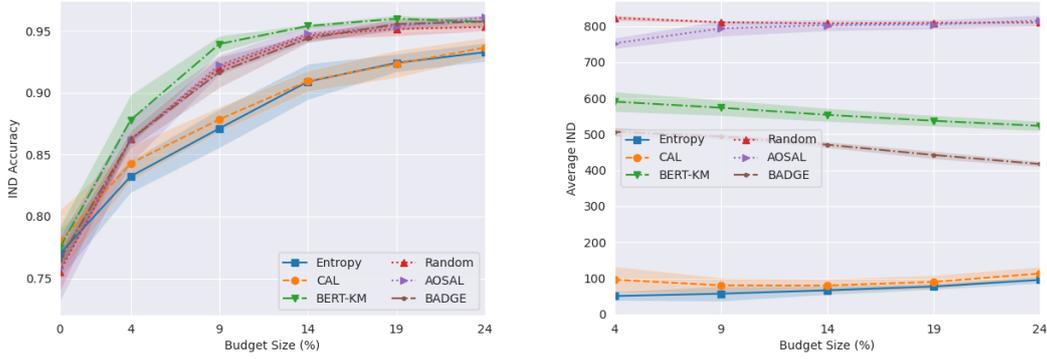
AL sampling follows a uniform distribution, it outperforms all baseline approaches when the amount of OOD instances in the unlabeled pool is considerably low (i.e., noise ratio at 30%). Despite this, the acquired IND performance does not always translate to high IND test accuracy, as indicated in Figure 4a (89.8% ( $\pm 0.01$ )). This is because the samples acquired may not always be informative in terms of uncertainty and diversity for effectively improving model performance.

Similar results on the consistency of AOSAL are shown in Figure 5. AOSAL shows comparable average performance to the baselines across all AL iterations in terms of IND test accuracy (91.2% ( $\pm 0.01$ )) in Figure 5a and acquired IND (681.20 ( $\pm 4.90$ )) in Figure 5b. Compared to the baselines such as Entropy (470.13 ( $\pm 61.09$ )), Random (678.52 ( $\pm 2.28$ )), and CAL (426.217 ( $\pm 35.91$ )), our AOSAL approach shows a higher amount of acquired IND averaged across all AL iterations. This in turn translates to comparable or higher accuracy on the IND test set. While the IND test accuracy results are comparable to other baselines such as BADGE (91.8% ( $\pm 0.01$ )) and BERT-KM (91.8% ( $\pm 0.00$ )), our AOSAL approach maintains a comparable accuracy as well as average acquired IND performance to BADGE (662.683 ( $\pm 10.31$ )) and BERT-KM (702.122 ( $\pm 4.98$ )) when encountering a variety of OOD instances in the unlabeled pool.

## 5 Analysis

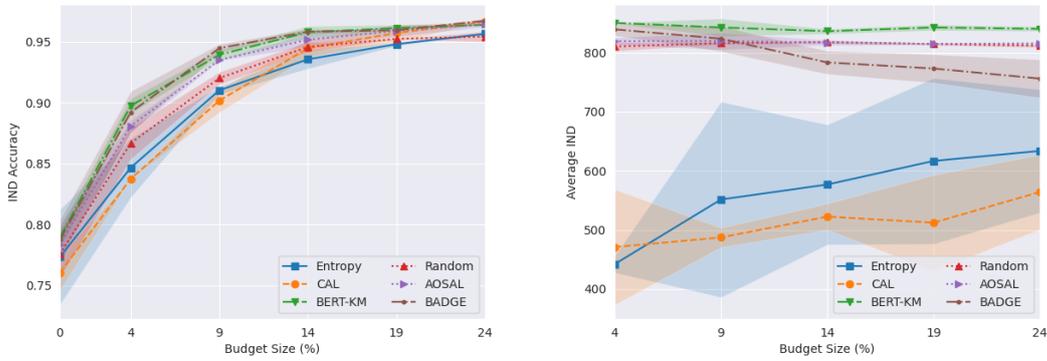
### 5.1 Ablation Study

We conduct an ablation study to check the AOSAL framework under varying budgets, analyzing how different distance metrics and OOD detection can influence IND accuracy. We compare six configurations of our framework, including AOSAL-CONST, which uses a constant threshold, and others such as AOSAL-NO-OOD, AOSAL-MAH-DIV, AOSAL-WAS-DIV, AOSAL-MAH-UNC, and AOSAL-WAS-UNC that use an adaptive FPR threshold, but differ in their application of Mahalanobis or Wasserstein distances and the incorporation of uncertainty and diversity metrics. To ensure fair and meaningful comparisons across all experimental settings, we utilize CLINC-Full as the IND data and SST-2 as the OOD data, with a fixed false positive rate (FPR) of 95%. This standardization helps maintain consistent experimental conditions throughout the study. Figure 6 shows that the con-



(a) CLINC (IND) and SST-2 (OOD) on IND test accuracy. (b) CLINC (IND) and SST-2 (OOD) on acquired IND.

Figure 4: Test accuracy results and averaged acquired IND on the each AL method over the CLINC-Full (IND) and SST-2 (OOD) dataset. Each method was ran 5 times with different seeds and the average accuracies were reported.



(a) CLINC (IND) and ROSTD (OOD) IND test accuracy. (b) CLINC (IND) and ROSTD (OOD) on acquired IND.

Figure 5: Test accuracy results and averaged acquired IND on the each AL method over the CLINC-Full (IND) and ROSTD (OOD) dataset. Each method was ran 5 times with different seeds and the average accuracies were reported.

figurations lacking OOD detection i.e., AOSAL-NO-OOD demonstrates a significant drop in the model’s performance, highlighting the crucial role of effective OOD detection mechanisms in enhancing the overall accuracy of the system. This analysis confirms the robustness and versatility of our AOSAL framework in adapting to different operational constraints and validates the utility of advanced distance measures for OOD detection in the AL environment.

## 5.2 Threshold Analysis

We conducted a detailed threshold analysis to evaluate the impact of various FPR thresholds on IND accuracy. Our study systematically explored the performance implications of different FPR levels including 90%, 95%, and 97% across multiple datasets. The dataset configurations, detailed in Table 1, included CLINC-Full as the IND dataset paired with ROSTD and SST-2 as OOD datasets. These combinations were selected to rigorously

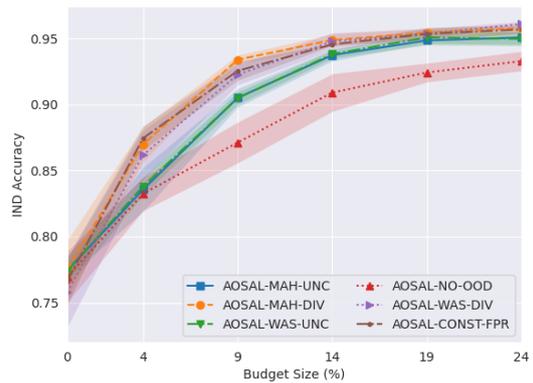


Figure 6: Ablation study on the IND test accuracy over the CLINC-Full (IND) and SST-2 (OOD) dataset using different AOSAL variants.

evaluate the robustness of our AOSAL approach across diverse scenarios. The results clearly indicate that the IND accuracy is sensitive to the FPR threshold set. For the CLINC-Full and ROSTD dataset configuration, the IND accuracy peaks at a

FPR (%)	IND ACC (%)	
	CLINC (IND) ROSTD (OOD)	CLINC (IND) SST-2 (OOD)
<b>90</b>	95.71 ( $\pm 0.00$ )	95.28 ( $\pm 0.01$ )
<b>95</b>	95.93 ( $\pm 0.01$ )	<b>96.14</b> ( $\pm 0.01$ )
<b>97</b>	<b>96.15</b> ( $\pm 0.00$ )	95.31 ( $\pm 0.00$ )

Table 1: IND accuracy at different FPR thresholds for CLINC-Full (IND), ROSTD (OOD), and SST-2 (OOD).

97% FPR setting, suggesting a balanced threshold that avoids excessive false positives while maintaining a high detection rate of in-domain samples. Conversely, tightening the FPR to 95% shows a slight dip in accuracy, which could imply an over-restriction misses some IND instances. A similar trend is observed in the CLINC-Full and SST-2 dataset configuration, reinforcing the importance of carefully calibrating the FPR threshold according to specific dataset characteristics and operational requirements. This analysis underscores the significance of the AOSAL’s adaptability to different operational scenarios. By systematically evaluating various FPR thresholds, we can identify the optimal setting that balances the trade-off between maintaining high in-domain accuracy and minimizing false positives.

## 6 Limitations

**Sensitivity to Hyperparameters.** One of the key challenges of our AOSAL framework is its dependence on hyperparameter settings. The choice of hyperparameters such as adaptive threshold for FPR and informative parameters is critical for achieving maximal learning efficiency. However, reaching this balance is by nature difficult since this directly affects the framework’s performance in correctly distinguishing OOD samples. Getting the wrong values for hyperparameters leads to either underconfidence or overconfidence in OOD instances and hence the model’s overall performance. Future works can be directed towards implementing more intelligent adaptive hyper-parameter tuning methods that are sensitive to changes in the data environment.

**Model Performance with Sparse Data.** Another critical limitation arises when there is a lack of data availability. With few input data points, our framework cannot generate and calibrate the right distance metrics for OOD detection. This can

hinder the accurate classification and enhancement of OOD detection, especially in the initial stages of training the model. There are potential ways to tackle these challenges, such as improving data augmentation methods and the use of synthetic data generation to help improve the model’s performance despite starting with minimal initial data.

**Scalability in Human-in-the-loop Setting.** While oracles enable AI models of TODS to train more efficiently with fewer samples via annotations, this process is not always scalable for annotators. This challenge in a human-in-the-loop setting is particularly evident when oracles provide a significant number of annotations for OOD samples within each AL round due to large unlabeled pools. Alternatively, previous works have created modeling approaches in other domains such as computer vision (Ning et al., 2022; Yang et al., 2024; Safaei et al., 2024) that train over both IND/OOD samples and AL sampling techniques for automatically extracting OOD samples in the unlabeled pool. Consequently, this effectively reduces the number of annotations the oracle provides.

## 7 Conclusion

In this paper, we presented AOSAL which is an AL framework that aims to improve the efficiency and effectiveness of TODS. AOSAL combines a distance-based OOD technique with an adaptive FPR threshold and an informativeness measure based on uncertainty and diversity. This integration enables AOSAL to improve the classification of IND and OOD samples and thus focuses primarily on the most useful IND examples from an unlabeled data pool for training. The experimental analysis we have conducted shows that AOSAL is highly effective for dealing with complex datasets in comparison to traditional active learning techniques. These real-world applications have demonstrated the practical usefulness and effectiveness of the framework in enhancing not only the robustness but also the accuracy of intent classification in TODS by the AOSAL framework’s ability to selectively acquire high-value IND training samples.

In future work, one can investigate advanced data augmentation and synthetic data approaches to facilitate training in data-deficient scenarios and design adaptive hyperparameter optimization of the system’s responsiveness to data variability.

## References

- Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. 2020. [Deep batch active learning by diverse, uncertain gradient lower bounds](#). In *Proceedings of ICLR*.
- Andleeb Aslam, Usman Qamar, Pakizah Saqib, Reda Ayesha, and Aiman Qadeer. 2020. A novel framework for sentiment analysis using deep learning. In *2020 22nd International Conference on Advanced Communication Technology (ICACT)*, pages 525–529. IEEE.
- Xun Deng, Wenjie Wang, Fuli Feng, Hanwang Zhang, Xiangnan He, and Yong Liao. 2023. Counterfactual active learning for out-of-distribution generalization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11362–11377.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of NAACL-HLT*.
- Pan Du, Suyun Zhao, Hui Chen, Shuwen Chai, Hong Chen, and Cuiping Li. 2021. [Contrastive coding for active learning under class distribution mismatch](#). In *Proceedings of IEEE/CVF*.
- Charlie Frogner, Chiyuan Zhang, Hossein Mobahi, Mauricio Araya, and Tomaso A Poggio. 2015. Learning with a wasserstein loss. In *Advances in neural information processing systems*, volume 28.
- Varun Gangal, Abhinav Arora, Arash Einolghozati, and Sonal Gupta. 2020. Likelihood ratios and generative classifiers for unsupervised out-of-domain detection in task oriented dialog. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7764–7771.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. [Deberta: Decoding-enhanced bert with disentangled attention](#). In *International Conference on Learning Representations*.
- Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. 2011. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*.
- Kuan-Hao Huang. 2021. [Deepal: Deep active learning in python](#). *arXiv preprint arXiv:2111.15258*.
- Ajay J Joshi, Fatih Porikli, and Nikolaos Papanikolopoulos. 2009. Multi-class active learning for image classification. In *2009 IEEE conference on computer vision and pattern recognition*, pages 2372–2379. IEEE.
- Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. 2019. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *Advances in neural information processing systems*, 32.
- Suraj Kothawade, Nathan Beck, Krishnateja Killamsetty, and Rishabh Iyer. 2021. Similar: Submodular information measures based active learning in realistic scenarios. *Advances in Neural Information Processing Systems*, 34:18685–18697.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*.
- Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. [An evaluation dataset for intent classification and out-of-scope prediction](#). In *Proceedings of EMNLP-IJCNLP*.
- David D Lewis. 1995. [A sequential algorithm for training text classifiers: Corrigendum and additional data](#). In *Proceedings of ACM SIGIR*.
- 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*.
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150.
- Katerina Margatina, Giorgos Vernikos, Loïc Barrault, and Nikolaos Aletras. 2021. [Active learning by acquiring contrastive examples](#). In *Proceedings of EMNLP*.
- Hieu T Nguyen and Arnold Smeulders. 2004. [Active learning using pre-clustering](#). In *Proceedings of ICML*.
- Kun-Peng Ning, Xun Zhao, Yu Li, and Sheng-Jun Huang. 2022. Active learning for open-set annotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 41–49.
- Dongmin Park, Yooju Shin, Jihwan Bang, Youngjun Lee, Hwanjun Song, and Jae-Gil Lee. 2022. Meta-query-net: Resolving purity-informativeness dilemma in open-set active learning. *Advances in Neural Information Processing Systems*, 35:31416–31429.
- Alexander Podolskiy, Dmitry Lipin, Andrey Bout, Ekaterina Artemova, and Irina Piontkovskaya. 2021. Revisiting mahalanobis distance for transformer-based out-of-domain detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 13675–13682.

- Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. 2021. A survey of deep active learning. *ACM computing surveys (CSUR)*, 54(9):1–40.
- Bardia Safaei, VS Vibashan, Celso M de Melo, and Vishal M Patel. 2024. Entropic open-set active learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4686–4694.
- Walter J Scheirer, Anderson de Rezende Rocha, Archana Sapkota, and Terrance E Boulton. 2012. Toward open set recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(7):1757–1772.
- Ozan Sener and Silvio Savarese. 2018. [Active learning for convolutional neural networks: A core-set approach](#). In *International Conference on Learning Representations*.
- Burr Settles. 2009. Active learning literature survey. *University of Wisconsin-Madison Department of Computer Sciences*.
- Changjian Shui, Fan Zhou, Christian Gagné, and Boyu Wang. 2020. [Deep active learning: Unified and principled method for query and training](#). In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, pages 1308–1318. PMLR.
- 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 EMNLP*.
- Yuxia Wu, Tianhao Dai, Zhedong Zheng, and Lizi Liao. 2024. Active discovering new slots for task-oriented conversation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- Yang Yang, Yuxuan Zhang, Xin Song, and Yi Xu. 2024. Not all out-of-distribution data are harmful to open-set active learning. *Advances in Neural Information Processing Systems*, 36.
- Changchang Yin, Buyue Qian, Shilei Cao, Xiaoyu Li, Jishang Wei, Qinghua Zheng, and Ian Davidson. 2017. Deep similarity-based batch mode active learning with exploration-exploitation. In *2017 IEEE International Conference on Data Mining (ICDM)*, pages 575–584. IEEE.
- Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-Graber. 2020. [Cold-start active learning through self-supervised language modeling](#). In *Proceedings of EMNLP*.
- Leihan Zhang and Le Zhang. 2019. An ensemble deep active learning method for intent classification. In *Proceedings of the 2019 3rd International Conference on Computer Science and Artificial Intelligence*, pages 107–111.
- Fedor Zhdanov. 2019. Diverse mini-batch active learning. *arXiv preprint arXiv:1901.05954*.

## A Appendix

### A.1 Generalization of Distance-based OOD Detection Method

As introduced in Section 3.2.1, the distance score function  $S(x_j)$  is designed to be adaptable to various distance metrics, accommodating different mathematical properties and score ranges. Specifically, for a  $\mathcal{K}$ -class classification problem, we maximize the selection of IND samples from  $\mathcal{D}_U$  by computing the minimum distance between an unlabeled sample  $x_j$  and each class in the labeled dataset  $\mathcal{D}_L$ . Herein, we demonstrate applicability of our generalized distance-based OOD detector to the Mahalanobis distance (Podolskiy et al., 2021) and Wasserstein distance (Frogner et al., 2015).

**Mahalanobis Distance.** We utilize the Mahalanobis distance has shown to be useful for classifying detecting OOD instances without the reliance accessing OOD instances for training (Podolskiy et al., 2021). This distance method is a way to determine the closeness of an data sample to a set of data samples that belongs to a class  $k$ .

Given an unlabeled sample  $x_j$  from  $\mathcal{D}_U$ , the Mahalanobis distance can be calculated as:

$$d(x_j) = \min_{k \in \mathcal{K}} (\Phi(x_j) - \mu_{x_k})^\top \Sigma^{-1} (\Phi(x_j) - \mu_{x_k}), \quad (7)$$

where  $\Phi(x_j)$  is the embedding of the unlabeled sample  $x_j$ ,  $\mu_{x_k}$  is the mean of the multivariate Gaussian distribution of class  $k \in \{1, \dots, \mathcal{K}\}$ , and  $\Sigma$  represents the covariance matrix. The calculations of  $\mu_k$  and  $\Sigma$  are computed as:

$$\mu_{x_k} = \frac{1}{N_k} \sum_k \Phi(x), \quad (8)$$

$$\Sigma = \frac{1}{N_L} \sum_k \sum_{i \in k} (\Phi(x_i) - \mu_{x_k})(\Phi(x_i) - \mu_{x_k})^\top, \quad (9)$$

where  $N_k$  is the number of training samples the class  $k$  and  $N_L$  is the total number of training samples in the labeled set. While the range of distances of the Mahalanobis distance is  $[0, \infty]$ , we transform the ranges Equation 7 to  $[0, 1]$  using Equation 3.

**Wasserstein Distance.** Similarly, the Wasserstein distance calculates the minimal cost of transporting

mass from the distribution of  $x_j$  to that of each class distribution  $k$  where the cost is defined by the ground distance between the distributions (Frogner et al., 2015). Given an unlabeled sample  $x_j$  from  $\mathcal{D}_U$ , the Wasserstein distance can be calculated as:

$$S(x_j) = \arg \min_{k \in \mathcal{K}} W(\Phi(x_j), \mu_{x_k}) \quad (10)$$

$$W(\Phi(x_j), \mu_k) = \inf_{\gamma \in \Gamma(P_{\Phi(x_j)}, P_{\mu_k})} \int \|\Phi(x_j) - \mu_k\|_2 d\gamma(\Phi(x_j), \mu_{x_k}) \quad (11)$$

Here,  $\Phi(x_j)$  is the feature vector of  $x_j$ ,  $P_{\Phi(x_j)}$  and  $P_{\mu_k}$  represent the empirical distributions of  $x_j$  and class  $k$ , respectively, and  $\Gamma(P_{\Phi(x_j)}, P_{\mu_k})$  contains all feasible joint distributions  $\gamma$  where the marginals are  $P_{\Phi(x_j)}$  and  $P_{\mu_k}$ . Wasserstein distances have a non-negative range  $[0, \infty]$ , where 0 represents perfect similarity between distributions. These distances can be normalized to the range  $[0, 1]$ , using a transformation similar to Equation 3.

### A.2 Dataset Details

In this section, we provide the dataset statistics of each NLP benchmark dataset shown in Table 2. In the following, we provide a brief description for each of the dataset as it related to intent classification.

**CLINC-Full.** The CLINC-Full dataset was introduced by (Larson et al., 2019) which is designed for intent classification across multiple domains such as banking, home, travel, and business. It contains a total of 23,700 queries, out of which 22,500 are in-distribution (IND) queries spanning 150 classes for intent classification tasks, and 1,200 are out-of-distribution (OOD) samples for out-of-scope prediction. This dataset is crucial for assessing the model’s capability to classify complex, real-world user intents and includes numerous OOD scenarios to evaluate robustness in model performance.

**ROSTD.** The Real Out-of-domain Sentences From Task-Oriented Dialog (ROSTD) dataset, proposed by (Gangal et al., 2020), is designed for training and evaluating intent classification models in task-oriented dialog systems with a focus on out-of-distribution robustness. It contains 34,059

Statistic	CLINC-Full	ROSTD	SST-2
Train	16950	25218	54577
Valid	2700	3537	6822
Test	4050	5304	6822
OOD samples	1200	4590	0
% of OOD samples in unlabeled pool	7.87%	20.22%	0%
IND classes	150	12	2

Table 2: Dataset statistics for CLINC-Full, ROSTD and SST-2.

queries across 12 classes, including in-distribution queries and an additional 4,590 out-of-distribution samples curated with human annotations. The dataset aims to facilitate the development of more robust dialog systems capable of handling out-of-distribution utterances effectively.

**SST-2.** The Stanford Sentiment Treebank (SST-2) (Aslam et al., 2020) is another well-established benchmark for sentiment analysis, particularly for tasks that involve considering sentence structure and sentiment polarity. It consists of 67,314 sentences for training, 855 for validation, and 1,821 for testing, all derived from movie review sentences on Rotten Tomatoes. Each sentence is labeled as positive, negative, or neutral.

### A.3 Model Implementation & Hyperparameters

In this section, we provide details of the model implementation and hyperparameters used in our experiments. We use a pre-trained BERT model (Devlin et al., 2019) from the HuggingFace library (Wolf et al., 2020) and integrated it in our Python environment using PyTorch 2.0 and PyTorch Lightning. We train BERT using a batch size of 32, learning rate of  $5e-5$ , AdamW optimizer epsilon  $1e-6$  and weight decay of 0.001, and embedding dimension of 768. For all datasets, we used a maximum sequence length of 256. We pre-trained the base model over 5 epochs on CLINC-Full and 1 epoch on ROSTD and SST-2. In the AL cycle, we use the newly acquired samples from  $\mathcal{D}_U$  to fine-tune BERT over the updated labeled set  $\mathcal{D}_L$ . We ensure fair comparison among each AL method by evaluating them 5 times using a different random seed. Each experiment is run on an Nvidia A100 80GB GPU. We use the open source materi-

als from (Huang, 2021; Ash et al., 2020; Margatina et al., 2021) to implement the baseline AL methods from their respective source code repository on GitHub.

# Dialogue Ontology Relation Extraction via Constrained Chain-of-Thought Decoding

Renato Vukovic, David Arps, Carel van Niekerk, Benjamin Matthias Ruppik,  
Hsien-Chin Lin, Michael Heck, Milica Gašić

Heinrich Heine University Düsseldorf

{renato.vukovic, david.arps, niekerk, ruppik, linh, heckmi, gasic}@hhu.de

## Abstract

State-of-the-art task-oriented dialogue systems typically rely on task-specific ontologies for fulfilling user queries. The majority of task-oriented dialogue data, such as customer service recordings, comes without ontology and annotation. Such ontologies are normally built manually, limiting the application of specialised systems. Dialogue ontology construction is an approach for automating that process and typically consists of two steps: term extraction and relation extraction. In this work, we focus on relation extraction in a transfer learning set-up. To improve the generalisation, we propose an extension to the decoding mechanism of large language models. We adapt Chain-of-Thought (CoT) decoding, recently developed for reasoning problems, to generative relation extraction. Here, we generate multiple branches in the decoding space and select the relations based on a confidence threshold. By constraining the decoding to ontology terms and relations, we aim to decrease the risk of hallucination. We conduct extensive experimentation on two widely used datasets and find improvements in performance on target ontology for source fine-tuned and one-shot prompted large language models.<sup>1</sup>

## 1 Introduction

State-of-the-art task-oriented dialogue (TOD) systems still rely on a fixed ontology to model their scope (Nguyen et al., 2023; Hudeček and Dusek, 2023). A *TOD ontology* comprises three levels of hierarchy: domains, slots and values. *Domains* are general topics of interest, *slots* are types of information about entities in a domain, and *values* are concrete instantiations of slots. Ontology thus forms a hierarchy: it is a directed graph where slots belong to domains and values in turn belong to slots. Note that slots can be shared across domains,

and so can values. An ontology is typically a prerequisite for generating API calls that access the underlying databases for entity retrieval. Further, the ontology defines the dialogue state, which is tracked by the system to determine the next actions given the evolving discourse.

The dependency on an ontology poses a significant challenge in transferring existing TOD systems to new domains and use cases. Although ontology-agnostic approaches do exist, their transfer capabilities are limited and their performance remains sub-par on novel data (Heck et al., 2022).

Large quantities of domain-specific TOD data, e.g. customer service recordings, are frequently available, but tend to come without annotation, rendering direct use for system development difficult (Brusco and Gravano, 2023). Manual labelling is error-prone, does not scale well and quickly becomes prohibitively expensive (Eric et al., 2020; Rosenbaum et al., 2022; Gung et al., 2023). Despite topical or domain mismatch, existing annotated datasets may provide information about TOD that can be leveraged to harness new data. For this reason, we are interested in utilising existing labelled TOD datasets to automatically generate a full ontology for new, yet-unlabelled, data.

Automatic dialogue *ontology construction* typically consists of two steps, dialogue term extraction (Vukovic et al., 2022) and hierarchy establishment. Although hierarchy establishment is often done via clustering (Hudeček et al., 2021; Yu et al., 2022) we approach it via *relation extraction* (RE), which is more similar to common information extraction pipelines (Genest et al., 2022; Xu et al., 2023). We call this task *dialogue ontology relation extraction* (DORE). A hierarchy is established by inferring in which level extracted terms lie, and by connecting terms across levels.

Although large language models (LLMs) have demonstrated considerable task transfer abilities (Brown et al., 2020; Ouyang et al., 2022),

<sup>1</sup>The code is available under <https://gitlab.cs.uni-duesseldorf.de/general/dsml/dialogue-ontology-relation-extraction-via-constrained-chain-of-thought-decoding>

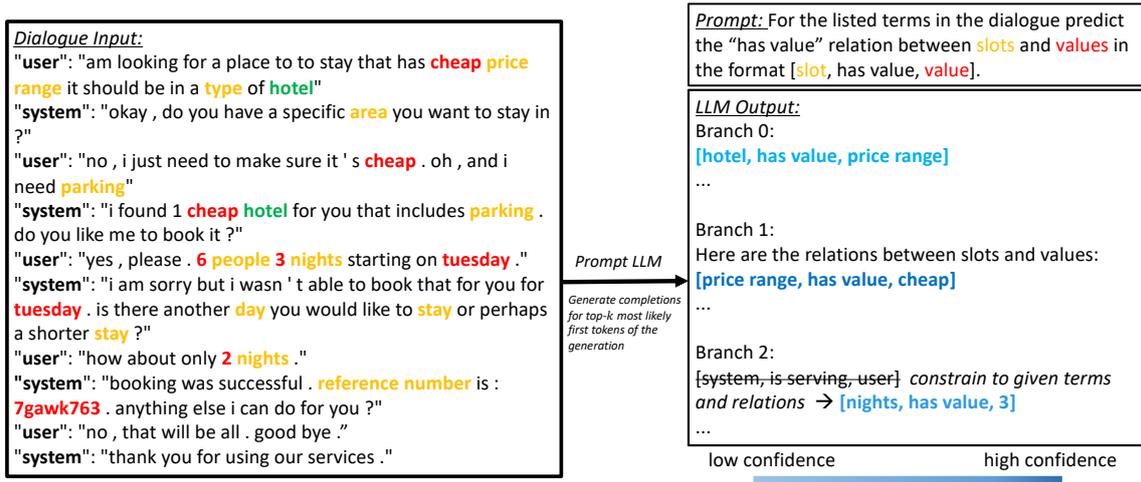


Figure 1: Example of constrained CoT-decoding for dialogue ontology extraction for a dialogue from MultiWOZ 2.1 (Eric et al., 2020). Domains are highlighted in green, slots in yellow and values in red. Branch 0 predicts an incorrect relation (*hotel* misclassified as slot) with lower confidence. Branch 1 has the highest confidence in the relation prediction, which is why it is selected as the final response. Also, it contains a form of reasoning that stresses the type of terms that are part of the relations to be predicted, i.e., slots and values. Branch 2 visualises constrained decoding, where the prediction of terms and relations is not possible if they are not present in the input.

they still lack behind specialised systems in TOD modelling when appropriate training data is available (Heck et al., 2023; Hudeček and Dusek, 2023).

In this work, we assume that some labelled out-of-domain source dialogue data is available to facilitate transfer learning. We examine two strategies of providing source data to an instruction-tuned LLM; 1) as one-shot examples in the prompt, and 2) as data for an additional round of supervised fine-tuning. We establish a challenging transfer setup by conducting experiments on two well-established medium to large scale multi-domain task-oriented dialogue benchmark datasets: MultiWOZ 2.1 (Budzianowski et al., 2018; Eric et al., 2020) and the Schema-Guided Dialogue (SGD; Rastogi et al., 2020) dataset. Since our focus is solely on DORE, we assume that the results of the first step of ontology construction, namely term extraction, are provided.

We propose to improve the decoding mechanism of an LLM in order to better leverage task-specific knowledge. Concretely, we constrain the generation to terms and relation types given in the model input to force the model to consider terms from the target data and output the desired format. We further adapt *chain-of-thought (CoT) decoding* (Wang and Zhou, 2024), which was recently proposed for logical reasoning, for DORE. Traditionally, CoT methods prompt or train the model to generate reasoning paths before giving the final answer (Wei

et al., 2022; Kim et al., 2023). CoT-decoding on the other hand exploits the observation that the presence of CoT-paths is correlated with higher confidence in the predicted answer in logical reasoning. We extend CoT-decoding to DORE by selecting the final answer based on the confidence of predicted relations in multiple generated model answer branches. Our final proposal, *constrained CoT-decoding for dialogue ontology extraction*, is the combination of our CoT-decoding approach to RE with constrained decoding, see Fig. 1. Empirically, this new decoding mechanism significantly outperforms both source one-shot and source fine-tuned baselines on the target data. Our contributions are as follows:

- We propose to induce an ontological hierarchy by accumulating ontology relation predictions from the dialogues in a TOD dataset.
- To the best of our knowledge, we are the first to apply *CoT-decoding* to dialogue ontology relation extraction.
- We develop an extension, called *constrained CoT-decoding*, for multi-relation extraction from task-oriented dialogues.
- Constrained CoT-decoding significantly improves the quality of relation predictions on the target dataset for both source one-shot and source fine-tuned baselines.

## 2 Related Work

**Dialogue Ontology Construction** We divide dialogue ontology construction into term extraction and relation extraction. Vukovic et al. (2022) improve out-of-domain generalisation of a dialogue term extraction model by making use of topological properties of the language model embedding space. Nguyen et al. (2023) improve phrasal segmentation of ontology terms via language model probing and contrastive learning. Since we evaluate the hierarchy on a global level based on relations, our approach is not directly comparable to clustering-based approaches such as Hudeček et al. (2021); Yu et al. (2022). In contrast to these methods, we view hierarchy establishment in isolation as a relation extraction task.

Yu et al. (2020) present DialogRE, a popular dataset for RE in short chit-chat dialogues. Closest to our approach, Albalak et al. (2022) jointly optimise RE and explanation generation to improve performance with a model-agnostic framework. Xu and Chen (2023) propose a zero-shot approach for extracting trigger words for dialogue relation extraction on DialogRE. However, these works focus on chit-chat dialogues, which do not include ontology relations.

**Relation Extraction with LLMs** LLMs show promising transfer capabilities out of the box (Laskar et al., 2023). Direct application to our task however is not promising, as it has been shown that aligned LLMs such as ChatGPT (OpenAI, 2022) do not perform well on extracting multiple relations at once (Lilong et al., 2024). This shortcoming has been linked to the influence of pre-training data distribution on downstream task performance (McCoy et al., 2023). RE data in particular amounts to a mere 0.5% of instruction-tuning datasets, and is hardly utilised for model selection (Wang et al., 2022; Zhang et al., 2023).

Traditionally, RE is performed in a pairwise manner (Zhang et al., 2023), resulting in quadratic complexity given the number of terms. This becomes intractable for generative LLMs when querying the LLM separately for each pair of terms. Alternatively, one may extract all relations present in a given input with a single LLM query, as is common in multi-relation extraction tasks such as document-level RE. For example, Lilong et al. (2024) extract relations by either predicting relations directly, or first predicting possible head entities in a docu-

ment. Zhang et al. (2023) align LLMs for zero-shot RE by transforming RE into a question answering (QA) task, which is more frequent in the instruction-tuning data.

**Constrained Decoding** Constrained decoding limits the tokens that can be generated. It is typically applied to LLMs to improve downstream task performance, reduce hallucination and ensure certain output formats. Bogoychev and Chen (2023) constrain decoding for translation to ensure that certain terminology is used. Roy et al. (2024) use constrained decoding with a lookahead heuristic to speed up adaptation of LLMs to plan generation according to a given API in TOD. We want to force the model to use its inherent task knowledge while transferring abilities to new data.

**Chain-of-Thought Reasoning** LLM performance on complex reasoning tasks improves when the model generates a chain of thought (CoT). Wei et al. (2022) include examples of multi-step reasoning in the prompt, and Kojima et al. (2022) prompt the model in a zero-shot fashion to “think step by step”. Reasoning capabilities can be further enhanced via specific training on CoT-data (Chung et al., 2024), or via teaching the model to reason (Zelikman et al., 2022). In contrast to this, we focus on eliciting model-inherent reasoning capabilities, without the need for specific prompts or training. As described in Sec. 3.2, we leverage the fact that a top- $k$  decoding beam usually contains a CoT (Wang and Zhou, 2024).

## 3 Constrained Chain-of-Thought Decoding for Ontology Relation Extraction

### 3.1 Problem Definition

Dialogue ontology relation extraction (DORE) aims at extracting all relations between different terms in a TOD dataset. As seen in Figure 1, for each dialogue paired with a list of ontology terms, the output is a set of relations similar to document-level relation extraction (Tan et al., 2022). However, we consider the joint relation prediction set accumulated from all dialogue-level predictions, rather than the dialogue-level performance. In the DORE task, the model receives as input a task-oriented dialogue  $D$  annotated with a list of ontology terms  $T$  present in this dialogue. The output are valid ontology relations  $R_{D,T}$  between the terms, which includes predicting whether a term

Relation	Verbaliser	Example
Domain-Slot	[Domain, has slot, Slot]	[hotel, has slot, price range]
Slot-Value	[Slot, has value, Value]	[price range, has value, cheap]
Value-Domain	[Value, has domain, Domain]	[cheap, has domain, hotel]
Equivalence	[Term1, refers to same concept as, Term2]	[cheap, refers to same concept as, low budget]

Table 1: Hierarchical dialogue ontology relation task definition with examples.

is a *domain*, *slot*, or *value*. A relation is denoted by a relational triplet with a head term, the relation and a tail term. Finally, the predicted relations for each dialogue are unified to form the final ontology relation set.

We consider 4 types of relation between ontology terms: *domain-slot*, *slot-value*, *value-domain* and *equivalent term* relations (see Table 1 for examples). Here, all relations except the equivalence relation are directed relations with a head and a tail term. *Domains* are general topics, such as *hotel* or *restaurant*, *slots* are types of information for entities in a *domain*, such as *price range* or *area* and *values* are concrete instantiations of slots, such as “cheap” or “west”. The equivalence relation connects terms from the same hierarchy level that point to the same ontological concept, e.g. “expensive” and “high-end” both represent a high price. In the prompt and labels, we denote the relation types through different verbalisers, shown in Table 1. Verbalisers are descriptions of task-specific labels in natural language. They align the task closer with the pre-training distribution of the LLM (Schick and Schütze, 2021; Mosbach et al., 2023).

Our hypothesis is that the general definitions of the ontology hierarchy relations enable seamless transfer to new data in order to construct a similarly structured ontology on the new data. Based on these relations, we focus on transferring the structural information about ontologies from a source dataset to a target dataset. Here, we consider a one-shot and a fine-tuning approach.

### 3.2 Chain-of-Thought Decoding

CoT reasoning in LLMs has demonstrated improved performance in various complex reasoning tasks (Sec. 2). The results of Wang and Zhou (2024) show that LLMs inherently possess reasoning capabilities, which can be elicited without explicit prompting through *Chain-of-Thought decoding*. Concretely, they experiment on pre-trained and instruction-tuned versions of PaLM 2 (Anil et al., 2023) and Mistral-7B (Jiang et al., 2023). They observe that although the greedily decoded

response might not always exhibit reasoning, one of the top- $k$  beams usually contains a CoT. This CoT not only shows higher confidence in the answer, but also exhibits greater accuracy. They propose to consider the top- $k$  probability tokens at the start of the predicted response. From there,  $k$  completions, called *branches*, are generated, resulting in  $k$ -times computational complexity during inference. The final response is chosen based on the confidence of the tokens that belong to the *answer* in each branch, i.e., the average confidence of the *answer tokens*. In logical reasoning, there is only one answer in each branch, which is a number. In that case, they identify the answer by prompting the model with “So the answer is:” at the end and match the following number to one in the preceding response. In our case, there are multiple answers per branch, which we identify based on the fact that relations are supposed to be predicted between brackets.

**CoT-Decoding for DORE** In this paper, we extend CoT decoding to handle the multi-answer scenario in the DORE task. We compute the confidence of answer tokens by utilising their structure, which, in our case, involves predicting relational triplets in the format  $[headterm, relation, tailterm]$  and the notion of disparity. The disparity of a probability distribution is the difference between the probability of the most likely outcome and the next most likely outcome. The confidence for each answer token for a given branch is measured by the average disparity of its tokens. Formally this is given by

$$\Delta_{i,a} = \frac{1}{n} \sum_{x_t \in a} p(x_t^{\text{top}} | x_{<t}) - p(x_t^{\text{next}} | x_{<t}), \quad (1)$$

where  $a$  is an answer (in our case the triplet),  $i$  is a branch,  $x_t$  are the answer tokens belonging to the answer in branch  $i$ ,  $x_t^{\text{top}}$  is the most likely token on position  $t$  and  $x_t^{\text{next}}$  the next most likely token on position  $t$ .  $x_{<t}$  are the tokens in branch  $i$  on positions preceding  $t$ , i.e. the context so far.

In DORE, answer tokens are those that form terms and relations in the predicted relational

triplets, which means there are three disparities per relation. This approach relies on detecting answer tokens in a generated response for confidence estimation, and we leave an extension to arbitrary answer structures to future work. The resulting triplet disparities are denoted as  $\Delta_a = [\Delta_h, \Delta_r, \Delta_t]$ . We explored mean, median, maximum, and minimum as aggregation strategies for relational triplet mentions, finding that all of them lead to similar results. For simplicity, we choose the mean to aggregate the disparity for a relational triplet in branch  $i$ , i.e.  $\Delta_{i,a} = \frac{1}{3}(\Delta_{h,i} + \Delta_{r,i} + \Delta_{t,i})$ .

We select the branch with the highest average disparity over the relations predicted in each branch to get the final set of relation predictions for a dialogue. The average disparity for branch  $i$  is given by

$$\bar{\Delta}_i = \frac{1}{n_{a,i}} \sum_{a \in R_i} \Delta_{i,a} \quad (2)$$

where  $a$  is a relational triplet,  $R_i$  is the set of relations and  $n_{a,i}$  is the number of relations in branch  $i$ . The final set of predicted relations is then given by

$$R_{\bar{\Delta}_{\max}} = \{R_i \mid i = \operatorname{argmax}\{\bar{\Delta}_0, \dots, \bar{\Delta}_k\}\} \quad (3)$$

We also experiment with a confidence threshold based approach for relation selection. Here, the average disparity of a relation is computed across occurrences in different branches:

$$\tilde{\Delta}_a = \frac{1}{n_a} \sum_{i \in \{1, \dots, k\}} \Delta_{i,a}, \quad (4)$$

where  $\Delta_{i,a}$  is the disparity of the answer  $a$  in the  $i$ -th branch and  $n_a$  is the number of occurrences of  $a$  across the different branches. The final set of predicted relations  $R_{\tilde{\Delta} > \Delta_{\text{threshold}}}$  is then

$$R_{\tilde{\Delta} > \Delta_{\text{threshold}}} = \{a \mid \tilde{\Delta}_a > \Delta_{\text{threshold}}\} \quad (5)$$

### 3.3 Constrained Decoding

We constrain the generation of the relation terms and relation types if the beginning of a relational triplet is predicted to ensure the structure and mitigate term and relation hallucination (see Figure 1). This means for a relational triplet,  $[h, r, t]$ , we ensure that  $h, t \in T$  and  $r \in R$ , where  $T$  is the set of terms for the current dialogue and  $R$  is the set of relation types given in the prompt. Note that we only constrain the generation when an opening bracket is predicted by the model, and resume to non-constrained generation after the generated relational triplet.

## 4 Experiments

### 4.1 Experimental Setup

We utilise the open-source Gemma 2B (Mesnard et al., 2024) instruction-tuned model with context size of 4096 for all experiments. In CoT-decoding we set  $k = 5$ . For a more thorough analysis of the impact of  $k$  in CoT decoding, resort to Wang and Zhou (2024). We always branch at the first token; branching at later tokens did not show improvements. For all CoT-decoding experiments, we select the relations from the branch with the highest disparity, as the threshold based method works worse and also adds a new hyperparameter. In the one-shot prompts, we use a combination of an instruction with simple natural language with a preceding example (Brown et al., 2020; Sahoo et al., 2024). For fine-tuning, we remove the example from the prompt.

**Datasets** For the source dataset, we employ the MultiWOZ 2.1 dataset (Eric et al., 2020). It has 7 domains and over 10,000 dialogues. We use the training set for training and select from it one random dialogue with relation annotation as one-shot exemplar. The target dataset is the schema-guided dialogue (SGD) dataset (Rastogi et al., 2020). It comprises more than 20,000 dialogues and 20 domains. We use the SGD test split for evaluation in the main results, which contains 4,201 dialogues and 18 domains. In the test set, there are 134 domain-slot relations, 6,162 slot-value relations, 8,233 value-domain relations and 330 equivalence relations. It is worth noting the SGD test set contains dialogues from different domains than the SGD training set, as well as a significant amount of unseen ontology relations. We use ConvLab-3 (Zhu et al., 2023b) for loading all the datasets.

**Training** For both fine-tuning and one-shot prompting, we utilise the original Gemma prompt template (Mesnard et al., 2024). For training, we utilise Low-rank adaptation (LoRA, Hu et al., 2022) with the default parameters in the peft library (Man-grulkar et al., 2022). We train the model on a single NVIDIA RTX8000 GPU and do inference with all models on one NVIDIA RTX6000 GPU.

We only consider a one-shot approach due to context size constraints, as the relational triplets in the exemplars contain brackets. Brackets are considered individual tokens, increasing the number of tokens significantly. Because of this a maximum of three exemplars fits in the context size,

Approach	F1-Score	Precision	Recall
<i>One-shot example from MultiWOZ</i>			
<i>Baseline</i> : Separate relation prediction	7.4	8.8	6.4
+ constrained decoding	<b>8.5*</b>	5.7	<b>17.3*†</b>
+ CoT decoding	<b>9.2*</b>	8.8	<b>9.6*</b>
+ constrained CoT decoding	<b>9.2*</b>	6.4	<b>15.9*</b>
<i>Fine-tuning on MultiWOZ</i>			
<i>Baseline</i> : Fine-tuning on MultiWOZ	10.9	6.8	28.8
+ constrained decoding	<b>12.0†</b>	7.4	32.3†
+ CoT decoding	10.6	7.6	17.6
+ constrained CoT decoding	<b>13.7†‡</b>	<b>9.8†</b>	23.0 ‡
<i>Upper Bounds using SGD Data</i>			
One-shot example from SGD + separate relation prediction	12.9	10.7	16.4
Fine-tuning on SGD	37.3	27.9	57.2

Table 2: Ontology Relation Prediction Results on the SGD test set. Results that are statistically significant over the baseline are highlighted in **bold**. Additionally, significant results based on dialogue-level evaluation for one-shot prompts are marked with \*. Significant results for fine-tuned models, evaluated globally based on five random seeds, are marked with †. Significant improvement over the one-shot model from the SGD upper bound on dialogue-level is marked with ‡. All significance tests are performed at a 5% level of significance.

which do not improve performance however, while increasing computational complexity. In the one-shot approach, we predict each relation type separately, since we found that the LLM struggles with jointly predicting all relation types. We also experimented with a zero-shot approach that performs significantly worse than one-shot.

We fine-tune the LLM via pattern-based fine-tuning (Schick and Schütze, 2021; Ma et al., 2023) with a prompt for all relation types on the MultiWOZ training split. We consider two upper bounds: an LLM trained on the SGD training split and a model utilising a one-shot exemplar from SGD.

## 4.2 Evaluation

In evaluation, we only consider relations within dialogues in the ground truth, i.e., both terms of a relation occur in the same dialogue. Relations from equivalent terms to other terms have to be found at least once. If  $[term_1, \textit{refers to same concept as}, term_2] \in R_{\textit{groundtruth}}$ , then  $[term_1, r, t] = [term_2, r, t]$ , where  $R_{\textit{groundtruth}}$  is the set of ground truth relations,  $r \neq \textit{refers to same concept as}$  is another relation type and  $t \in T$  is a third related term. E.g., the relations [price range, has value, high-end] and [price range, has value, expensive] are equivalent, since [expensive, refers to the same concept as, high-end]. Thus, the prediction of the former relation counts as a prediction for the latter and vice versa.

To compute the global micro F1 score, we compare the accumulated set of relations predicted from all the dialogues with the ground truth ontology relations. Note that we only consider exactly matching terms in relations to be correct.

For significance tests on the one-shot prompted models, we employ a pairwise  $t$ -test on dialogue level. For fine-tuned models, we use 5 random seeds for training and an independent  $t$ -test.

## 4.3 Results

Table 2 shows the full results on the target test set, see Appendix A for results for each relation type.

**Source One-Shot Approach** We found that when predicting all relations at once in a one-shot fashion the model is completely unable to fulfil the task, so we resort to predicting one relation at a time. The one-shot approach is mainly improved through constrained decoding, although the combination with CoT-decoding is also significantly better than the baseline. Note that the source one-shot model is able to get closer to the performance of a model with a one-shot example from the target data with constrained CoT-decoding.

**Source Fine-tuning Approach** For the source fine-tuned model, *constrained CoT-decoding* leads to significant improvements over the baseline. Furthermore, it significantly outperforms a model using a one-shot exemplar from the target data on all metrics. Constraining CoT-decoding helps per-

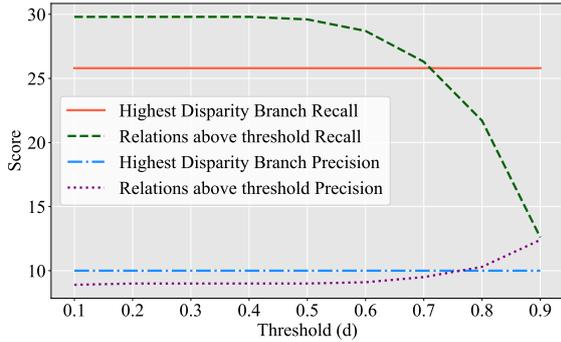


Figure 2: Different relation confidence thresholds across branches compared to the *highest disparity branch approach* for mean aggregation. Displayed are recall and precision for the MWOZ fine-tuned constrained CoT decoding model.

formance, since the constraints mitigate overconfidence on the source data after fine-tuning.

Interestingly, although the target fine-tuned model is the best model, it is not able to find all relations on the test set. As mentioned in Section 4.1, the SGD test set contains domains different to the SGD training set, which makes this task particularly difficult. In contrast to the excellent performance of LLMs on a variety of tasks, there is a lot of room for improvement on this task.

#### 4.4 Calibration Analysis

In Figure 2, we see that an absolute confidence threshold is not as meaningful and adds the problem of choosing the correct threshold as hyperparameter. Moreover, a high threshold leads to only a small increase in precision, while losing a significant amount of recall. Our results are in line with recent findings about instruction-tuned LLMs (Kapoor et al., 2024) being overconfident. We find that the model’s confidence on predicted relations is generally on a high level, indicating overconfidence, as the significant changes in performance happen at high confidence thresholds. For lower thresholds, the performance remains unchanged, as most confidences are quite high and hence the set of predicted relations stays the same. Although this shows that the thresholds are less meaningful, the relative confidence of the branches is meaningful, since choosing the highest disparity branch leads to good performance.

#### 4.5 How useful are predictions from the additional branches?

In line with the findings from Wang and Zhou (2024), we find that for the instruction-tuned

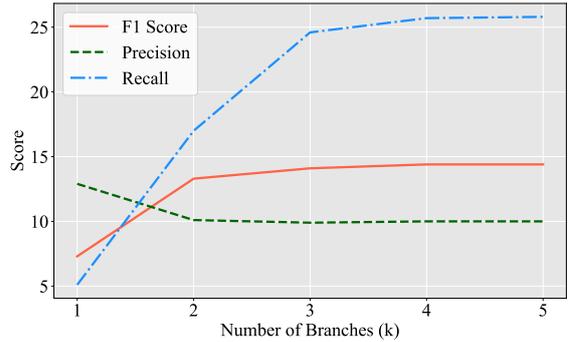


Figure 3: MWOZ fine-tuned constrained CoT decoding model with different choices for the number of branches  $k$  for the *highest disparity branch method*. Shown are SGD test set F1 score, precision, and recall.

Gemma model, the gain in performance can be mainly attributed to the first additional branch (see Figure 3). While the F1 score is increased slightly up to  $k = 3$ , the jump in recall from  $k = 2$  to  $k = 3$  is more significant. This also shows that the branches from lower ranked first tokens lead to responses with higher total confidence across the relation predicted in the respective branch, which is why they are chosen in the highest disparity branch selection method.

#### 4.6 Qualitative Analysis

Wang and Zhou (2024) found that LLMs struggle to generate CoTs for less frequent tasks in the pre-training data. In our analysis, we found that higher confidence completions often follow a recap of the type of terms and relations that should be predicted. Illustrated in Figure 4 is an example of a response to the one-shot equivalence prediction prompt with constrained CoT-decoding (see Appendix B for completions of the other decoding approaches). Here, branches 0 and 1 contain a repetition of the information given in the prompt. In branch 1 however, the focus on the equivalence mentioned in the prompt is followed by a response that does not use the proper format for the answer to be parsed correctly. The last branch has the highest confidence and is chosen ultimately. Here, the focus on the task relation and the provided dialogue is part of the generated introduction to the response.

As seen in Table 3, for the fine-tuned models, there are no such reasonings observed, although the constrained CoT decoding significantly improves performance. We hypothesise that the additional branches together with the constraints in decoding force the model to use task knowledge from

```

PROMPT: <EXAMPLE> Predict the relations between equivalent
terms in the form [term1, refers to same concept as, term2] as in
the example given above.

RESPONSE:
BRANCH 0:
**Domain:** Hotel
**Slot:** Type of hotel
**Value:** Price range
- [price per night, refers to same concept as, 135 carlingview drive
etobicoke]
....
BRANCH 1:
## Synonyms/Semantically Equivalent Terms:
**Domain:** Hotel
- best western plus toronto airport carlingview hotel
- ....
BRANCH 2:
1. **price range** [1, refers to same concept as, price]...
BRANCH 3:
• **price range** [1, refers to same concept as, price]
• ...
BRANCH 4 (highest confidence):
Here are the predicted synonyms or semantically equivalent terms
from the provided dialogue:
**Price range:**
- [price per night, refers to same concept as, cost]
....
**Day:**
- [tuesday next week, refers to same concept as, check in date]

```

Figure 4: Example of constrained CoT decoding for one-shot equivalence relation prediction. Branch 4 adds focus on the relation type. It has the highest average confidence in the predicted relations and is chosen as the final response. Some response parts are left out for illustration purposes (“...”).

fine-tuning, rather than what it has learned about the source data distribution. This can be observed when comparing CoT-decoding with constrained CoT-decoding, where the unconstrained version mainly generates terms it has seen on MultiWOZ, such as the “reference number” slot that is not present in SGD. The constrained version on the other hand forces the model to use task knowledge instead of distributional knowledge, leading to a much better coverage of the terms mentioned in the dataset, if the correct branch is chosen based on confidence. When observing completions to other dialogues, we found that the qualitatively best branches are not necessarily those with the highest confidence, indicating that a more sophisticated branch selection strategy might boost performance further. We leave such an improvement to future work. When comparing constrained decoding with vanilla greedy decoding, it becomes apparent that constraining the generation greatly improves the output structure and the utilisation of mentioned terms in the target dataset.

## 5 Discussion

Although the performance of the fine-tuned model is improved by constrained CoT-decoding, it is not clear where the improvement comes from based on qualitative analysis alone, as this model generates no reasoning paths. The workings and explainability of constrained CoT-decoding for fine-tuned models have not been investigated yet, but are relevant future research directions.

Our results imply that for tasks that are not frequently used in instruction-tuning data, it can be useful to utilise existing TOD data for training an LLM rather than annotating a few examples on the target data for the DORE task (see Section 4.3). Due to the length of examples in the DORE task, the amount of annotated examples that can fit in the prompt is highly limited, aggravating the applicability of few-shot approaches.

The results strengthen the finding that annotated data should be used if available (Zhu et al., 2023a). Constrained CoT-decoding only improves performance on the target data, where task knowledge is more useful than distributional knowledge about the source data. Mahowald et al. (2024) state that for learning tasks where language is used in a functional way to accomplish certain goals, such as DORE, self-supervised next-token prediction is not sufficient. Instead, the model needs to learn task-specific information via specialised fine-tuning to solve such tasks, which are not frequently present in pre-training data and involve task-specific reasoning. DORE can only be solved by knowing the specific relationship definition provided in the task descriptions, which the model cannot handle if it was not trained on the task. In summary, the presented results reinforce the observation that fine-tuning and specialised architectures are still needed to solve functional language-based tasks that cannot be solved by pattern matching alone.

The computational complexity of CoT-decoding is  $k$  times higher than regular greedy decoding, however, decoding of the different branches can be done in parallel. Nonetheless, one should note that only one run of ontology construction is needed, as the ontology can be repeatedly used for other dialogue modelling tasks such as state tracking. It is worth stressing that CoT-decoding only increases inference cost, while training cost remains the same. Compared to few-shot prompting, which also increases inference cost due to longer input context, there is no additional annotation cost.

Approach	Response
Baseline	- [price per night, has domain], [best western plus toronto north york hotel & suites, has domain], '\$ 63', 'hotel', 'has slot', 'has domain', 'hotel' ] nessunize "\$ 63" \$ 117' - [best western plus toronto north york hotel & suites, has domain], 'has slot', 'hotel', 'has domain', 'attraction' ] - [hotel, has slot], 'refers to same concept as', 'stay' ] - [hotel, has slot], 'has domain', 'hotel' ] ...
Constrained Decoding	- [price per night,has value,\$ 117] - [hotel,has slot,price per night] - [best western plus toronto north york hotel, has domain, hotel] - [hotel,has slot,address] - [best western plus toronto north york hotel, has domain, restaurant]... [hotel,has slot,date] ...
CoT-Decoding	[1, has domain, hotel] [best western plus toronto north york hotel & suites, has domain, hotel] [1, has domain, address] [best western premier toronto airport carlingview hotel, has domain, hotel] [hotel, has slot, day] [best western plus toronto north york hotel, has domain, hotel] [hotel, has slot, name] [hotel, has slot, room] [hotel, has slot, area] [hotel, has slot, reference number] ...
Constrained CoT-Decoding	[1, has domain, hotel] ... [hotel,has slot,price per night] [best western plus toronto north york hotel, has domain, restaurant] [best western plus toronto north york hotel, has domain, hotel] [hotel,has slot,address] [hotel,has slot,check in date] [hotel,has slot,street address] [hotel,has slot,date] [hotel,has slot,tuesday next week] [hotel,has slot,check in date] ...

Table 3: MWOZ fine-tuned model example response excerpts for the different decodings on SDG test dialogue 100. For CoT-decoding methods, only the chosen branch is displayed.

## 6 Conclusion

We propose constrained chain-of-thought (CoT)-decoding, a new decoding mechanism for dialogue ontology generation (DORE) in a transfer set-up. An LLM using a one-shot example from the source data is significantly improved using the proposed constrained CoT-decoding mechanism. Fine-tuning an LLM on the source data and using constrained CoT-decoding for inference on the target data outperforms a one-shot target data model significantly.

The results warrant further research into DORE in particular, and into eliciting reasoning in LLMs by adapting the decoding mechanism in general. Moreover, we offer a method for applying LLMs to tasks that are underrepresented in pre-training and where the vanilla LLMs perform poorly. Our method is appealing as it does not necessitate labelling new examples. Future research directions include explainability of constrained CoT-decoding in fine-tuned LLMs and including CoT-decoding during fine-tuning.

## 7 Limitations

In this work we assume a pipeline approach, however with the raise of LLMs, end-to-end solutions tend to be more accurate. We leave the task of jointly extracting dialogue terms and relations for future investigation. Due to constraints in computational infrastructure, we were not able to run open-source LLMs with the size of ChatGPT, which might be promising however. We abstained from utilising proprietary models, such as ChatGPT, for

increased transparency and reduced risk of training data contamination.

Furthermore, the need for an annotated source dataset limits the application to low-resource languages and tasks. The reliance on a specific answer structure for confidence estimation limits application to less structured tasks.

Finally, what we consider the upper bound, which was trained on the target dataset, can be argued to be a low bar too, reaching only an F1 of 37. This warrants more research on this task also on the same data setting.

## 8 Acknowledgements

RV and BMR are supported by funds from the European Research Council (ERC) provided under the Horizon 2020 research and innovation programme (Grant agreement No. STG2018 804636) as part of the DYMO project. CVN and HL are supported by the Ministry of Culture and Science of North Rhine-Westphalia within the framework of the Lamarr Fellow Network. MH is supported by funding provided by the Alexander von Humboldt Foundation in the framework of the Sofja Kovalevskaja Award endowed by the Federal Ministry of Education and Research. Computational infrastructure and support were provided by the Centre for Information and Media Technology at Heinrich Heine University Düsseldorf and Google Cloud. We want to thank the anonymous reviewers whose comments improved the quality of our paper.

## References

- Alon Albalak, Varun Embar, Yi-Lin Tuan, Lise Getoor, and William Yang Wang. 2022. [D-REX: Dialogue relation extraction with explanations](#). In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 34–46, Dublin, Ireland. Association for Computational Linguistics.
- 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](#). *Preprint*, arXiv:2305.10403.
- 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.
- 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.
- Pablo Brusco and Agustín Gravano. 2023. [Automatic offline annotation of turn-taking transitions in task-oriented dialogue](#). *Computer Speech & Language*, 78:101462.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. 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, et al. 2024. [Scaling instruction-finetuned language models](#). *Journal of Machine Learning Research*, 25(70):1–53.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. [MultiWOZ 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking base-lines](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 422–428, Marseille, France. European Language Resources Association.
- Pierre-Yves Genest, Pierre-Edouard Portier, Elöd Egyed-Zsigmond, and Laurent-Walter Goix. 2022. [PromptORE - A Novel Approach Towards Fully Un-supervised Relation Extraction](#). In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, CIKM '22*, page 561–571, New York, NY, USA. Association for Computing Machinery.
- James Gung, Raphael Shu, Emily Moeng, Wesley Rose, Salvatore Romeo, Arshit Gupta, Yassine Benajiba, Saab Mansour, and Yi Zhang. 2023. [Intent induction from conversations for task-oriented dialogue track at DSTC 11](#). In *Proceedings of The Eleventh Dialog System Technology Challenge*, pages 242–259, Prague, Czech Republic. Association for Computational Linguistics.
- Michael Heck, Nurul Lubis, Benjamin Ruppik, Renato Vukovic, Shutong Feng, Christian Geisshauer, Hsien-chin Lin, Carel van Niekerk, and Milica Gašić. 2023. [ChatGPT for zero-shot dialogue state tracking: A solution or an opportunity?](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 936–950, Toronto, Canada. Association for Computational Linguistics.

- Michael Heck, Nurul Lubis, Carel van Niekerk, Shutong Feng, Christian Geishausser, Hsien-Chin Lin, and Milica Gašić. 2022. **Robust dialogue state tracking with weak supervision and sparse data**. *Transactions of the Association for Computational Linguistics*, 10:1175–1192.
- 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*.
- Vojtěch Hudeček and Ondřej Dušek. 2023. **Are large language models all you need for task-oriented dialogue?** In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 216–228, Prague, Czechia. Association for Computational Linguistics.
- Vojtěch Hudeček, Ondřej Dušek, and Zhou Yu. 2021. **Discovering dialogue slots with weak supervision**. 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 2430–2442, Online. 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, Léo 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.
- Sanyam Kapoor, Nate Gruver, Manley Roberts, Arka Pal, Samuel Dooley, Micah Goldblum, and Andrew Wilson. 2024. **Calibration-tuning: Teaching large language models to know what they don’t know**. In *Proceedings of the 1st Workshop on Uncertainty-Aware NLP (UncertainLP 2024)*, pages 1–14, St Julians, Malta. Association for Computational Linguistics.
- Seungone Kim, Se Joo, Doyoung Kim, Joel Jang, Seonghyeon Ye, Jamin Shin, and Minjoon Seo. 2023. **The CoT collection: Improving zero-shot and few-shot learning of language models via chain-of-thought fine-tuning**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12685–12708, Singapore. Association for Computational Linguistics.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. **Large Language Models are Zero-Shot Reasoners**. In *Advances in Neural Information Processing Systems*.
- Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Huang. 2023. **A systematic study and comprehensive evaluation of ChatGPT on benchmark datasets**. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 431–469, Toronto, Canada. Association for Computational Linguistics.
- Xue Lilong, Zhang Dan, Dong Yuxiao, and Tang Jie. 2024. **AutoRE: Document-Level Relation Extraction with Large Language Models**. *Preprint*, arXiv:2403.14888.
- Bolei Ma, Ercong Nie, Helmut Schmid, and Hinrich Schuetze. 2023. **Is Prompt-Based Finetuning Always Better than Vanilla Finetuning? Insights from Cross-Lingual Language Understanding**. In *Proceedings of the 19th Conference on Natural Language Processing (KONVENS 2023)*, pages 1–16, Ingolstadt, Germany. Association for Computational Linguistics.
- Kyle Mahowald, Anna A. Ivanova, Idan A. Blank, Nancy Kanwisher, Joshua B. Tenenbaum, and Evelina Fedorenko. 2024. **Dissociating Language and Thought in Large Language Models**. *Trends in Cognitive Sciences*, 28(6):517–540.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. **PEFT: State-of-the-art Parameter-Efficient Fine-Tuning methods**. <https://github.com/huggingface/peft>.
- R. Thomas McCoy, Shunyu Yao, Dan Friedman, Matthew Hardy, and Thomas L. Griffiths. 2023. **Members of Autoregression: Understanding Large Language Models Through the Problem They are Trained to Solve**. *Preprint*, arXiv/2309.13638.
- 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 Crepey, 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 Mikula, 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 Klimentko, 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*.
- Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. 2023. [Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12284–12314, Toronto, Canada. Association for Computational Linguistics.
- Hoang Nguyen, Chenwei Zhang, Ye Liu, and Philip Yu. 2023. [Slot induction via pre-trained language model probing and multi-level contrastive learning](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 470–481, Prague, Czechia. Association for Computational Linguistics.
- OpenAI. 2022. [ChatGPT: Optimizing language models for dialogue](#). Accessed 2024-05-23.
- 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.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8689–8696.
- Andy Rosenbaum, Saleh Soltan, Wael Hamza, Yannick Versley, and Markus Boese. 2022. [LINGUIST: Language model instruction tuning to generate annotated utterances for intent classification and slot tagging](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 218–241, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Shamik Roy, Sailik Sengupta, Daniele Bonadiman, Saab Mansour, and Arshit Gupta. 2024. [FLAP: Flow-adhering planning with constrained decoding in LLMs](#). In *NAACL 2024*.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications. *Preprint, arXiv:2402.07927*.
- Timo Schick and Hinrich Schütze. 2021. [Exploiting cloze-questions for few-shot text classification and natural language inference](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 255–269, Online. Association for Computational Linguistics.
- Qingyu Tan, Lu Xu, Lidong Bing, Hwee Tou Ng, and Sharifah Mahani Aljunied. 2022. [Revisiting DocRED - addressing the false negative problem in relation extraction](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8472–8487, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Renato Vukovic, Michael Heck, Benjamin Ruppik, Carel van Niekerk, Marcus Zibrowius, and Milica Gašić. 2022. [Dialogue term extraction using transfer learning and topological data analysis](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 564–581, Edinburgh, UK. Association for Computational Linguistics.
- Xuezhi Wang and Denny Zhou. 2024. [Chain-of-Thought Reasoning Without Prompting](#). *Preprint, arXiv:2402.10200*.
- 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, Krma 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, Sujay 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, 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 *Advances in Neural Information Processing Systems*.
- Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong Xu, Xiangyu Zhao, Xian Wu, Yefeng Zheng, and Enhong Chen. 2023. Large language models for generative information extraction: A survey. *Preprint, arXiv:2312.17617*.
- Ze-Song Xu and Yun-Nung Chen. 2023. [Zero-shot dialogue relation extraction by relating explainable triggers and relation names](#). In *Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)*, pages 123–128, Toronto, Canada. Association for Computational Linguistics.

Dian Yu, Kai Sun, Claire Cardie, and Dong Yu. 2020. [Dialogue-based relation extraction](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4927–4940, Online. Association for Computational Linguistics.

Dian Yu, Mingqiu Wang, Yuan Cao, Izhak Shafran, Laurent Shafey, and Hagen Soltau. 2022. [Unsupervised slot schema induction for task-oriented dialog](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1174–1193, Seattle, United States. Association for Computational Linguistics.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. [STaR: Bootstrapping Reasoning With Reasoning](#). In *Advances in Neural Information Processing Systems*.

Kai Zhang, Bernal Jimenez Gutierrez, and Yu Su. 2023. [Aligning instruction tasks unlocks large language models as zero-shot relation extractors](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 794–812, Toronto, Canada. Association for Computational Linguistics.

Dawei Zhu, Xiaoyu Shen, Marius Mosbach, Andreas Stephan, and Dietrich Klakow. 2023a. [Weaker than you think: A critical look at weakly supervised learning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14229–14253, Toronto, Canada. Association for Computational Linguistics.

Qi Zhu, Christian Geishauer, Hsien-chin Lin, Carel van Niekerk, Baolin Peng, Zheng Zhang, Shutong Feng, Michael Heck, Nurul Lubis, Dazhen Wan, Xiaochen Zhu, Jianfeng Gao, Milica Gašić, and Minlie Huang. 2023b. [ConvLab-3: A flexible dialogue system toolkit based on a unified data format](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 106–123, Singapore. Association for Computational Linguistics.

## **A Results for Different Relation Types**

See Table 4 for results on the different ontology relation types.

## **B Example Completions for One-shot Equivalence Relation Prompt**

See Table 5 for example completions for the one-shot equivalence prompt.

Approach	Relation Type	F1-Score	Precision	Recall
<i>One-shot example from MultiWOZ</i>				
<i>Baseline: Separate relation prediction</i>	all	7.4	8.8	6.4
	domain-slot	3.9	2.1	41.0
	slot-value	9.4	18.8	6.3
	value-domain	8.0	12.7	5.9
	equivalence	1.8	1.0	6.9
+ constrained decoding	all	8.5	5.7	17.3
	domain-slot	1.1	0.5	50.7
	slot-value	9.3	7.0	13.8
	value-domain	14.4	11.5	19.5
	equivalence	1.4	0.7	16.1
+ CoT decoding	all	9.2	8.8	9.6
	domain-slot	2.5	1.3	50.7
	slot-value	16.0	18.4	14.2
	value-domain	7.7	12.7	5.5
	equivalence	1.8	1.0	7.3
+ constrained CoT decoding	all	9.2	6.4	15.9
	domain-slot	1.2	0.6	50.7
	slot-value	13.4	10.3	19.2
	value-domain	12.3	11.6	13.1
	equivalence	1.2	0.7	10.3
<i>Fine-tuning on MultiWOZ</i>				
<i>Baseline: Fine-tuning on MultiWOZ</i>	all	10.8	6.7	28.4
	domain-slot	5.9	3.2	49.3
	slot-value	6.5	3.7	25.3
	value-domain	20.5	15.3	30.8
	equivalence	4.3	2.5	18.8
+ constrained decoding	all	11.1	6.8	30.4
	domain-slot	4.7	2.5	52.9
	slot-value	6.8	3.9	25.5
	value-domain	19.8	13.9	34.1
	equivalence	3.9	2.1	21.5
+ CoT decoding	all	9.3	6.3	17.4
	domain-slot	3.9	2.0	49.3
	slot-value	5.7	3.5	15.5
	value-domain	16.7	14.8	18.9
	equivalence	2.3	4.4	1.5
+ constrained CoT decoding	all	14.4	10.0	25.8
	domain-slot	3.2	1.6	64.9
	slot-value	12.1	9.0	18.4
	value-domain	19.2	13.9	30.9
	equivalence	4.7	2.8	16.7
<i>Upper Bounds using SGD Data</i>				
One-shot example from SGD + separate relation prediction	all	12.9	10.7	16.4
	domain-slot	3.4	1.8	46.3
	slot-value	17.7	20.9	15.4
	value-domain	13.4	11.0	17.1
	equivalence	2.5	1.4	9.1
Fine-tuning on SGD	all	36.8	26.4	60.9
	domain-slot	7.3	3.8	88.8
	slot-value	30.6	20.7	58.9
	value-domain	53.5	46.9	62.1
	equivalence	12.3	6.9	57.6

Table 4: Ontology Relation Prediction Results for each relation type on the SGD test set.

Approach	Response
Baseline Greedy Decoding	<p>**Domain:** Hotel  **Slot:** Type of hotel  **Value:** Price range  - [price range, refers to same concept as, price] - [reference number, refers to same concept as, reference]  **Domain:** Price  **Slot:** Cost per night  **Value:** \$ 63  - [cost per night, refers to same concept as, price]  **Domain:** Location  **Slot:** Street address  **Value:** 135 carlingview drive etobicoke, ontario m9w 5e7, canada  - [street address, refers to same concept as, location] ...</p>
Constrained Decoding	<p>**Domain:** Hotel  **Slot:** Type of hotel  **Value:** Price range  - [price per night,refers to same concept as,135 carlingview drive etobicoke, ontario m9w 5e7, canada]  - [best western plus toronto north york hotel,refers to same concept as,best western plus toronto north york hotel &amp; suites] ...</p>
CoT-Decoding	<p>Here are the predicted synonyms or semantically equivalent terms from the provided dialogue:  **Price range:**  - [price range, refers to same concept as, price]  **Reference number:**  - [reference number, refers to same concept as, reference]  **Hotel:**  - [hotel, refers to same concept as, best western plus toronto north york hotel &amp; suites]  **Stay:**  - [stay, refers to same concept as, nights]  **Price:**  - [price, refers to same concept as, cost]&lt;eos&gt;</p>
Constrained CoT-Decoding	<p>Here are the predicted synonyms or semantically equivalent terms from the provided dialogue:  **Price range:**  - [price per night,refers to same concept as,cost]  **Reference number:**  - [price per night,refers to same concept as,check in date]  **Hotel:**  - [hotel,refers to same concept as,best western premier toronto airport carlingview hotel]  **Day:**  - [tuesday next week,refers to same concept as,check in date]&lt;eos&gt;</p>

Table 5: MWOZ one-shot only equivalence model example response excerpts for the different encodings on SDG test dialogue 100. For CoT-decoding methods, only the chosen branch is displayed.

# InteLLA: Intelligent Language Learning Assistant for Assessing Language Proficiency Through Interviews and Roleplays

Mao Saeki<sup>1,2</sup>, Hiroaki Takatsu<sup>1,2</sup>, Fuma Kurata<sup>1</sup>, Shungo Suzuki<sup>1</sup>,  
Masaki Eguchi<sup>1</sup>, Ryuki Matsuura<sup>1</sup>, Kotaro Takizawa<sup>1</sup>, Sadahiro Yoshikawa<sup>2</sup>,  
and Yoichi Matsuyama<sup>1,2</sup>,

<sup>1</sup>Waseda University, <sup>2</sup>Equemenopolis, Inc.

Correspondence: saeki@equ.ai

## Abstract

The primary challenge in utilizing dialogue systems for reliable language assessment for interactional skills lies in obtaining ratable speech samples that demonstrate the user's full range of ability. We thus developed a multimodal dialogue system that employs adaptive sampling strategies and enables a mixed initiative interaction through extended interview and roleplay dialogues. The interview is a system-led dialogue aimed at evaluating the user's overall proficiency. The system dynamically adjusts the question difficulty based on a real-time assessment to induce linguistic breakdowns, which provides evidence of the user's upper limits of proficiency. The roleplay, on the other hand, is a mixed-initiative, collaborative conversation intended to assess interactional competence such as turn management skills. Two experiments were conducted to evaluate our system in assessing oral proficiency. In the first experiment, which involved an interview dataset of 152 speakers, our system demonstrated high accuracy in automatically assessing overall proficiency. However, we observed that linguistic breakdowns were less likely to occur among high-proficiency users, indicating some room for further enhancing the ratability of speech samples. In the second experiment based on a role-play dataset of 75 speakers, the speech samples elicited by our system was found to be as ratable for interactional competence as those elicited by experienced teachers, demonstrating our system's capability in conducting interactive conversations. Finally, we report on the deployment of our system with over 10,000 students in two real-world testing scenarios.

## 1 Introduction

Language testing plays a critical role in ensuring effective language learning, as it provides valuable feedback on learners' proficiency levels and guides instructional planning (Fulcher, 2010). Assessment of oral proficiency is particularly important, as speaking and listening skills are essential

for effective communication in a second language. Traditional methods of oral proficiency assessment, however, face several challenges, including the subjectivity of human raters and the difficulty of creating standardized, scalable testing environments (Galaczi and Taylor, 2018).

To address these challenges, several studies have explored automated systems for oral assessment. For example, Ockey and Chukharev-Hudilainen (2021) evaluated the potential of spoken dialogue systems (SDS) for paired oral discussion tasks, concluding that a standardized assessment may favor SDS over human interlocutors due to its systematic behavior. Recent advancements in large language models (LLMs) have further simplified the implementation of such dialogue tasks. However, a significant challenge remains in obtaining "ratable" speech samples that accurately represent the full extent of a learner's language capabilities. Assessment of oral proficiency requires not only measuring linguistic competence, such as grammar and vocabulary, but also evaluating interactional behaviours, including turn-taking, topic management, and repair strategies (McNamara, 1996). Additionally, to provide a reliable assessment, it is crucial to observe the upper linguistic limits of the user (Liskin-Gasparro, 2003). Therefore, an effective dialogue system must be capable of engaging users in a manner that naturally reveals these competencies while also being scalable as a testing tool.

To this end, we developed the Intelligent Language Learning Assistant, *InteLLA*, a multimodal dialogue system designed to elicit spontaneous speech samples from second language learners through a combination of a 15-minute interview and a 10-minute roleplay session. By dynamically adjusting the topic difficulty based on real-time assessments, the system aims to provoke linguistic breakdowns that serve as evidence of a learner's upper proficiency limits. Additionally, the mixed-initiative roleplay component is designed to evalu-



Figure 1: The InteLLA system for oral proficiency assessment. The user connects to an online video call with InteLLA from their web browser on PC, tablet or smartphone.

ate the user's interactional competence in a collaborative setting.

To ensure the functionality and potential limitations of our system for large-scale real-world implementation, this paper reports two experiments: Chapter 4 evaluates how well the system can assess oral proficiency through various experiments designed to test its efficacy; Chapter 5 reports on the field testing results of our system used in real-world testing scenarios with university and high school students in Japan. We also discuss our first year operation of our system in terms of practicality and social impacts.

## 2 Related Work

### 2.1 CEFR

The Common European Framework of Reference for Languages (CEFR) serves as a comprehensive foundation for the development of language syllabi and curricula, as well as the evaluation of foreign language proficiency (Council of Europe, 2020). According to the CEFR, the key competencies for effective language communication include range (vocabulary richness), accuracy (grammatical correctness), fluency (smoothness and flow of speech), interaction (ability to engage in conversational exchange), and coherence (Engaging in effective conversational exchange). These competencies are defined across six proficiency levels: A1, A2, B1, B2, C1, and C2, with A1 representing the beginner level and C2 indicating proficient or near-native speaker capabilities.

The CEFR outlines specific communicative activities referred to as "Can-Do" statements, which

articulate what learners at each proficiency level should be able to achieve. These "Can-Do" serve as guidelines to determine the appropriate level for a learner based on their demonstrated abilities in a certain social situation. For instance, at the B1 level, learners should be able to handle most situations likely to arise while traveling in an area where the language is spoken.

This standardization is particularly valuable in the development of dialogue systems for language testing, as it offers an established baseline for designing tasks, including the interlocutor's behaviors, and evaluating user performance in a reliable and valid manner.

### 2.2 Oral Proficiency Interview

In many computerised speaking assessments, the user is given a reading script or situational explanation and is then required to record their speech. Such monologue-based speaking score, however, only have moderate correlations with those elicited in interactive dialogue tasks (Roever and Ikeda, 2022). On the other hand, due to their dynamic and co-constructive nature, dialogic tasks inevitably introduce variability in examiner behaviours and thus affect the learner's performance in the test (Galaczi and Taylor, 2018). This inherent variability poses a challenge for maintaining consistent and reliable assessments in dialogue-based tasks.

To draw out such interactive abilities, interview-based assessments of speaking proficiency conducted by trained professionals have long been considered, a representative implementation being the ACTFL-OPI (Liskin-Gasparro, 2003). The ACTFL-OPI interview consists of several phases. It begins with a "warm-up" where the interviewer asks questions or engages in small talk to familiarize the examinee with the test. Through this warm-up, the interviewer conducts a preliminary evaluation to decide the difficulty level of the first main topic. Next, the main part of the assessment, the "iterative process" takes place. The interviewer alternates between questions that are perceived as comfortably easy and challengingly difficult for the examinee to induce signs of "breakdown". Typically, breakdowns are indicated by hesitation, stumbling, lack of response, or rephrasing. This iterative process continues until sufficient information is obtained to assess the examinee's proficiency accurately. Automated assessment systems that mimic this interview strategy, such as the ACTFL Oral Proficiency Computer, exist. However, these sys-

tems do not rely on the user’s previous responses but rather output a predefined list of questions sequentially (Isbell and Winke, 2019). Although some measures are taken such as adjusting the difficulty of questions based on self-assessment before the interview, dynamic level adjustments during the interview, as performed by human experts, are not conducted.

Research into systems that conduct interview or counseling-like dialogues has been extensive in the domains other than language testing (Morbini et al., 2014; Inoue et al., 2020). These systems aim to elicit user speech through natural listening and question generation, but few explicitly evaluate user performance. Additionally, there is considerable research on using dialogue systems for speaking proficiency assessment (Ramanarayanan et al., 2019; Litman et al., 2016), but these studies generally assign the same tasks to all users from the perspective of test fairness and avoiding dialogue breakdowns.

### 2.3 Roleplay Dialogue

While structured interaction tasks such as the ACTFL-OPI have been used extensively to elicit ratable samples to assess linguistic competence (e.g., vocabulary, grammar, pronunciation), it falls short in assessing a full range of interactional competence. As such, language assessment researchers attempt to incorporate roleplay tasks in their tests to simulate authentic social settings for the examinees to demonstrate their abilities to enact simulated social roles by maintaining interpersonal relationships and managing turn-taking in a collaborative and cooperative manner (Kasper and Youn, 2018). By design, such roleplay dialogues should involve mixed-initiative interactions where both the system and the user can take the lead in conversation. This requirement is essential to making it possible to evaluate how well the learner handles unexpected turns and engages in collaborative communication.

### 2.4 System Requirements

Based on the aforementioned considerations, our system needs to effectively assess oral proficiency through both structured interviews and collaborative roleplay interactions. To achieve this, we have established the following requirements for the conversational agent being developed in this project:

1. **Adaptive speech sampling strategy:** The system should ask relevant questions and provide

responses tailored to the user’s language level, efficiently sampling ratable speech data for assessment. Multimodal interaction, including non-verbal gestures, is needed to elicit authentic speech, ensuring that scores are generalizable to real-world communication.

2. **Mixed-initiative interaction:** The system should enable collaborative, mixed-initiative dialogues, wherein both the system and the user can dynamically control the conversation. This will allow users to demonstrate their interactional competence, including aspects such as turn-taking and topic development.
3. **Scalability:** To ensure the test is accessible and fair for a diverse user base, the system must be usable across different locations and operable on low-end devices.

## 3 System Design

The InteLLA system is a multimodal dialogue system where the user connects to an online video call from their personal device, as shown in Figure 1. We adopted a modular architecture, wherein multiple modules, each responsible for specific dialogue capabilities such as ASR, operate concurrently to enable fully-duplex communication (Figure 2). For the ASR module, we employ the Google Text-to-Speech service. The details of the other modules will be discussed in subsequent sections.

### 3.1 Video Communication Module

To enable users to access the system via video call directly from their web browser, the system is hosted on a server, with agent audio and visuals streamed to the user through a Web Real-Time Communication (WebRTC) solution. This setup leverages server-side GPU resources for machine learning and rendering, ensuring a rich conversational experience even for users with low-end devices. This configuration is crucial for maintaining equitable and consistent testing.

### 3.2 Dialogue Management Module

LLMs have greatly simplified the design and management of dialogues by enabling the specification of conversation rules through prompts (Brown et al., 2020). However, these models often struggle to maintain coherence when the input (i.e. prompts and dialogue history) becomes excessively long. This poses a particular challenge in our use case,

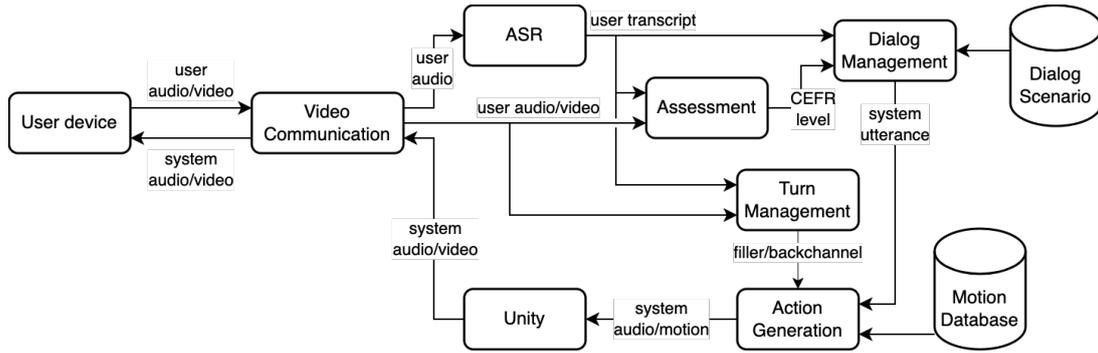


Figure 2: System architecture of IntelLA, comprising video communication, automatic speech recognition, dialogue management, turn management, action generation and assessment modules.

where a single conversation may extend from 20 to 30 minutes. Additionally, altering dialogue content based on real-time assessments for adaptive testing remains an issue.

To address these challenges, we employed a hybrid approach that combines LLMs with scenario-based dialogue management. Specifically, we segmented the interviews and roleplays into multiple topics, with each topic having a sub-goal such as "asking about hobbies" or "conducting a roleplay to borrow a PC from the user," designed to be completed within a 3 to 5-minute timeframe. These topics are tailored for each CEFR level.

Following the OPI framework described in Section 2.2, the conversation initiates with a warm-up topic, designed to make the user comfortable with the system. During the conversation, the assessment module, explained in Section 3.5, evaluates the user’s proficiency. Based on this assessment, users are assigned a topic that matches or slightly exceeds their proficiency level. This aims to induce linguistic breakdowns, thereby efficiently observing the user’s upper proficiency limits, as shown in Appendix A.1.

When a topic changes, the prompt for the LLM is updated, and the dialogue history is reset. By compartmentalizing conversations in this manner, the LLM can adhere to strict instructions for each individual topic, ensuring coherent and controlled dialogue over a total duration of 20 to 30 minutes. To maintain the memory of previous topics, we summarize earlier dialogue segments and incorporate these summaries into the updated prompts.

To enhance the ratatability of speech samples, a panel of applied linguistics researchers and experienced teachers carefully designed and piloted prompts. Following the literature on the correspondence between representative linguistic functions and CEFR levels (O’Sullivan et al., 2002), top-

ics were decided in terms of how likely learners are to use target linguistic functions in response to the system’s question. For instance, B2-level learners are expected to have the ability to produce a longer, coherent utterance, and thus the topics, for instance, should require them to compare and contrast multiple ideas. These are combined with generic prompts such as persona of the agent, guideline for the interview and summarized history, and fed to a LLM to generate the next system utterance. We use OpenAI’s GPT for the utterance generation.

### 3.3 Turn Management Module

In spoken dialogue systems, knowing when to speak is as crucial as knowing what to say for maintaining smooth interaction (Skantze, 2021). This is particularly important in the context of oral proficiency testing, where users often produce long pauses between sentences as they formulate their responses, increasing the likelihood of system interruptions.

During these pauses, it is often discernible whether the user intends to continue speaking or has finished based on grammatical completeness, prosody, and eye gaze. To utilize such multimodal cues, we trained an end-of-turn detection model that incorporates text, audio, and image data to predict whether the user has finished their turn, as proposed by Kurata et al. (2023). However, turn overlaps are inevitable, even in human conversation. Not all overlaps are detrimental; for instance, the user may simply be providing backchannel feedback to the system. To determine whether the system should continue speaking or pause when turn overlaps occur, we implemented a barge-in detection system based on the overlap resolution model by Gervits and Scheutz (2018).

This module is also responsible for generating

backchannels and fillers. Backchannels are necessary cues to indicate the system is listening to the user, thereby encouraging the user to speak more. Verbal and non-verbal backchannels are generated at the end of clauses. Fillers signify the system’s intention to speak and avoid awkward pauses between turns, which may happen due to latency introduced by utterance and action generation. A filler utterance is generated when the user’s end-of-turn is detected and the system’s next utterance does not begin immediately after.

### 3.4 Action Generation Module

While text-to-speech (TTS) has been extensively studied, body and facial motion generation have received comparatively less attention. Although early linguistic-inspired rule-based gesture generation approaches were proposed (Cassell et al., 2001), few end-to-end models exist that use audio and text input to generate body gesture data (Kucherenko et al., 2020). However, the end-to-end models are not fast enough for real time communication. Additionally, while such models can create smooth movements synchronized with speech rhythm, they often struggle to generate semantic gestures that are essential for making conversations engaging.

To achieve real-time generation of natural body facial motions, we employed a database-driven approach. First, we constructed a database of actions performed by a motion actor, with each action mapped to corresponding text descriptions. When generating a motion, the input text is compared to the texts in the database to calculate the cosine similarity of embedded texts. The action most similar to the input text is then selected. Speech is generated using a TTS model, and mouth movements are generated based on vowel sounds estimated from the synthesized speech.

The combined data for speech, body and facial motions are then sent to a game engine for the agent animation to be rendered. Specifically, we used Sentence-BERT (Reimers and Gurevych, 2019) for text embedding, Google Text-to-Speech for TTS, and Unity for rendering the agent.

### 3.5 Assessment Module

We propose a speaking proficiency assessment model that takes multimodal dialogue data obtained during the conversation with the user, and simultaneously predicts proficiency levels across one holistic criterion (overall) and five analytic criteria: range, accuracy, fluency, phonology, and co-

herence. The model has multiple encoder modules to consider a wide range of multimodal features theoretically important in language assessment, such as vocabulary richness (Eguchi and Kyle, 2020), grammatical accuracy (Murakami and Ellis, 2022), fluency (Matsuura et al., 2022; Suzuki et al., 2021), goodness of pronunciation (Saito and Plonsky, 2019), and coherence of discourse (Qin, 2022). To capture these linguistic features, each encoder module has a model as a feature extractor that have been pre-trained in various natural language processing tasks such as grammatical error correction (Omelianchuk et al., 2020), coreference resolution (Otmazgin et al., 2023), and pronunciation scoring (Zhang et al., 2021). The inputs of the model are the user’s audio and video, speech-recognized text, and the system’s utterance text. After various linguistic features are extracted from these input data by the encoder modules, the outputs of each encoder module are blended by the transformer encoder (Vaswani et al., 2017). Then, the vector sequences, in which the influence of the interaction of the various linguistic features is embedded by the transformer encoder, are input to each network specialized for proficiency assessment of each CEFR category. The output layers for each CEFR category with softmax as activation function output the likelihood of each level. The probabilities are converted to a continuous value score  $x$  by the following equation:  $x = \sum_{c=1}^6 c \times p_c$  where  $p_c$  represents the probability of level  $c$  (1:A1, 2:A2, ..., 6:C2) in a category ( $\sum_{c=1}^6 p_c = 1$ ). After computing the discrete level boundaries of A1-C2 so as to maximize Quadratic Weighted Kappa (QWK) in the validation dataset based on  $x$ , a normalized score  $x'$  is fed back to the learner so that the boundaries of each level are evenly spaced: A1:[0, 1.0], A2:(1.0, 2.0), ..., C2:(5.0, 6.0).

The model was trained on 232 interview dialogues previously collected, and rated for the CEFR score by trained raters. Figure 3 shows an example of the assessment presented to the user. Rationales for the assessment are provided for each category and proficiency, based on the CEFR.

## 4 Experiments

To evaluate the system in terms of its capability of eliciting ratable speech samples, we conducted two experiments. The first experiment was designed to test the system’s adaptive speech sampling strategy in system-led interview dialogues in terms of

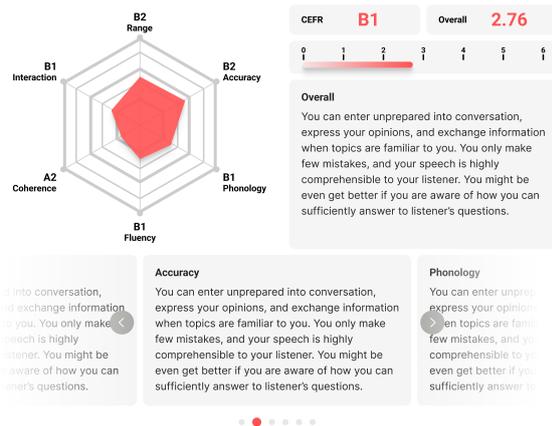


Figure 3: Example of assessment result, including the six core competencies defined by CEFR and the overall score, along with the rationale for these scores.

scoring accuracy as well as the frequency of target phenomenon, that is, linguistic breakdown. The second experiment was set up to gauge the quality of mixed-initiative interaction in roleplay tasks. Given the multifaceted nature of interactional features, the second experiment aims at holistically evaluating the system using human experts' ratings, comparing the scoring reliability between human-interlocutors and the current system.

#### 4.1 Ethical Statement

All data collection for this study, including field testing were reviewed and approved in advance by the ethical review committee ("Ethics Review Procedures concerning Research with Human Subjects") of Waseda University. Prior to all experiments, a consent form outlining the experimental procedures and the use of data (specifically that the recorded audio and video data would be used exclusively for research purposes) was explained to the participants. For high school participants, the procedure was explained to both them and their parents or guardians. Consent was obtained through a detailed consent form, ensuring all parties were fully informed before participation.

#### 4.2 Interview Experiment

We recruited 152 university students with varying levels of oral proficiency to participate in an interview session with our system. Among the participants, 94 participants were female, and 58 were male, with an average age of 20. Each user were given 4 topics, and the whole interview lasted around 15 to 20 minutes. The recordings from these interviews were assessed for CEFR levels by the three trained raters, all of whom hold MA de-

		Pred					
		A1	A2	B1	B2	C1	C2
True	A1	4	3	0	0	0	0
	A2	0	25	6	0	0	0
	B1	0	2	32	7	0	0
	B2	0	0	3	49	1	0
	C1	0	0	0	2	11	3
	C2	0	0	0	0	0	4

Table 1: Confusion Matrix of automatic assessment ("Pred") and the gold standard by the trained human raters ("True").

grees in TESOL or equivalent as well as more than 5 years of teaching experience, and completed a rater training program conducted by researchers in Applied Linguistics. The inter-rater reliability for the CEFR assessment was measured using QWK, which ranged from 0.800 to 0.835, indicating high consistency among raters. In instances of disagreement between raters, the true label was determined through discussion. We then compared the final scores from the assessment module to the human raters' scores (gold standard). The QWK between our system and the gold standard was 0.929, demonstrating very high reliability. The confusion matrix, comparing the model's predictions with human ratings, is shown in Table 1. As evident from the confusion matrix, all model predictions were within one level of the human scoring.

Next, recordings were evaluated for linguistic breakdowns by the same raters. A breakdown was defined as "failure to manage to maintain their speech or respond to the question sufficiently," following the criteria established in (Isbell and Winke, 2019). The occurrence of breakdowns observed in a recording for each proficiency level was observed as follows: A1 and A2 – 100%, B1 – 79.4%, B2 – 42.9%, C1 – 20.9%. C2 proficiency level participants were excluded from this analysis since, theoretically, they would not exhibit breakdowns. These results indicate that students with higher proficiency experienced fewer breakdowns. This trend is expected, as higher proficiency learners, particularly those at B2 or higher levels, may employ a range of linguistic repertoires to strategically navigate around breakdowns (Council of Europe, 2020). However, such strategic behavior can influence other aspects of utterances, including lexical richness and circumlocution. Therefore, the system's adaptive sampling strategy should be evaluated with these considerations in mind. Given the consistency of ratings across levels, it is plausible to argue that despite some room for improve-

ment especially for advanced learners, the current adaptive sampling strategy can elicit ratable speech samples from learners at various proficiency levels.

### 4.3 Roleplay Experiment

We recruited a total of 75 university students for the roleplay data collection. Among the participants, 54 were female, 20 were male, and one participant did not answer, with an average age of 20. Each participant completed two roleplay sessions with a one week interval: one with a human examiner and one with our system. The order of interlocutor conditions was counterbalanced across participants. Five experienced English tutors were randomly assigned to each student to complete the roleplay in the human session. We adapted a roleplay task used previously in the context of second language assessment literature (Al-Gahtani and Roever, 2018), shown in Appendix A.2. Upon completion of the data collection, four experienced tutors (recruited from the same pool of the examiners) rated each session recording in terms of interactional competence (IC) (Galaczi and Taylor, 2018). Since there was no established rating scale for the assessment of IC, we developed our own CEFR-inspired IC scale. Given our focus on mixed initiatives in interaction, we decided to include two relevant components of IC: **Turn-management** and **Topic-management**. Turn-management is defined as the ability to sustain a cooperative and collaborative conversation through appropriate turn-taking, Topic-management pertains to developing ideas collaboratively toward the intended interactional outcome. The detailed descriptors are shown on Table 7 and 6 in Appendix B.3.

Using a spiral rating design (Eckes, 2015), students' performances were evaluated by alternating pairs of two raters, and each rater assessed only one of the student's videos to mitigate bias such as halo effects. This resulted in a total of 528 raw data points in a 6-level ordinal scale from A1 to C2 (i.e., 66 students  $\times$  2 interlocutor types  $\times$  2 raters  $\times$  2 rating criteria). IC dimensions that could not be observed in the video were marked as unratable.

To evaluate the extent to which our system elicited speech samples that are informative for IC assessment (i.e., ratability), we compared the scoring reliability of IC ratings between the interlocutor conditions (human tutors vs. the system). To systematically control for the effects of rater severity and examinees' proficiency levels, the reliability index was estimated through a series of

Many-Facet Rasch Modeling (MFRM; for details, see Appendix B) (Eckes, 2015). Results revealed the comparable level of reliability between the interlocutor conditions of human tutors (0.767) and our system (0.771). See Appendix B.1 for details.

The infit/outfit statistics based on the Rasch model indicates that the AI-based roleplay followed more closely with the assumption of the Rasch measurement model (see Table 4 in Appendix B.2). Taken together, these findings suggest that speech samples elicited through our system are as ratable as human interlocutors for IC assessment, and the system yields psychometrically more consistent data for assessing IC components related to mixed initiatives than human tutors.

### 4.4 Discussion

The interview experiment demonstrated that the InteLLA system can elicit ratable speech samples for oral proficiency assessment, evidenced by high inter-rater reliability both among human raters and between human and our system. However, we also found the low rate of linguistic breakdowns among high-level participants. This could be attributed to their problem-solving strategies. This suggests that there should be some room for enhancing the ratability of speech samples. Future work, for instance, will need to engage with the accuracy of real-time assessment mechanisms that can operate effectively with fewer samples.

Conversely, the roleplay experiment showed that our system can sufficiently elicit interactional competence for human ratings, specifically turn management and topic management, on par with human interlocutors. Future work includes extending the assessment model with the capability to automatically evaluate interactional competence.

## 5 Field Testing

We report on the system's performance and stakeholders' satisfactions in the real-world scenarios with university and high school students in Japan.

### 5.1 Field Testing with University Students

Over the past year (AY2023), the InteLLA system has been deployed to provide assessments to over 10,000 Japanese university students. The system served as a middle-stakes test, where the results were used to determine the appropriate English class level for each student. Tests were administered remotely, with students using their personal computers from home.

To evaluate the system's performance, we randomly selected 300 recordings for detailed analysis. These recordings were scored according to the CEFR level by three trained raters, with the final score determined by majority vote in cases of disagreement. The reliability of the automatic assessment, when compared to human ratings, was found to be 0.869, demonstrating a high level of reliability. However, three out of 300 recordings were deemed unratable, indicating they could not be scored reliably due to technical problems. The user's audio input was too small for ASR to recognize and for the dialogue management to keep the conversation coherently. These issues should have caused significant delays of the system responses due to network problems and consequently interfere with ratable speech elicitation. The results of this field test demonstrate that our system can provide accurate oral proficiency assessments even in uncontrolled, real-world scenarios.

## 5.2 Field Testing with High School Students

As another field study, a total of 97 students in Chiba prefecture in Japan, all aged 16, participated in eight English conversation sessions over a period of one month. The first and last sessions served as a pretest and posttest and were conducted using the interview scenarios. The other sessions in between engage them with daily conversations similar to the roleplay format in the second experiment. After each session, students completed a brief questionnaire assessing their learning motivation.

The pretest and posttest scores were compared using a linear mixed-effects model to estimate the group-level improvement, including the random-effect variable of participants to controlling for individual variability in the pretest scores. The analysis revealed a significant increase of 0.30 points ( $p < 0.001$ ) out of 6.0. Among various patterns of score changes, we found A2-level students at the pretest significantly improved and reached B1-levels at the time of the posttest. Students who exhibited notable improvement in this category also showed a positive trend in survey responses over time. These responses included "Enjoyment of the conversations", "Feeling of being able to express themselves", "Comfort and relaxation while speaking", "Desire to speak more in English." Notably, we adopted intact classes for this field-testing study, meaning that these improvements may not solely be attributed to interactions with IntelLLA but also

to the students' regular English classes during the period of the study. these findings may indicate the potential of using our system as learning materials for English speaking skills. However, this study demonstrates the potential of using multimodal dialogue systems such as IntelLLA for developing English speaking skills and language learning motivation.

## 6 Conclusions and Future Directions

In this paper, we presented IntelLLA, a multimodal dialogue system designed for the assessment of oral proficiency. IntelLLA is designed to elicit ratable speech samples that display the user's full range of interactional skills. To enhance the ratability of speech elicited, the system is required to adaptively change the difficulty levels of questions to collect learners' linguistic breakdowns as the evidence of their upper limit of proficiency. To capture learners' ability to maintain collaborative conversations, the system is expected to enable mixed initiative interaction where learners need to engage with turn-taking management and topic development. To evaluate IntelLLA's usefulness in oral proficiency assessment, we conducted two experiments using interview and roleplay conversations. The results from the interview conversations demonstrated that our system consistently elicited ratable samples especially for lower-level learners, and automatically estimated scores based on those samples exhibited a high agreement with experts' ratings. In roleplay conversations, the ratability of the speech samples elicited through IntelLLA were comparable to experienced English teachers in terms of scoring reliability by human raters. The current study takes an interdisciplinary approach to integrating research on multimodal dialogue systems into real-world problems in language learning and testing. We demonstrated that IntelLLA is suitable for middle-stake assessments, effectively scaling to accommodate a large number of users. The scalability and accuracy position IntelLLA as a valuable tool for oral proficiency assessment in varied contexts.

Although our system showed the potential for reliable standardised assessment of oral proficiency, several challenges should be acknowledged to minimize possible negative consequences of its further real-world implementations (Dai et al., 2024). In the large-scale deployment of IntelLLA, Quality of Service (QoS) (Aurrecochea et al., 1998) factors,

such as latency derived from the inference model and communication architecture, as well as the resolution of IntelLLA's rendering engine, are critical from the perspectives of Quality of Experience (QoE) (Möller and Raake, 2014) and cost. Future work is expected to evaluate how the configuration of the developed large-scale dialogue system impacts user engagement (Kurata et al., 2024), considering a dialogue quality management system that can efficiently facilitate system improvements.

## 7 Acknowledgement

The research presented in this study was achieved through funding and support from the project "Technological Development for Next-Generation Artificial Intelligence Evolving with Humans (JPNP20006) / Development of an Online Language Learning Support AI System that Grows with Humans" by New Energy and Industrial Technology Development Organization (NEDO), and "Beyond 5G Seed Creation Program / Research and Development of an XR Communication Infrastructure for Realizing High-Immersion Interaction Experiences with Conversational AI Agents (JPJ012368C06301)" by the National Institute of Information and Communications Technology (NICT).

## References

- Saad Al-Gahtani and Carsten Roever. 2018. Proficiency and preference organization in second language refusals. *Journal of Pragmatics*, 129:140–153.
- Cristina Aurrecochea, Andrew T Campbell, and Linda Hauw. 1998. A survey of qos architectures. *Multimedia systems*, 6:138–151.
- Trevor G Bond and Christine M Fox. 2013. *Applying the Rasch model: Fundamental measurement in the human sciences*. Psychology Press.
- 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.
- Justine Cassell, Hannes Högni Vilhjálmsson, and Timothy Bickmore. 2001. Beat: the behavior expression animation toolkit. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 477–486.
- Council of Europe, editor. 2020. *Common European framework of reference for languages: learning, teaching, assessment ; companion volume*. Council of Europe Publishing, Strasbourg.
- David Wei Dai, Shungo Suzuki, and Guanliang Chen. 2024. [Generative ai for professional communication training in intercultural contexts: where are we now and where are we heading?](#) *Applied Linguistics Review*.
- R.J. de Ayala. 2022. *The Theory and Practice of Item Response Theory*. Methodology in the Social Sciences Series. Guilford Publications.
- Thomas Eckes. 2015. *Introduction to many-facet Rasch measurement: Analyzing and evaluating rater-mediated assessments*.
- Masaki Eguchi and Kristopher Kyle. 2020. Continuing to explore the multidimensional nature of lexical sophistication: The case of oral proficiency interviews. *Modern Language Journal*, 104(2):381–400.
- Glenn Fulcher. 2010. [Practical language testing](#). *Practical Language Testing*, pages 1–352.
- Evelina Galaczi and Lynda Taylor. 2018. Interactional competence: Conceptualisations, operationalisations, and outstanding questions. *Language Assessment Quarterly*, 15(3):219–236.
- Felix Gervits and Matthias Scheutz. 2018. [Pardon the interruption: Managing turn-taking through overlap resolution in embodied artificial agents](#). In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*, pages 99–109, Melbourne, Australia. Association for Computational Linguistics.
- Koji Inoue, Kohei Hara, Divesh Lala, Kenta Yamamoto, Shizuka Nakamura, Katsuya Takanashi, and Tatsuya Kawahara. 2020. A Job Interview Dialogue System That Asks Follow-up Questions: Implementation and Evaluation with an Autonomous Android. *Transactions of the Japanese Society for Artificial Intelligence*, 35(5):D–K43 1–10.
- Dan Isbell and Paula Winke. 2019. Actfl oral proficiency interview—computer (opic). *Language Testing*, 36(3):467–477.
- Gabriele Kasper and Soo Youn. 2018. [Transforming instruction to activity: Roleplay in language assessment](#). *Applied Linguistics Review*, 9:589–616.
- Taras Kucherenko, Patrik Jonell, Sanne van Waveren, Gustav Eje Henter, Simon Alexanderson, Iolanda Leite, and Hedvig Kjellström. 2020. Gesticulator: A framework for semantically-aware speech-driven gesture generation. In *Proceedings of the ACM International Conference on Multimodal Interaction*.

- Fuma Kurata, Mao Saeki, Masaki Eguchi, Shungo Suzuki, Hiroaki Takatsu, and Yoichi Matsuyama. 2024. Development and validation of engagement and rapport scales for evaluating user experience in multimodal dialogue systems. In *Proceedings of the 14th International Workshop on Spoken Dialogue Systems Technology, Hokkaido, Japan*, pages 1–14.
- Fuma Kurata, Mao Saeki, Shinya Fujie, and Yoichi Matsuyama. 2023. [Multimodal Turn-Taking Model Using Visual Cues for End-of-Utterance Prediction in Spoken Dialogue Systems](#). In *Proc. INTERSPEECH 2023*, pages 2658–2662.
- Judith E. Liskin-Gasparro. 2003. [The ACTFL Proficiency Guidelines and the Oral Proficiency Interview: A brief history and analysis of their survival](#). *Foreign Language Annals*, 36(4):483–490.
- Diane Litman, Steve Young, Mark Gales, Kate Knill, Karen Ottewell, Rogier van Dalen, and David Vandyke. 2016. Towards using conversations with spoken dialogue systems in the automated assessment of non-native speakers of english. In *Proceedings of the 17th annual meeting of the special interest group on discourse and dialogue*, pages 270–275.
- Ryuki Matsuura, Shungo Suzuki, Mao Saeki, Tetsuji Ogawa, and Yoichi Matsuyama. 2022. Refinement of utterance fluency feature extraction and automated scoring of L2 oral fluency with dialogic features. In *Proceedings of the 14th Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*, page 1312–1320.
- Tim McNamara. 1996. [Measuring second language performance](#).
- Sebastian Möller and Alexander Raake. 2014. *Quality of experience: advanced concepts, applications and methods*. Springer.
- Fabrizio Morbini, David Devault, Kallirroi Georgila, Ron Artstein, David Traum, and Louis-Philippe Morency. 2014. A Demonstration of Dialogue Processing in SimSensei Kiosk. In *SigDial*, pages 254–256. Association for Computational Linguistics.
- Akira Murakami and Nick C. Ellis. 2022. Effects of availability, contingency, and formulaicity on the accuracy of english grammatical morphemes in second language writing. *Language Learning*, 72(4):899–940.
- Gary J Ockey and Evgeny Chukharev-Hudilainen. 2021. [Human versus Computer Partner in the Paired Oral Discussion Test](#). *Applied Linguistics*, 42(5):924–944.
- Kostiantyn Omelianchuk, Vitaliy Atrasevych, Artem Chernodub, and Oleksandr Skurzhanyskiy. 2020. GECToR – grammatical error correction: Tag, not rewrite. In *Proceedings of the 15th Workshop on Innovative Use of NLP for Building Educational Applications*, page 163–170.
- Shon Otmazgin, Arie Cattan, and Yoav Goldberg. 2023. LingMess: Linguistically informed multi expert scorers for coreference resolution. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, page 2752–2760.
- Barry O’Sullivan, Cyril J Weir, and Nick Saville. 2002. Using observation checklists to validate speaking-test tasks. *Language testing*, 19(1):33–56.
- Jie Qin. 2022. Potential contribution of SFL to task-based research: An examination of planning effects using genre-based theme analysis. *System*, 104:1–12.
- Vikram Ramanarayanan, Matthew Mulholland, and Yao Qian. 2019. Scoring interactional aspects of human-machine dialog for language learning and assessment using text features. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 103–109.
- 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.
- Carsten Roever and Naoki Ikeda. 2022. What scores from monologic speaking tests can (not) tell us about interactional competence. *Language Testing*, 39(1):7–29.
- Kazuya Saito and Luke Plonsky. 2019. Effects of second language pronunciation teaching revisited: A proposed measurement framework and meta-analysis. *Language Learning*, 69(3):652–708.
- Emanuel A Schegloff and Harvey Sacks. 1973. Opening up closings.
- Gabriel Skantze. 2021. Turn-taking in conversational systems and human-robot interaction: a review. *Computer Speech & Language*, 67:101178.
- Shungo Suzuki, Judit Kormos, and Takumi Uchihara. 2021. The relationship between utterance and perceived fluency: A meta-analysis of correlational studies. *Modern Language Journal*, 105(2):435–463.
- 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 *Proceedings of the 31st Conference on Neural Information Processing Systems*, pages 5998–6008.
- Junbo Zhang, Zhiwen Zhang, Yongqing Wang, Zhiyong Yan, Qiong Song, Yukai Huang, Ke Li, Daniel Povey, and Yujun Wang. 2021. [speechocean762: An open-source non-native English speech corpus for pronunciation assessment](#). In *Proceedings of the 22rd Annual Conference of the International Speech Communication Association*, page 3710–3714.

## A Dialogue Samples

### A.1 Interview Example

An example of an interview dialogue is shown below. The first to third parts of the adjacency pairs (Schegloff and Sacks, 1973) are referred to as 1st, 2nd, and 3rd, respectively.

**Warm-up phase :** First, in the warm-up phase, relatively simple conversations such as greetings and small talk are conducted to ease tension and get a general sense of the user's level. In this example, the proficiency assessment model estimates the user to be at the A1 level, and then it moves on to the level-check phase to confirm this estimation.

InteLLA: What is your favorite season? (1st - **topic initiation**)

User: My favorite season is winter. (2nd)

InteLLA: Are there any activities you like to do in winter? (1st - **additional question**)

User: Uh ... Ski and making snowman. (2nd)

InteLLA: That sounds like a lot of fun. (3rd - **acknowledgement**) Could you tell me more about it? (1st - **continuation request**)

User: I like skiing with family. I go every year. (2nd)

**Level Check :** In this example, the level check was deemed successful because fluency was maintained during the A1 level conversation. Next, the process moves to the "probe" phase, where it is tested whether the user can maintain fluency with a slightly more difficult topic.

InteLLA: Alright. What did you eat for breakfast this morning? (1st - **topic initiation**)

User: I ate uh... Sandwich it is chicken and salad it is very delicious.(2nd)

InteLLA: Do you usually eat breakfast? (1st - **additional question**)

User: Uh yes I always eat breakfast. (2nd)

InteLLA: Oh, that's a healthy diet. (3rd - **acknowledgement**)

User: Yeah. (3rd)

InteLLA: I see what time do you usually eat breakfast. (1st - **additional question; A1 level**)

User: Uh seven A.M. I wake up and I go to kitchen and I eat breakfast. (2nd)

**Probe Phase :** In this example, a breakdown in fluency occurred during the A2 level conversation, indicating that the user's English conversation ability is estimated to be at the A1 level. In actual dialogues, this probe and level-check routine is repeated several times. Finally, the user's level is confirmed, and a closing dialogue is conducted during the cool-down phase.

InteLLA: Have you ever been to a foreign country? (1st - **topic initiation**)

User: Uh no. I never go to foreign country. (2nd)

InteLLA: Ok. which country would you like to visit in the future? (1st - **additional question**)

User: I would like visit ... Singapore. (2nd)

InteLLA: Why is that? (1st - **continuation request**)

User: Because I want visit ... I like go to nice ... ah nice ... (2nd - **breakdown detected**)

InteLLA: That's ok. Let's move on. (1st - **topic transition**)

### A.2 Roleplay Example

A conversation-analytic (CA) transcription of a roleplay dialogue from the experiment is shown in table 2 ("Lending a computer" scenario). Table 3 shows CA notations used in the transcript. In this scenario, InteLLA plays a role of a university friend of the user, asking to borrow a PC from them because hers broke down while she was working on the term paper. This roleplay scenario is an adapted version from previous study on second language learning and assessment (Al-Gahtani and Roever, 2018). This roleplay attempts to assess the following Can-Do statements for the user role:

- **(Conversation; B2+)** Can indicate reservations and reluctance, state conditions when agreeing to requests or granting permission, and ask for understanding of their own position.
- **(Conversation; B2)** Can sustain relationships with users of the target language without unintentionally amusing or irritating them or requiring them to behave other than they would with another proficient language user.

The example roleplay card based on (Al-Gahtani and Roever, 2018) is as follows.

### Roleplay card:

Read the following instructions carefully. You have 3 mins at maximum to prepare for this role play.

**Situation** You are a student. It's 11 pm. You're working on a term paper that is due tomorrow morning at 8 am. You are planning to spend all night finishing the paper. Now, you decided to take a quick 10-min break. You opened your phone and noticed a text message from A. Friend A lives two floors above you, but you have only known for a month. So you wondered what happened to A. The text message says that he/she wants to borrow your PC because theirs broke down. Because you have ONLY one computer (which you are using for writing the paper), you think it is inconvenient for you to lend it to another person. Now, friends A rings your doorbell and you are answering it.

### Task

- Explain your situation and first try to decline the request.
- Then negotiate for a solution that works for both of you.
- You can lend it to them but make sure that you secure enough time to finish your term paper.
- Do NOT show irritation or annoyance to the friend A.

- Examinee Interactional Competence (66 persons)
- Rater severity (4 raters)
- Criteria (Turn and Topic-management)
- Interlocutor difficulty (AI and 5 human tutors)

Following Eckes (Eckes, 2015), a full MFRM can be expressed in the following formula:

$$\ln \left[ \frac{p_{nljk}}{p_{nljk-1}} \right] = \theta_n - \beta_l - \eta_v - \alpha_j - \tau_k$$

where

- $p_{nljk}$  = probability of person  $n$  receiving a rating of  $k$  on criteria  $l$  from rater  $j$  when the interlocutor is  $v$ ,
- $p_{nljk-1}$  = probability of person  $n$  receiving a rating of  $k - 1$  on criteria  $l$  from rater  $j$  when the interlocutor is  $v$ ,
- $\theta_n$  = ability (= IC) of person  $n$ ,
- $\beta_l$  = difficulty of criteria  $l$ ,
- $\eta_v$  = difficulty of interlocutor  $v$ ,
- $\alpha_j$  = severity of rater  $j$ ,
- $\tau_k$  = difficulty of receiving a rating of  $k$  relative to  $k - 1$ .

This allows the estimation of locations of each constituent from each facet on a latent logit scale.

## B Many-Facet Rasch Modeling

Many-facet Rasch Modeling is a psychometric approach often used in performance assessment (i.e., type of assessment involving a set of raters evaluating performances of the test-takers on pre-determined criteria for their skill mastery) (Eckes, 2015). In performance assessments such as the roleplay in Experiment 2 (Section 4.3), multiple factors can add "noises" to the raw score, including but not limited to choice of raters, tasks used, and interlocutors. Many-facet Rasch Model attempts to account for the different sources of variation in the raw score (i.e., facets) and to transform the raw score into a latent logit score. In so doing, it simultaneously computes the harshness of rater and the difficulty of interlocutors on the same logit scale.

The design of experiment 2 yielded the following facets:

### B.1 Person separation reliability

After fitting a Rasch model, the reliability of ratings can be calculated to indicate the consistency of person's ability estimate (i.e., location on the latent logit scale) that is beyond the influence of other facets. This is calculated by dividing the amount of variation in Expected A Posteriori (EAP) estimates of person's abilities based on the Rasch model after considering other facets over the total amount of variation in persons' abilities. This reliability estimate ranges from 0 to 1, a high score indicating a high level of reliability. As presented in Section 4.3, the person separation reliability was .767 for human-based roleplays and .771 for AI-based roleplays. For more information about person separation reliability estimate see (Eckes, 2015; Bond and Fox, 2013).

Table 2: Conversation-analytic transcript of a roleplay dialogue ("Lending a computer" scenario)

Speaker	CA transcription
IntelLLA:	Ok in this roleplay, I will start talking. Are you ready?
User:	Yes I am.
IntelLLA:	Alright three two one. (0.6)
IntelLLA:	Oh hi. (0.5)
User:	[Hi: [Sorry to bother you. (0.3)
User:	°Hmm°, =
IntelLLA:	=Did you look at my message?
User:	.hh yes, I just opened it. I was actually working on my paper and I didn't have time tuh check [your message.
IntelLLA:	[Hmm. (2.8)
IntelLLA:	I see. (0.6) I'm in a bit of bind. [My computer died and I need=
User:	[Hmm,
IntelLLA:	=to finish my paper. (0.5)
User:	Mm-huh=
IntelLLA:	=Could I borrow your PC? (0.8)
User:	.hhh I really wish I could. but it's just that I'm still working on my term paper, a:nd [it's due] tomorrow mornin' at eight. a:nd
IntelLLA:	[hmm, ]
User:	.hh I was planning to spend all day:: finishing the pa- all night sorry (.) finishing the paper a:nd .hh I'm not really sure if I could lend it to you before that. (4.0)
IntelLLA:	Oh (.) I understand. (0.5) Your paper is important too. (0.6) Could I borrow it after you finish? (1.1)
User:	°U::m,° we:ll I wish I could if I could finish it by eight [.hhh] >I could lend it to you bu'< how many hours do you=
IntelLLA:	[Hmm,]
User:	=need (.) for that. (2.5)
IntelLLA:	I think need about three hours. (0.8)
User:	[three hours. °right° ]
IntelLLA:	[Could I have it from eight] to eleven? (0.8)
User:	Oh from eight to eleven? Sure that's after my (.) uh: term paper is done so::=
IntelLLA:	=hmm=
User:	=I::m >hoping I could finish it by the deadline,< so: I will work on it all night, a:nd >I'll finish it by eight< so: (0.3) is that fine? Is your- when's your [paper due?
IntelLLA:	[Hmm, (4.2)
IntelLLA:	My paper is due at 9 in the morning. (0.7) But I need some time to finish it. (0.7)
User:	Hmmm right but you're not >goin' tuh be able to turn it in< by (0.9) [the deadline, right?
IntelLLA:	[hmm, (4.0)
IntelLLA:	I understand your concern. (0.6)
IntelLLA:	Could I have it a bit earlier (.) maybe at 7? (0.4)
User:	Yeah, sure I could finish- try to finish it by seven and .hhh maybe you could work on [that] fo::r two hours and then you could=
IntelLLA:	[Hmm,
User:	=hopefully turn it in by the deadline? (2.5)
IntelLLA:	Oh wow! (0.5) That would be great. (0.6) Thank you so much. ((Roleplay continues))

Table 3: CA notations used in the example transcript

CA notation	Description
Period (.)	Falling intonation
Comma (,)	Slightly rising intonation
Question (?)	Rising intonation
(ss.mm)	Silence in seconds
(.)	A brief silence (shorter than 250ms)
Colon (:)	Lengthening of previous sounds
Dash (-)	A cut-off of speech
Equal sign (=)	Latching (i.e., transition spaces minimized between turns)
Opening bracket ([)	Overlap onset
Closing bracket (])	Overlap offset (optional)
.hh	Inhalation
Hhh	Exhalation
Degree sign (°)	whispering; smaller voice
>WORD<	words pronounced at a faster pace
<WORD>	words pronounced at a slower pace

## B.2 Infit/Outfit statistics

In the context of educational measurement, a good assessment instrument should be able to "discriminate" among persons with different ability levels. One important assumption made by Rasch-family models is that score distributions from a good measurement instrument roughly follows a logistic regression with a slope of 1 (de Ayala, 2022). With such a strong assumption on the underlying patterns of data, it is impossible to obtain a perfect fit to the empirical data (Bond and Fox, 2013). Put differently, it is possible to obtain statistics on how well each constituent from each facet performs in relation to this model assumption. Two fit statistics (Infit and Outfit statistics) are commonly used to assess the amount of deviations of persons, raters, interlocutors, etc.

Outfit statistics is an unweighted average of squared standardized residuals (de Ayala, 2022; Bond and Fox, 2013). As such, it tends to emphasize the unexpected scoring patterns that are located far from the person's (or rater's) estimated scores. On the other hand, Infit statistics is a weighted average, which underscores misfit that are close to the persons' (, raters', or interlocutors') location estimates.

An ideal infit and outfit statistics is considered to be close to 1 (Bond and Fox, 2013). Infit/outfit statistics over 1.3 may indicate underfitting, suggesting some erratic scoring patterns. Infit/outfit statistics smaller than 0.7 may indicate overfitting

Table 4: Fit statistics for roleplay interlocutors.

Interlocutor	Outfit	Infit
AI	0.980	0.986
Tutor A	1.253	1.327
Tutor B	1.231	1.144
Tutor C	0.823	0.813
Tutor D	0.884	0.882
Tutor E	0.762	0.774

and too deterministic pattern of rating scores. As shown in Table 4, our system showed a good fit to the data according to both infit and outfit statistics. Some variations in misfit patterns were observed for individual human tutors. Two of them (A and B) slightly underfit (although mostly acceptable range) while the other three tutors tended to overfit (which was less problematic in this context).

## B.3 CEFR Descriptors

In this section, we introduce the descriptors we adopted for the rating described in Section 4. Table 5 shows the descriptors of the overall oral interaction defined by (Council of Europe, 2020).

Based on the definition of the interactional competence by (Galaczi and Taylor, 2018) describing "the ability to co-construct interaction in a purposeful and meaningful way, taking into account socio-cultural and pragmatic dimensions of the speech situation and event," Table 6 and 7 shows our extended descriptors of turn management and topic management respectively.

**Table 5: Overall Oral Interaction:** The ability to engage in spoken communication, managing and participating in conversations with fluency and spontaneity, while effectively responding to and understanding various contexts.

Level	Descriptor
C2	<ul style="list-style-type: none"> <li>- Can take part effortlessly in any conversation or discussion and have a good familiarity with idiomatic expressions and colloquialisms.</li> <li>- Can express fluently and convey finer shades of meaning precisely. If a problem arises, can backtrack and restructure around the difficulty so smoothly that other people are hardly aware of it.</li> </ul>
C1	<ul style="list-style-type: none"> <li>- Can express fluently and spontaneously without much obvious searching for expressions.</li> <li>- Can use language flexibly and effectively for social and professional purposes.</li> <li>- Can formulate ideas and opinions with precision and relate contributions skilfully to those of others.</li> </ul>
B2	<ul style="list-style-type: none"> <li>- Can interact with a degree of fluency and spontaneity that makes regular interaction with users of the target language quite possible.</li> <li>- Can take an active part in discussion in familiar contexts, accounting for and sustaining views.</li> </ul>
B1	<ul style="list-style-type: none"> <li>- Can deal with most situations likely to arise while travelling in an area where the language is spoken.</li> <li>- Can enter unprepared into conversation on topics that are familiar, of personal interest or pertinent to everyday life (e.g. family, hobbies, work, travel and current events).</li> </ul>
A2	<ul style="list-style-type: none"> <li>- Can communicate in simple and routine tasks requiring a simple and direct exchange of information on familiar topics and activities.</li> <li>- Can handle very short social exchanges, even though understanding enough to keep the conversation going oneself is not usually possible.</li> </ul>
A1	<ul style="list-style-type: none"> <li>- Can interact in a simple way provided the other person is prepared to repeat or rephrase things at a slower rate and help formulate what is being tried to express.</li> <li>- Can ask and answer simple questions in areas of immediate need or on very familiar topics.</li> </ul>

**Table 6: Turn Management:** The ability to keep the conversation cooperative and collaborative, in relation to the expected balance of contributions to the interaction among participants by means of socioculturally and pragmatically appropriate turn-taking.

Level	Descriptor
C2	<ul style="list-style-type: none"> <li>- Can interact with ease by (skillfully) interweaving his/her contributions into the conversation.</li> </ul>
C1	<ul style="list-style-type: none"> <li>- Can initiate, respond appropriately, and balance conversations, linking contributions to those of other speakers.</li> </ul>
B2	<ul style="list-style-type: none"> <li>- Can initiate discourse appropriately, actively invite the partner, take their turn when appropriate, and end conversation when they need to, though they may not always do this elegantly.</li> <li>- Can gain time and keep the turn while formulating what they want to express (e.g. "That's a difficult question to answer").</li> <li>- Can maintain and balance a natural and collaborative flow to the interaction (no long pauses within/between turns, no dominating interruptions).</li> <li>- Can make prompt and relevant responds appropriately, linking contributions to those of other speakers.</li> </ul>
B1	<ul style="list-style-type: none"> <li>- Can start up a conversation and help keep it going by asking people relatively spontaneous questions about a special experience or event, expressing reactions and opinions on familiar subjects.</li> <li>- Can intervene in a discussion on a familiar topic, using a suitable phrase to get the floor.</li> </ul>
A2	<ul style="list-style-type: none"> <li>- Can ask and answer questions about habits and routines, pastimes and past activities, and plans and intentions.</li> <li>- Can participate in short conversations in routine contexts on topics of interest.</li> </ul>
A1	<ul style="list-style-type: none"> <li>- Can ask and answer simple questions, initiate and respond to simple statements in areas of immediate need or on very familiar topics, including the factual information of themselves and other people (e.g. their home country, family, school).</li> </ul>

**Table 7: Topic Management:** The ability to develop ideas collaboratively, as opposed to extending their own speech, in relation to the communicative purpose and outcome and the topic of the interaction

Level	Descriptor
C2	<ul style="list-style-type: none"> <li>- Can advise on or discuss sensitive issues without awkwardness, understanding colloquial references, dealing diplomatically with disagreement and criticism.</li> <li>- Can link contributions skilfully to those of others, widen the scope of the interaction and help steer it towards an outcome.</li> </ul>
C1	<ul style="list-style-type: none"> <li>- Can develop others'/own ideas and relate own contribution skilfully to that of others.</li> </ul>
B2	<ul style="list-style-type: none"> <li>- Can take the initiative to introduce and contribute relevant new ideas in a discussion, extending the partner's thoughts and working towards joint decisions.</li> <li>- Can effectively summarize the discussion at key stages, evaluate the main points within their area of expertise, and propose the next steps to advance the interaction.</li> <li>- Can enhance the interaction by providing comments and asking questions that deepen collective understanding.</li> </ul>
B1	<ul style="list-style-type: none"> <li>- Can ask others to explain their ideas, give or seek personal views and opinions, and summarize the opinions or the points reached in an interaction.</li> <li>- Can help focus the argument and keep the development of ideas on course.</li> </ul>
A2	<ul style="list-style-type: none"> <li>- Can exchange what to do in the evening or at the weekend / what to do, where to go and make arrangements to meet.</li> </ul>
A1	<ul style="list-style-type: none"> <li>- Can exchange likes and dislikes for sports, foods, etc., using a limited repertoire of expressions, when addressed clearly, slowly and directly.</li> </ul>

# Curriculum-Driven Edubot: A Framework for Developing Language Learning Chatbots Through Synthesizing Conversational Data

Yu Li<sup>\*†</sup>, Shang Qu<sup>\*‡</sup>, Jili Shen<sup>§</sup>, Shangchao Min<sup>§</sup>, Zhou Yu<sup>†</sup>

<sup>†</sup>Columbia University      <sup>§</sup>Zhejiang University

<sup>‡</sup>University of Science and Technology of China

{yl15016, zy2461}@columbia.edu    qushang@mail.ustc.edu.cn

{22105040, msc}@zju.edu.cn

## Abstract

Chatbots have become popular in educational settings, revolutionizing how students interact with material and how teachers teach. We present Curriculum-Driven EduBot, a framework for developing a chatbot that combines the interactive features of chatbots with the systematic material of English textbooks to assist students in enhancing their conversational skills. We begin by extracting pertinent topics from textbooks and using large language models to generate dialogues related to these topics. We then fine-tune an open-source model using our generated conversational data to create our curriculum-driven chatbot. User studies demonstrate that EduBot outperforms ChatGPT in leading curriculum-based dialogues and adapting its dialogue to match the user’s English proficiency level. By combining traditional textbook methodologies with conversational AI, our approach offers learners an interactive tool that aligns with their curriculum and provides user-tailored conversation practice. This facilitates meaningful student-bot dialogues and enriches the overall learning experience within the curriculum’s pedagogical framework.

## 1 Introduction

The emergence of conversational agents has significantly impacted educational technology, changing how students interact with material and how teachers impart knowledge (Zhang and Aslan, 2021; Okonkwo and Ade-Ibijola, 2021; Cunningham-Nelson et al., 2019). These agents, more commonly known as “chatbots,” have shown usefulness in various educational settings, from teaching computer programming (Chinedu and Ade-Ibijola, 2021) to strengthening conversational skills (Li et al., 2022). However, its application comes with inherent challenges, especially in conversational skill development. Most chatbots primarily respond to user queries and follow the instruc-

tions provided. This approach contrasts with traditional language learning, which commonly follows a structured, textbook-based curriculum. As students progress through educational materials, they expect coherent and consistent content. Unfortunately, conventional chatbots may engage in generic conversations that include language or content unsuitable for a student’s level of proficiency, potentially impeding their learning progress.

To address these challenges, we propose a framework called Curriculum-Driven EduBot for developing a chatbot based on a specific curriculum. Our chatbot will focus on predetermined topics and use vocabulary from the curriculum to better align with the English proficiency of the users. It will act as a conversational practice partner, combining the interactive features of chatbots with the organized content of English textbooks. First, we extract relevant topics from textbooks and use large language models (LLMs) to synthesize fixed-format personas for both participants in the dialogue. Then, we use LLMs to synthesize dialogues based on these topics and personas, incorporating the vocabulary provided in the textbook. Subsequently, we fine-tune an open-source model with our generated conversational data to construct our chatbot. Our chatbot is more than just a responsive tool, it is an academic companion that guides students through coherent and friendly dialogues tailored to their English proficiency level. As illustrated in Figure 1, existing chatbots, such as ChatGPT, are not based on a curriculum. Instead of being conversational learning partners, they mainly act as AI-driven Q&A systems, and their content may not align with the student’s educational objectives. In contrast, our chatbot is constructed from synthesized dialogues that include clearly defined characters, curriculum-appropriate topics, and textbook-based vocabularies, thus providing an interactive and user-tailored conversational experience.

We conducted a thorough user study to evaluate

\* Both authors contributed equally to the work.

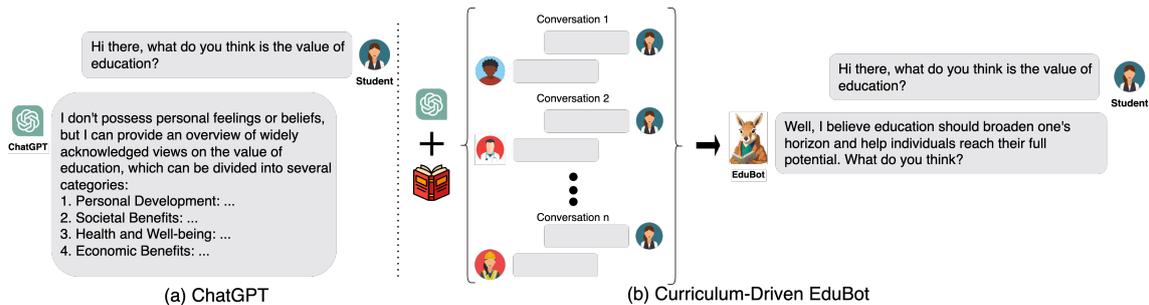


Figure 1: Comparison between ChatGPT vs. Our Curriculum-Driven EduBot. ChatGPT operates as an AI-powered Q&A tool, delivering comprehensive responses from a broad knowledge base. The Curriculum-Driven EduBot is fine-tuned with synthesized conversations, offering an interactive and adaptive learning experience through conversational practice.

Curriculum-Driven EduBot, using a high-quality college English textbook intended for English learners as a benchmark. Our findings indicate that our chatbot outperforms ChatGPT in various metrics. Specifically, 75% of students found EduBot to be particularly effective in facilitating interactive conversations, and they believed it was better suited to their English proficiency. The results and conversation examples from the user study clearly demonstrate that our chatbot is more closely aligned with the role of a language-learning companion. Furthermore, 83.3% of students expressed willingness to recommend EduBot to others, and 87.5% of students believe that interactions with EduBot can help students improve their conversational skills. In summary, our main contributions are as follows:

- We introduce a novel framework for curriculum-driven chatbots. Our approach involves synthesizing dialogues that incorporate fixed-format personas, curriculum topics, and relevant vocabularies. Subsequently, we fine-tune an open-source model to develop the chatbot, effectively integrating interactive chatbot features with structured educational content.
- We applied our framework to a specific curriculum. User studies reveal that EduBot outperforms ChatGPT. 87.5% of students believe that EduBot can help them improve their conversational skills.

## 2 Related Work

Many studies have shown that Artificial Intelligence (AI) can be utilized in educational settings (Chen et al., 2020b; Hinojo-Lucena et al., 2019; Chen et al., 2020a). For example, Rodrigues

and Oliveira (2014) created a formative assessment system capable of creating and assessing tests and tracking learners' progress. Similarly, Lan et al. (2014) proposed a machine learning-based approach to learning analytics, highlighting its potential to assess student knowledge. Recent advances in LLMs (Komeili et al., 2022; Shuster et al., 2022; OpenAI, 2023; Ouyang et al., 2022; et al., 2022) have had a major impact on the use of chatbots in educational settings (Cunningham-Nelson et al., 2019; Okonkwo and Ade-Ibijola, 2021; Kuhail et al., 2023). These conversational agents provide personalized learning experiences, engage learners, and help them retain knowledge. For example, Vasconcelos and dos Santos (2023) investigated the capabilities of ChatGPT<sup>1</sup> and Bing Chat<sup>2</sup> as resources that foster critical thinking and understanding of concepts to improve STEM education. Moreover, Li et al. (2022) used chatbots as conversational practice partners, providing learners with automatic grammar error feedback for language learning. Our research builds on these advancements by utilizing advanced open-source language models, enabling students to participate in discussions aligned with their curriculum.

Language learning, traditionally dependent on static resources such as textbooks and structured courses, has benefited greatly from curriculum-aligned approaches that combine consistency with adaptability. Krashen (1982) highlighted the importance of customized content delivery in language learning, suggesting that when learners engage with material that aligns with a structured curriculum, they often experience better comprehension and retention. Many researchers have ad-

<sup>1</sup><https://chat.openai.com>

<sup>2</sup><https://www.bing.com/new>

vocated systematically integrating curriculum content into new learning platforms to provide contextually relevant language exposure (Murphy et al., 2020; Clark, 2016; Andrade, 2014). For example, Rodríguez-Castro (2018) explored the potential of digital tools, such as virtual reality simulation, that map their content to official language learning curricula, ensuring that learners stay on track while taking advantage of interactive digital experiences. Furthermore, Ho et al. (2011); Holden and Sykes (2011) demonstrated the potential of curriculum-based gamification in language learning. Connecting game elements with curriculum milestones can motivate and engage learners longer. Qian et al. (2023) applied lexically constrained decoding to a dialog system, encouraging it to use curriculum-aligned words and phrases, resulting in better understanding and increased interest in practicing English. Our chatbot is the first to generate conversations from curricula and be trained on an open-source model.

The use of pre-trained language models (PLMs) (Roberts et al., 2019; Wang, 2021; et al., 2020; OpenAI, 2023; Zhang et al., 2022; Touvron et al., 2023; et al., 2023; Penedo et al., 2023) has enabled the generation of synthetic conversational data to enrich limited datasets, particularly in privacy-sensitive domains such as the medical domain (Varshney et al., 2023). Previous research has used PLMs to augment various conversational datasets (Chen et al., 2023a; Zheng et al., 2023a; Chen et al., 2022; Kim et al., 2022a; Chen et al., 2023b). For example, Zheng et al. (2023a) and Chen et al. (2022) used GPT-J (Wang, 2021) to generate responses tailored for emotional support dialogues and comprehension tasks, respectively. Kim et al. (2022b) proposed a collaborative human-AI paradigm in which a human operator and GPT-3 alternate in conversation. Chen et al. (2023a) generated dyadic and multiparty dialogues grounded on specific topic words, demonstrating outputs comparable to human-crafted ones. Our approach generates in-depth conversations based on educational curricula, allowing us to shape personas, focus on specific topics, and make lexical choices during data synthesis.

### 3 Method

We propose a framework for building a curriculum-based chatbot that can converse on topics derived from a given curriculum while aligning its re-

sponses to the user’s English proficiency level. As shown in Figure 2, our development process is divided into two parts. First, we use ChatGPT to generate simulated human-to-human dialogues based on textbook topics. Then, we fine-tune an open-source model to create our chatbot.

#### 3.1 Conversational Data Augmentation

The art of synthesizing human-human dialogues relies on two main factors: the topics being discussed and the personalities of the people involved in the conversation (Chen et al., 2023a; Kim et al., 2022a; Chen et al., 2022). To synthesize dialogues based on a curriculum, we propose a three-step approach. We start by extracting the main topics from the textbook and generating associated subtopics. Second, we develop a variety of personalities for the participants in the synthetic dialogues. Last, we synthesize dialogues based on the topics and personas obtained in the previous steps.

##### 3.1.1 Augment Topics

The range of topics covered in each curriculum unit is often limited. To broaden our synthetic dialogues to include a wide range of topics, we first extract the primary topics of the curriculum and then use ChatGPT to generate associated subtopics for each primary topic. For example, in our application, the primary topic of the first unit is “The True Value of Education”. We expand it to topics such as “The importance of education in personal and professional development” and “The role of education in promoting social justice and equity”. This process ensures that our dialogues are comprehensive and varied. The prompt given to ChatGPT in the augmentation process is detailed in Appendix A.1.1. Further information on this step and sample input-output pairs can be found in Appendix B.1.

##### 3.1.2 Create Personas

To enrich the conversational context, we also prompt ChatGPT to create personas for two dialogue participants: Person 1 and Person 2. These personas are crafted to reflect diverse backgrounds, including demographic characteristics (e.g., gender and race), socioeconomic status, cultural backgrounds, Myers-Briggs Type Indicator (MBTI) personality profiles, and personal experiences. Since the dialogue occurs between our chatbot and a student, and the model is trained to take on the role of Person 1 in the dialogue, we specify that Person 2’s background information consistently represents

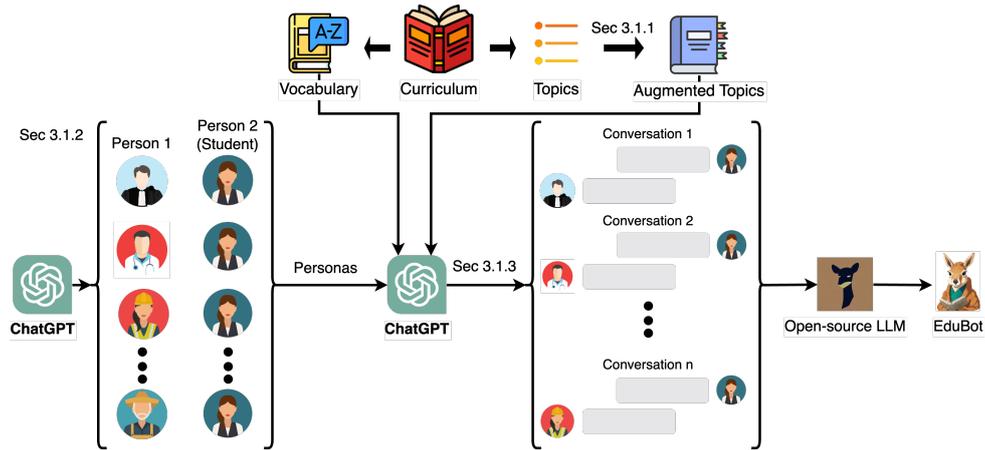


Figure 2: The initial step of the Curriculum-Driven EduBot Development is to enhance textbook topics (Sec.3.1.1). Following this, personas are created for synthetic conversation participants (Sec.3.1.2). Then dialogues are constructed based on vocabulary, topics and personas (Sec. 3.1.3). After this, an open-source model Vicuna is fine-tuned to get the EduBot ready for deployment (Sec. 3.2)

a typical student for the textbook we choose. In contrast, we randomly generate Person 1’s background information. Adopting this ‘fixed-random’ strategy offers two primary benefits: 1. It enables our chatbot to be trained with the student persona acting as the user and the alternate persona as the chatbot. Thus, the chatbot is ready to anticipate that its user is a student. 2. This encourages ChatGPT to generate conversations about topics commonly discussed by students, such as college life, which increases the chatbot’s appeal to users from this background. A detailed description of the prompts for this step can be found in Appendix A.1.2.

### 3.1.3 Compose Dialogues

We now instruct ChatGPT to generate synthetic dialogues using the generated personas and topics. To tailor the dialogue to the user’s English proficiency level and ensure that the dialogue aligns with the vocabulary that students are familiar with, we follow (Qian et al., 2023) and extract a subset of words from the vocabulary list of the relevant textbook unit to integrate into the conversation. We instruct ChatGPT to use a pair of personas generated in Step 2, one fixed as a student and the other with randomized characteristics. Participants with these personas will use the words in the vocabulary and converse on a topic chosen from our extended topic list in Step 1. To engage users and improve user experience, we also follow previous work and instruct the chatbot to actively lead the dialogue (University of California, 2019). Therefore, Person 1, representing the chatbot in the synthetic dialogues, is prompted to guide the dialogue. This

allows our chatbot to take the conversational lead with students, providing direction and guidance. The prompt given to ChatGPT in this step can be found in Appendix A.1.3.

## 3.2 Fine-Tuning An Open-Source Language Model With Synthesized Conversational Data

We use the synthesized dialogues to fine-tune an open-source large language model. Using open-source models offers several advantages: First, we can customize the underlying architecture and parameters according to our needs. In addition, we can synthesize additional data as required and improve the model through successive iterations. Last, open-source models are usually free, which significantly reduces costs.

We choose Vicuna-13B<sup>3</sup>, a cutting-edge open-source language model, for our specific application. We use it to build our chatbot since it possesses impressive language understanding capabilities. We fine-tune a single Vicuna-13B model using topics taken from all the units in the textbook. This approach ensures that our chatbot has a comprehensive knowledge base for all topics in the textbook. Following (Bao et al., 2023), during training, the chatbot takes on the role of Person 1, while the student takes on the role of Person 2. The prompt given to Vicuna during training can be found in Appendix A.2.

<sup>3</sup><https://huggingface.co/lmsys/vicuna-13b-v1.3>

### 3.3 Deployment of the Fine-tuned EduBot

Before using EduBot, students select a unit from the curriculum. We then assign a persona to the chatbot and choose a topic from the selected unit's topic list. Next, we randomly pick a set of words from the "New Words" vocabulary list of the unit. Finally, we use this information to create a specialized prompt for the fine-tuned EduBot. The implementation details can be found in Appendix A.3.

## 4 Experiments

### 4.1 Curriculum Source

In our evaluation, we selected the widely-used "New College English" (3rd edition) textbook, specifically the "Audiovisual Said Tutorial" from the third level, used in advanced English courses. This tutorial consists of eight units, each with a list of conversation topics. We generated ten associated topics for each main topic, as outlined in Section 3.1.1, and included ten randomly selected words from each unit's word list in the dialogues described in Section 3.1.3. This method produced 7,687 dialogues across the eight units for further development. Detailed statistics on our synthetic data are available in Appendix C.

### 4.2 Baseline

To evaluate our chatbot's performance, we use ChatGPT as our baseline because it generates meaningful, contextually appropriate responses. It has been trained on various datasets, allowing it to respond to diverse topics. We do not employ zero-shot prompted Vicuna as our baseline because it often fails to follow prompts, producing lengthy, hard-to-understand responses due to its smaller size and weaker instruction-following capability. Our fine-tuning process improves the Vicuna model and resolves this issue.

We observed that the length of responses has a significant impact on the user experience, in that some students prefer longer responses from the chatbot. This preference may be due to the text-based format of our chatbot. In comparison to speech-based chatbots, users may be more accepting of longer responses when using text-based chatbots because the repetition of information is less noticeable. However, long replies may hinder the development of conversational skills, as users might read the material and provide short responses rather than engage actively. To ensure fair assessment, we limit ChatGPT's response length using

the following prompt:

- As a social chatbot, please engage in a conversation about <Topic>. Share interesting anecdotes, facts, and experiences related to <Topic> Each response should be either one or two sentences. Please make all responses short and concise. Follow the above rules for all your utterances.

This prompt generally ensures concise replies from ChatGPT, though occasionally it produces lengthy responses, especially when users request detailed explanations.

### 4.3 Experimental Settings

#### 4.3.1 Participants

For our user study, we recruited 24 students from a renowned university in China via student discussion forums and in class. All participants had taken "College English 4," corresponding to the "New College English (3rd edition)" textbook, within the past year. To register, participants completed a background survey.

A total of 48 students completed the survey, among which 24 participated in the entire experiment and provided valid results. Of these students, 19 were in their second year, 4 in their third year, and 1 in their fourth year. Participants had an average age of 19.26 years, were from 20 majors, and had studied English for between 8 to 15 years (averaging 11.65 years). Their final grades for "College English 4" ranged from 2.1/5.0 to 5.0/5.0, averaging 4.06/5.0.

#### 4.3.2 Procedures

We conducted experiments in which participants were assigned either Unit 1 or Unit 2 of the textbook. Each participant had four conversations, two with EduBot and two with ChatGPT, each containing at least 20 utterances. To prevent bias, we randomly labeled bots A and B for each session and had participants converse first with Bot A and then with Bot B.

Participants completed a questionnaire immediately after interacting with the two chatbots. They first summarized each of their four conversations. The main questionnaire consisted of 20 criteria divided into six categories: Consistency with the Curriculum, English Proficiency Level, Role Identification, Quality of Conversation Language, Quality of

Conversation Content, and General Usefulness. For each criterion, participants chose whether Bot A was better, Bot B was better, or both were the same. All questions and instructions were in both Chinese and English, and participants could refer to their conversation records and textbook content. Each study took 20-30 minutes, and participants received \$5 compensation, adhering to China's minimum wage standards<sup>4</sup>. We excluded one submission for incorrect dialogue summaries and three for self-conflicting answers. Appendix F presents the user interface of our experiments, while the complete background survey and questionnaire are provided in Appendix E.1 and Appendix E.2.

## 5 Results and Discussion

The full results of the user study are shown in Table 1. We show the win rates for each questionnaire criterion. The results indicate that EduBot outperforms ChatGPT in several aspects.

**EduBot's language quality was on par with ChatGPT.** Similar percentages of participants preferred EduBot (29.2%) and ChatGPT (25.0%) regarding the coherence and fluency of the chatbots' utterances. This shows that EduBot produces responses of high language quality.

**EduBot offers a diverse range of relevant dialogue topics.** Through topic augmentation based on the curriculum, we aim for EduBot to center its conversations around topics that are relevant to but not directly listed in the textbook. Significantly more participants chose EduBot (50.0%) over ChatGPT (16.7%) when asked which chatbot mentioned topics and content that were not directly covered in the textbook and course. This shows that EduBot is capable of discussing diverse topics, compared with ChatGPT, which was only prompted with topics taken directly from the textbook.

At the same time, EduBot's conversation content remains in line with the curriculum. When asked which chatbot's conversation topics were more related to the course, student opinions were almost evenly divided. EduBot does not perform as well as ChatGPT in bringing up anecdotes, examples, questions, etc., related to the course. We believe that this is because ChatGPT gives longer statements that provide more material, while EduBot's

answers are more concise and concentrated on inquiring and engaging the user. This contrast is discussed in greater detail in Section 5.

**EduBot's conversations align better with students' English proficiency levels.** An equal percentage of participants (37.5%) chose EduBot and ChatGPT regarding which chatbot provided more vocabulary from their English course. We believe that this is because, without extra guidance, outputs produced by ChatGPT are generally close to the textbook in language difficulty. This makes it difficult to highlight EduBot's alignment with the students' English proficiency level. However, 37.5% of students found ChatGPT used many vocabulary words they did not understand, compared to 20.8% for EduBot. This shows that ChatGPT's conversations were sometimes too challenging for our target users. The varied English proficiency levels among participants led to mixed results in this section. We investigate the different preferences of students with varying English levels in Appendix G.

**EduBot's conversations are more natural and realistic.** Participants found their conversations with EduBot more natural and similar to real-life interactions. This distinction arose because, during the fine-tuning stage, EduBot has access to synthetic dialogues that emulated real-life conversations of Chinese college students. A higher percentage of students thought that EduBot was concise and accurate (50% vs. 12.5% for ChatGPT), natural and realistic (62.5% vs. 4.2% for ChatGPT). On the other hand, most participants found ChatGPT's responses too long and repetitive. Furthermore, results show that EduBot was better at guiding the conversation. 75% of students agreed that EduBot asked questions to guide the conversation, compared to only 4.2% for ChatGPT. Using EduBot, users found it easier to follow the dialogue without needing to introduce new topics to keep the conversation going.

**EduBot acknowledges the personas of both dialogue participants.** When conversing with EduBot, a larger proportion of participants felt that the chatbot was aware that they were Chinese college students (41.7%) compared to when they were talking to ChatGPT (29.2%). EduBot showed its knowledge of the user's identity by customizing its answers to the user's role. When it brought up common experiences of college students, participants

<sup>4</sup><https://take-profit.org/en/statistics/wages/china/>

Section	Question	EduBot (%)	ChatGPT (%)	Same (%)
Consistency With Curriculum	1. The main topics of my conversations with the chatbot were closely related to what I learned in English class.	41.7	50.0	8.3
	2. The chatbot brought up anecdotes, examples, questions, etc., related to what I learned in English class.	25.0	41.7	33.3
	3. The chatbot mentioned topics and content that were not directly covered in the textbook and course.	50.0	16.7	33.3
English Proficiency Level	1. During our conversations, the chatbot mentioned some vocabulary words that I learned in my English course.	37.5	37.5	25.0
	2. The chatbot used many vocabulary words that I didn't understand.	20.8	37.5	41.7
	3. I didn't find the conversations too easy to be helpful.	16.7	29.2	54.2
Role Identification	1. During conversations, I felt that the chatbot recognizes that I am a Chinese college student.	41.7	29.2	29.2
	2. During the two conversations with the chatbot, I felt like I was talking with two different people.	20.8	12.5	66.7
Language Quality	1. The utterances provided by the chatbot were coherent and fluent.	29.2	25.0	45.8
	2. The chatbot's responses were concise and accurate.	50.0	12.5	37.5
	3. Unlike in real everyday conversations, the chatbot's responses were long and redundant at times.	8.3	66.7	25.0
	4. Interactions with the bot were similar to natural, realistic conversations and not overly formal.	62.5	4.2	33.3
Content Quality	1. The chatbot acknowledged what I said and provided reasonable responses.	37.5	41.7	20.8
	2. The chatbot provided unique and personal perspectives regarding the selected topic.	45.8	37.5	16.7
	3. The chatbot used personal experiences to support its opinions.	33.3	16.7	50.0
	4. The chatbot actively raised questions to guide the course of the conversation.	75.0	4.2	20.8
	5. The chatbot didn't output offensive or hurtful responses.	0.0	8.3	91.7
General Usefulness	1. I would find it useful to use the chatbot to review what I learned in class.	16.7	25.0	58.3
	2. I would recommend the chatbot to other students.	37.5	16.7	45.8
	3. I believe that continuing to use the chatbot will help me improve my English conversation skills.	25.0	12.5	62.5

Table 1: Questionnaire Results (Percentage of Responses)

could easily build upon these topics and continue the conversation.

EduBot was also able to showcase its assigned personas. Slightly more participants felt like they were talking to two different people in their two conversations when using EduBot. EduBot performed better than ChatGPT in providing unique perspectives (45.8% compared to 37.5%) and personal experiences (33.3% compared to 16.7%). The incorporation of diverse personas makes EduBot's conversations more realistic and engaging. Examples of the above two characteristics are presented and further discussed in Section 6.

**EduBot is more useful for developing English conversation skills.** EduBot outperformed ChatGPT in two out of three criteria for general usefulness. A greater proportion of participants preferred

EduBot in terms of recommending it to other students and using it to enhance their conversational skills. On the contrary, slightly more participants thought that ChatGPT would be beneficial for reviewing what they had learned in class. As discussed in previous sections, the probable explanation is that ChatGPT presents relevant information in more length and detail.

In summary, the results show that EduBot offers diverse conversation topics, takes on specific personas, acknowledges the persona of its target user, and provides natural and realistic utterances. On the other hand, the advantage of ChatGPT mainly lies in its ability to provide elaborate content surrounding the conversation topic. Generally, EduBot is a more effective conversational partner that better cultivates the user's English conversation skills.

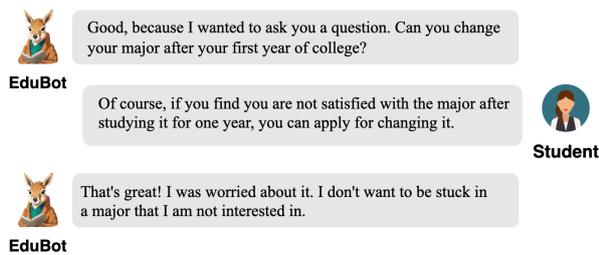


Figure 3: Case 1

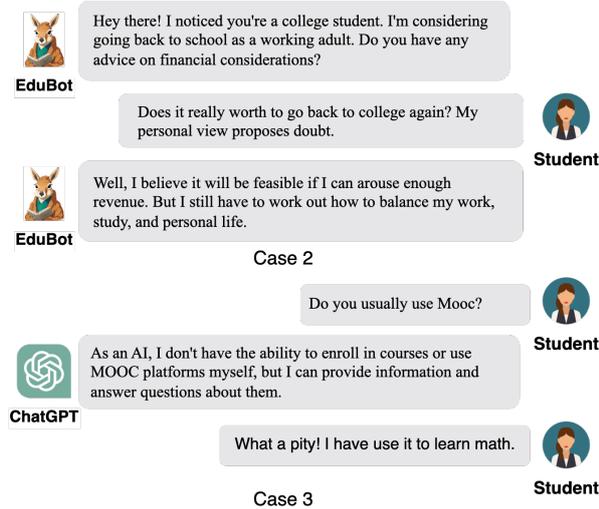


Figure 4: Case 2 from EduBot and Case 3 from ChatGPT-Bot.

## 6 Case Study

We present several dialogue snippets from our user study to illustrate the advantages and limitations of EduBot.

**EduBot recognizes the role of the user.** EduBot demonstrated the ability to recognize that the user was a college student and introduced familiar, relevant topics or questions. For instance, in Case 1 (Figure 3), EduBot asked about the user's university policy. By tailoring its questions to the user's background, EduBot created more personalized and meaningful interactions, which are crucial for maintaining engagement and ensuring effective learning.

**EduBot provides personal opinions and experiences.** EduBot formulated opinions and experiences consistent with the persona specified in the prompt, making conversations more realistic and engaging. In Case 2 (Figure 4), EduBot took on the persona of a working adult and provided personal experience on continuing education after starting work. In contrast, ChatGPT often did not provide realistic answers when asked about personal experiences,

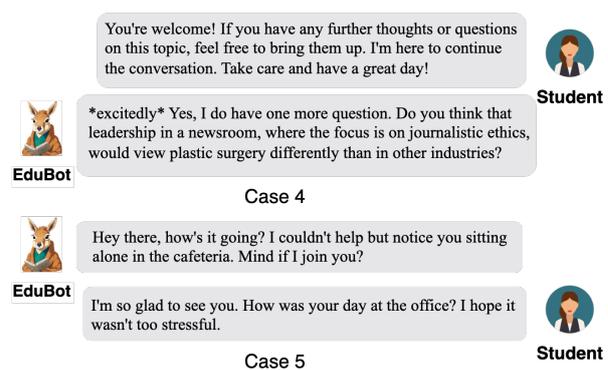


Figure 5: Case 4 and Case 5.

disrupting the natural flow of the conversation. For instance, in Case 3, ChatGPT struggled to offer a suitable response regarding its opinion on MOOC, an online learning platform.

**Limitations of EduBot.** We observed two phenomena that limited the quality of EduBot's conversations in several user study cases. First, EduBot occasionally included descriptions of its emotions or actions that should not appear in everyday conversations, as shown in Case 4 (Figure 5). Second, EduBot sometimes makes incorrect assumptions about the user's feelings or the conversation context. For example, in Case 5, EduBot hallucinated that the user was alone in the cafeteria. These issues stemmed from ChatGPT generating such scenarios in the data used to fine-tune EduBot. In the future, to address these issues, we plan to refine our data synthesis process and implement stricter post-processing methods to filter out unnatural content.

## 7 Conclusion and Future Work

In this work, we present Curriculum-Driven EduBot, a framework for developing a curriculum-based chatbot that combines the structured nature of English textbooks with the dynamic nature of chatbot interactions. We extract relevant topics from textbooks, then use LLMs to synthesize conversations around these topics. We fine-tune an open source model using these conversational data. Our user studies show that EduBot is more effective than ChatGPT in facilitating curriculum-related discussions, and is also able to adjust the chatbot to match the user's English proficiency. These results demonstrate EduBot's ability to provide a contextually appropriate conversational platform to develop conversation skills. In the future, there are opportunities to expand the content spectrum, incorporate multimedia elements, and introduce real-time

feedback mechanisms. As we incorporate these improvements, we hope to see EduBot evolve into an indispensable learning companion.

## References

- Maria de Lourdes Andrade. 2014. Role of technology in supporting english language learners in today's classrooms.
- Jianzhu Bao, Rui Wang, Yasheng Wang, Aixin Sun, Yitong Li, Fei Mi, and Ruifeng Xu. 2023. [A synthetic data generation framework for grounded dialogues](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10866–10882, Toronto, Canada. Association for Computational Linguistics.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. A survey on evaluation of large language models. *arXiv preprint arXiv:2307.03109*.
- Lijia Chen, Pingping Chen, and Zhijian Lin. 2020a. [Artificial intelligence in education: A review](#). *IEEE Access*, 8:75264–75278.
- Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Seokhwan Kim, Andy Rosenbaum, Yang Liu, Zhou Yu, and Dilek Hakkani-Tur. 2023a. [PLACES: Prompting language models for social conversation synthesis](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 844–868, Dubrovnik, Croatia. Association for Computational Linguistics.
- Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Andy Rosenbaum, Seokhwan Kim, Yang Liu, Zhou Yu, and Dilek Hakkani-Tür. 2022. [Weakly supervised data augmentation through prompting for dialogue understanding](#). In *NeurIPS 2022 Workshop on SyntheticData4ML*.
- Maximillian Chen, Xiao Yu, Weiyan Shi, Urvi Awasthi, and Zhou Yu. 2023b. [Controllable mixed-initiative dialogue generation through prompting](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 951–966, Toronto, Canada. Association for Computational Linguistics.
- Xieling Chen, Haoran Xie, Di Zou, and Gwo-Jen Hwang. 2020b. [Application and theory gaps during the rise of artificial intelligence in education](#). *Computers and Education: Artificial Intelligence*, 1:100002.
- Okonkwo Chinedu and Abejide Ade-Ibijola. 2021. Python-bot: A chatbot for teaching python programming. *Engineering Letters*, 29:25–34.
- Megan Clark. 2016. [The use of technology to support vocabulary development of english language learners](#). Samuel Cunningham-Nelson, Wageeh Boles, Luke Trouton, and Emily Margerison. 2019. *A Review of Chatbots in Education: Practical Steps Forward*, page 299–306. Engineers Australia, Brisbane, Queensland.
- Hugo Touvron et al. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Tom Brown et al. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Yuntao Bai et al. 2022. [Constitutional ai: Harmlessness from ai feedback](#). *Preprint*, arXiv:2212.08073.
- Francisco-Javier Hinojo-Lucena, Inmaculada Aznar-Díaz, María-Pilar Cáceres-Reche, and José-María Romero-Rodríguez. 2019. [Artificial intelligence in higher education: A bibliometric study on its impact in the scientific literature](#). *Education Sciences*, 9(1).
- Caroline M.L. Ho, Mark Evan Nelson, and Wolfgang Müeller-Wittig. 2011. [Design and implementation of a student-generated virtual museum in a language curriculum to enhance collaborative multimodal meaning-making](#). *Computers & Education*, 57(1):1083–1097.
- Christopher L Holden and Julie M Sykes. 2011. [Leveraging mobile games for place-based language learning](#). *International Journal of Game-Based Learning (IJGBL)*, 1(2):1–18.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2022a. [Soda: Million-scale dialogue distillation with social commonsense contextualization](#). *ArXiv*, abs/2212.10465.
- Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. 2022b. [ProsocialDialog: A prosocial backbone for conversational agents](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4005–4029, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2022. [Internet-augmented dialogue generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8460–8478, Dublin, Ireland. Association for Computational Linguistics.
- Stephen Krashen. 1982. Principles and practice in second language acquisition.
- Mohammad Amin Kuhail, Nazik Alturki, Salwa Alramlawi, and Kholood Alhejori. 2023. [Interacting with educational chatbots: A systematic review](#). *Education and Information Technologies*, 28(1):973–1018.

- Andrew S. Lan, Andrew E. Waters, Christoph Studer, and Richard G. Baraniuk. 2014. [Sparse factor analysis for learning and content analytics](#). *Journal of Machine Learning Research*, 15:1959 – 2008. Cited by: 96.
- Yu Li, Chun-Yen Chen, Dian Yu, Sam Davidson, Ryan Hou, Xun Yuan, Yinghua Tan, Derek Pham, and Zhou Yu. 2022. [Using chatbots to teach languages](#). In *Proceedings of the Ninth ACM Conference on Learning @ Scale, L@S '22*, page 451–455, New York, NY, USA. Association for Computing Machinery.
- V Murphy, Henriette Arndt, Jessica Briggs Baffoe-Djan, Hamish Chalmers, Ernesto Macaro, Heath Rose, Robert Vanderplank, and Robert Woore. 2020. Foreign language learning and its impact on wider academic outcomes: A rapid evidence assessment.
- Chinedu Wilfred Okonkwo and Abejide Ade-Ibijola. 2021. [Chatbots applications in education: A systematic review](#). *Computers and Education: Artificial Intelligence*, 2:100033.
- OpenAI. 2023. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- 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](#). *Preprint*, arXiv:2203.02155.
- 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*.
- Kun Qian, Ryan Shea, Yu Li, Luke Kutszik Fryer, and Zhou Yu. 2023. User adaptive language learning chatbots with a curriculum. In *Artificial Intelligence in Education. Posters and Late Breaking Results, Workshops and Tutorials, Industry and Innovation Tracks, Practitioners, Doctoral Consortium and Blue Sky*, pages 308–313, Cham. Springer Nature Switzerland.
- Adam Roberts, Colin Raffel, Katherine Lee, Michael Matena, Noam Shazeer, Peter J. Liu, Sharan Narang, Wei Li, and Yanqi Zhou. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. Technical report, Google.
- Fátima Rodrigues and Paulo Oliveira. 2014. [A system for formative assessment and monitoring of students' progress](#). *Computers & Education*, 76:30–41.
- Mónica Rodríguez-Castro. 2018. [An integrated curricular design for computer-assisted translation tools: developing technical expertise](#). *The Interpreter and Translator Trainer*, 12:1–20.
- Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju, Eric Michael Smith, Stephen Roller, Megan Ung, Moya Chen, Kushal Arora, Joshua Lane, et al. 2022. Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage. *arXiv preprint arXiv:2208.03188*.
- 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](#). *Preprint*, arXiv:2302.13971.
- Davis University of California. 2019. [Gunrock 2.0: A user adaptive social conversational system](#). In *Alexa Prize SocialBot Grand Challenge 3 Proceedings*.
- Deeksha Varshney, Aizan Zafar, Niranshu Kumar Behera, and Asif Ekbal. 2023. Knowledge grounded medical dialogue generation using augmented graphs. *Scientific Reports*, 13(1):3310.
- Marco Antonio Rodrigues Vasconcelos and Renato P. dos Santos. 2023. [Enhancing STEM learning with ChatGPT and Bing Chat as objects to think with: A case study](#). *Eurasia Journal of Mathematics, Science and Technology Education*, 19(7):em2296.
- Ben Wang. 2021. Mesh-Transformer-JAX: Model-Parallel Implementation of Transformer Language Model with JAX. <https://github.com/kingoflolz/mesh-transformer-jax>.
- 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](#). *Preprint*, arXiv:2303.04048.
- Ke Zhang and Ayse Begum Aslan. 2021. [AI technologies for education: Recent research & future directions](#). *Computers and Education: Artificial Intelligence*, 2:100025.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. [Opt: Open pre-trained transformer language models](#). *Preprint*, arXiv:2205.01068.
- Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng Zhang, and Minlie Huang. 2023a. [AugESC: Dialogue augmentation with large language models for emotional support conversation](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1552–1568, Toronto, Canada. Association for Computational Linguistics.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang,

Joseph E. Gonzalez, and Ion Stoica. 2023b. [Judging llm-as-a-judge with mt-bench and chatbot arena. Preprint](#), arXiv:2306.05685.

## A Prompts

### A.1 Data Augmentation Prompts

#### A.1.1 Augment Topics

- Given an input topic, generate a list of  $\langle n \rangle$  closely related topics that could be explored further.

Input topic:  $\langle \text{Topic} \rangle$

#### A.1.2 Create Personas

- Please provide me with one individual Person 1 with different backgrounds, including information about their demographic, socio-economic status, culture, MBTI personality type, and personal experiences, no need to show names. Then provide me with one individual Person 2 who is a  $\langle \text{student role information} \rangle$  but with different information.

We can substitute the  $\langle \text{student role information} \rangle$  with a comprehensive and detailed description of the students who actually use the textbook we select. More information about this step, along with an example of input and output, can be found in [B.2](#).

#### A.1.3 Compose Dialogues

- Generate a single conversation between these two people as Person 1 and Person 2 about the topic  $\langle \text{Topic} \rangle$ . Please take into account their distinct personalities and their backgrounds. Begin the conversation with Person 1.  
Please include the following keywords in Person 1's utterances:  $\langle \text{Vocab} \rangle$   
Person 1 should guide the conversation by asking more questions.

More details about this step, as well as examples of input and output, are provided in [B.3](#).

### A.2 Vicuna Prompt

We design the prompt structure for Vicuna as follows:

- As a social chatbot, please engage in a conversation while adopting the following personas:

$\langle \text{Person 1 Persona} \rangle$ .

Engage in a conversation about  $\langle \text{Topic} \rangle$  by showcasing your personas. Share interesting anecdotes, facts, and experiences related to  $\langle \text{Topic} \rangle$ . The English level of the conversation should be at CEFR  $\langle \text{English Proficiency Level of Textbook} \rangle$ .

To ensure that our bot is compatible with the English proficiency level of the textbook, we use The Common European Framework of Reference for Languages (CEFR) to control the difficulty level of language in our training process. CEFR is a widely used method to classify the difficulty level of texts. It defines six levels that represent increasing levels of difficulty or proficiency: A1, A2, B1, B2, C1 and C2. We include the CEFR level of the textbook in our system prompt. More information on our implementation can be found in [Appendix D](#).

### A.3 Fine-tuned EduBot Prompt

- As a social chatbot, please engage in a conversation while adopting the following personas:

$\langle \text{Persona} \rangle$ .

Engage in a conversation about  $\langle \text{topic} \rangle$  by showcasing your personas. Share interesting anecdotes, facts, and experiences related to  $\langle \text{Topic} \rangle$ . Include the following words in your utterances:  $\langle \text{Vocab} \rangle$ .

The English level of the conversation should be at CEFR  $\langle \text{Textbook English Proficiency} \rangle$ .

EduBot's deployment prompt is similar to its training prompt, except it incorporates the vocabulary of new words. This consistency helps EduBot generate responses related to the topic, draw on the persona's background information, and keep the language difficulty level consistent, similar to Person 1 in the synthetic dialogues. Additionally, by including new words, EduBot helps students revisit terms in the corresponding material.

### A.4 ChatGPT Prompt

- As a social chatbot, please engage in a conversation about  $\langle \text{Topic} \rangle$ .  
Share interesting anecdotes, facts,

and experiences related to <Topic>  
Each response should be either one or two sentences. Please make all responses short and concise. Follow the above rules for all your utterances.

We use the above prompt when using ChatGPT for our baseline model. We select a curriculum topic in the desired conversation unit and instruct ChatGPT to provide succinct responses. This guides ChatGPT's responses to be similar in length to those of EduBot, ensuring fair evaluation.

## B Data Augmentation Examples

### B.1 Example of Topic Augmentation Using ChatGPT

An example is presented below to demonstrate how ChatGPT is used to generate related subtopics from a main topic in the textbook.

- Primary Topic:  
The true value of education
- Prompt to ChatGPT:  
Given an input topic, generate a list of 10 closely related topics that could be explored further.  
Input topic: The true value of education
- Generated Sub-topics:  
“The importance of education in personal and professional development”,  
“The impact of education on social and economic mobility”,  
“The relationship between education and individual well-being”,  
“The role of education in promoting social justice and equity”,  
“The benefits of a well-rounded education that includes arts, humanities, and social sciences”,  
“The value of education in fostering critical thinking and problem-solving skills”,  
“The potential of education in promoting innovation and entrepreneurship”,  
“The relationship between education and global competitiveness”,

“The significance of lifelong learning in today’s rapidly changing world”,  
“The need for education reform to address societal challenges and prepare students for the future”

This example demonstrates the capability of ChatGPT to take a given topic and generate a variety of associated sub-topics that expand the breadth of our conversational database.

### B.2 Example of Synthesizing Personas of Participants Using ChatGPT

An example is presented below to demonstrate how ChatGPT is used to generate personas of participants in a conversation.

- Prompt to ChatGPT:  
Please provide me with one individual Person 1 with different backgrounds, including information about their demographic, socio-economic status, culture, MBTI personality type, and personal experiences, no need to show names. Then provide me with one individual Person 2 who is a college student originating from China but with different information.
- Generated Personas:  
Person 1:  
Gender: Male  
Demographic: African American  
Socio-economic status: Working class  
Culture: Baptist  
MBTI personality type: ENFP  
Personal experience: He grew up in a low-income neighborhood, and as the oldest child, he felt responsible for taking care of his siblings. He struggled with academics in high school but managed to graduate. He is now working as a bartender and uses his outgoing personality to make friends with his customers.  
  
Person 2:  
Gender: Female  
Demographic: Chinese  
Socio-economic status: Upper middle class  
Culture: Confucianism

MBTI personality type: INTP

Personal experience: She comes from a wealthy family and has had access to quality education throughout her life. She is currently studying computer science in college and hopes to start her own tech company after graduation. Growing up in a traditional Chinese household, she feels pressure from her parents to succeed academically and make the family proud. She struggles with balancing her ambitious career goals with her desire for independence and freedom.

This example demonstrates that ChatGPT can generate personas in a fixed format and follows our instructions that Person 2 is a college student with other random personas.

### B.3 Example of Synthesizing a Dialogue Using ChatGPT

We demonstrate the use of ChatGPT in our approach by providing an example. We use the personas in Appendix B.2 and the topic augmented in Appendix B.1 to create a dialogue. The words used are taken from the vocabulary related to the topic “The true value of education” in the textbook.

- Prompt to ChatGPT:

Generate a single conversation between these two people as Person 1 and Person 2 about the topic “The importance of education in personal and professional development”.

Please take into account their distinct personalities and their backgrounds. Begin the conversation with Person 1.

Please include the following keywords in Person 1’s utterances: “due”, “get down to”, “get away with”, “slam dunk”, “the easy way out”, “hand out”, “supposedly”, “revenue”, “industrial”, “rustle”

Person 1 should guide the conversation by asking more questions

- Generated Dialogue:

Person 1: Hey, have you ever thought about the significance of education in personal and

professional development? People sometimes believe that they can get away with not studying and still become successful, but that’s not always the case.

Person 2: Yes, I completely agree with you on that. Education paves the way for a successful career and personal growth. It provides the necessary skills and knowledge to make informed decisions and take on challenges in life.

Person 1: I see that you’re a computer science major. What made you choose this field?

Person 2: Well, I’ve always been interested in technology and innovation. I think computer science is an industry that is constantly growing and provides a lot of opportunities for revenue and growth.

Person 1: That’s interesting. Do you think attending college is a requirement for success in computer science, or do you think there are other routes to succeed?

Person 2: I think college is definitely a great way to learn about the industry and gain practical skills. However, there might be alternative ways of gaining knowledge and experience in the industry too. What do you think, Person 1?

Person 1: I believe education is important in any field, whether it’s through traditional schooling or hands-on experience. Sometimes it’s easier to take the easy way out and just hope someone will hand you an opportunity, but usually, you have to get down to work and rustle up some opportunities yourself. Sometimes people think there’s a magic formula to success, but there’s not a slam

dunk for anyone. You need to put in the effort and invest in yourself. What are your thoughts on that?

Person 2: I couldn't agree more. Education provides a solid foundation for professional development, but it's also important for personal growth. Knowing how to learn and adapt to change is essential in today's fast-paced and industrial world. Education is not just about getting a degree or a job, it's about being a lifelong learner and constantly upgrading yourself. What do you think about that?

Person 1: I think that's a smart way of seeing things! Education is supposed to be an investment in ourselves, not just for our future jobs or salaries. It gives us the tools to think critically, challenge ourselves, and grow in all areas of life. It's like learning a language - just like in school, you study grammar first before you can speak more fluently. It's the same with learning about other subjects. You're not going to jump into the field without basics. But once you have a foundation, you can build up from there.

This example illustrates how ChatGPT can generate conversations between Person 1 and Person 2 based on our instructions. We can observe that the dialogue revolves around the topic we specify, and most of the words we provide are used in the dialogue. Furthermore, both participants incorporate their individual experiences of their personas into the conversation.

## C Conversational Data Statistics

Using our chosen curriculum as the basis, we synthesized 880 to 1,210 dialogues per unit, averaging 1,058.76 dialogues for each. These dialogues comprise an average of 11.77 utterances, on average containing 28.71 words each. This section analyzes the statistical characteristics of our synthesized dialogues. To ensure the quality of our conversation data and its alignment with our objectives, we em-

ployed three attributes in our data synthesis process: curriculum topics, fixed-format personas, and relevant vocabularies. We first examine our generated personas for diversity and breadth in Sec. C.1. Then we evaluate the distribution of target words within dialogues in Sec. C.2. Moreover, to ascertain the congruence of our dialogues with the English proficiency standards of the textbook, we leveraged ChatGPT to assess word difficulty levels in both our synthesized dialogues and the curriculum in Sec. C.3.

### C.1 Persona Trait Distribution

As elaborated in Section 3.1.2, including conversation personas is important for ensuring diverse, engaging conversation content and styles. We first examine the range of personality traits represented in the generated personas. We use keyword string matching to extract the persona trait values from the generated persona descriptions. Figures 6 and 7 show the gender and MBTI personality type distributions of the personas, respectively. Synthetic dialogues include nearly equal proportions of both genders. The personality type distribution is not uniform, but all 16 types are represented in the synthetic dataset.

In addition, we verify the nationalities in the persona descriptions of Person 2. 8,000 of the total 8,470 persona descriptions explicitly specify "China" or "Chinese". This indicates that in most cases, ChatGPT successfully followed the additional instructions regarding Person 2, mentioned in Section 3.1.2.

### C.2 Target Word Distribution

During synthetic conversation generation, we included 10 target words in each prompt to be included in Person 1's utterances. Therefore, for each synthetic dialogue created, we compute the number of times the target words in the prompt are used in each dialogue turn. The first graph in Figure 8 displays the distribution of dialogues based on the total number of target words included by Person 1 and Person 2, respectively. Most of the words are included in Person 1's utterances, and in the majority of dialogues, Person 1 mentions at least half of the 10 vocabulary words. The second graph in Figure 8 shows the total number of vocabulary words included in each dialogue turn for each person.

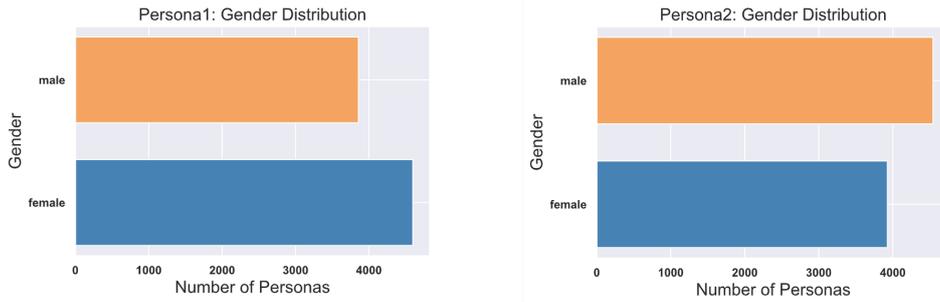


Figure 6: Distribution of gender in personas

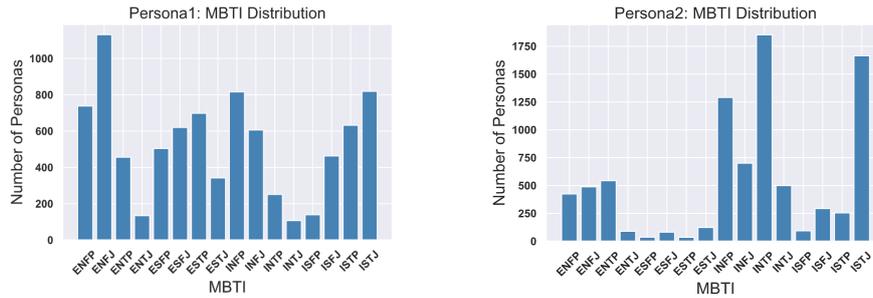


Figure 7: Distribution of MBTI personality types in personas

### C.3 English Proficiency Level

We evaluate whether the English proficiency level of the generated dialogues is similar to that of the curriculum. We use ChatGPT as an evaluator, as it has demonstrated its prowess in various language evaluation tasks Zheng et al. (2023b); Wang et al. (2023); Chang et al. (2023). We follow Zheng et al. (2023b) and utilize ChatGPT to automatically classify dialogues according to the CEFR scale using the following prompt:

- Evaluate the English proficiency of the given conversation according to the CEFR scale.

Provide one of the following six answers: A1, A2, B1, B2, C1, C2.

Output the CEFR level of the following conversation: <conversation>

<conversation> corresponds to the complete synthetic dialogue to be evaluated.

We then use the same method to evaluate the English proficiency level of “New College English” (3rd Edition), the original textbook we choose, by replacing the last sentence of the prompt with:

- Output the CEFR level of the following paragraph: <paragraph>

We assess each paragraph in the sample texts from “New College English”. The results of our evaluation for Unit 1 are shown in Figure 9. We found that

synthetic dialogues are comparable to those found in textbooks, yet they are slightly more challenging. This indicates that our method of synthesizing dialogues effectively ensures that our dialogues match the English proficiency level of the original textbook.

## D Implementation Details

To train a model for our application, we choose the 13B Vicuna model<sup>5</sup>. During the training phase, we carefully match each turn of our generated dialogues with the corresponding training turn in Vicuna format. As mentioned in Section 3.1.2, Person 1’s persona represents the chatbot’s side, while Person 2’s persona represents the students’. Therefore, we use utterances from Person 1 as the system’s responses and those from Person 2 as user requests throughout our training process. We train the Vicuna model for 3 epochs, beginning with a learning rate of  $2e-5$ . We use a batch size of 1 on each GPU and a gradient accumulation step of 16. We utilize 8 A100 GPUs and the training process takes three hours to complete.

<sup>5</sup><https://lmsys.org/blog/2023-03-30-vicuna/>

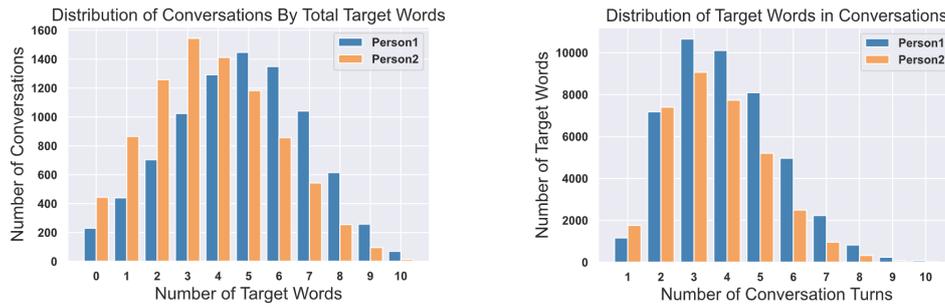


Figure 8: Distribution of target words

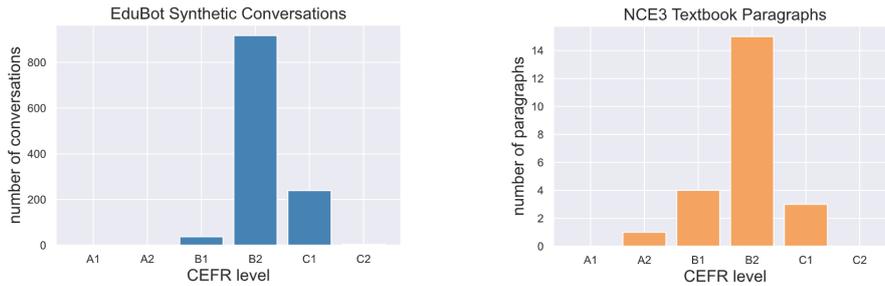


Figure 9: English proficiency levels of synthetic conversations and textbook paragraphs

## E Background Survey and Questionnaire

### E.1 Background Survey

Table 2 shows the full background survey we used for recruiting participants. “College English 4” uses the “New College English” (3rd edition) textbook and is a mandatory course for student participants of our user study. CET-4 and CET-6 are standardized English proficiency exams for Chinese college students.

Table 2: Background Survey for User Study Participants

Number	Question
1	Student ID
2	WeChat ID
3	Gender
4	Age
5	Grade
6	Major
7	Duration of English Learning
8	Overall Grade for <i>College English 4</i>
9	CET-4 Total Score
10	CET-4 Examination Date
11	CET-6 Total Score
12	CET-6 Examination Date
13	Available Time Slots

### E.2 Questionnaire

Table 3 presents the questionnaire we used to compare the quality of EduBot and ChatGPT from various aspects.

## F User Interface

We used the following user interface for both EduBot and ChatGPT. The user first selects a unit from the textbook (Figure 10) as the main topic of conversation, then proceeds to chat with the bot (Figure 11).

## G Analysis of Participants’ English Proficiency Levels

In this section, we analyze the influence of participants’ English proficiency levels on their perception of the two chatbots. We divided the participants into the following three groups according to their overall grade for the course “College English 4”: Group A consists of 8 students with scores between 2.1 and 3.6, Group B of 10 students with scores between 3.9 and 4.5, and Group C of 6 students with scores between 4.8 and 5.0. We reached the following conclusions.

Table 3: Questionnaire

Section	Number	Question
Participant Information	1	Student ID
Dialogue Summarization	2	Please summarize the main content of your first conversation with chatbot A.
	3	Please summarize the main content of your second conversation with chatbot A.
	4	Please summarize the main content of your first conversation with chatbot B.
	5	Please summarize the main content of your second conversation with chatbot B.
Consistency with Curriculum	6-1	The main topics of my conversations with the chatbot were closely related to what I learned in English class.
	6-2	The chatbot brought up anecdotes, examples, questions, etc., related to what I learned in English class.
	6-3	The chatbot mentioned topics and content that were not directly covered in the textbook and course.
English Proficiency Level	7-1	During our conversations, the chatbot mentioned some vocabulary words that I learned in my English course.
	7-2	The chatbot used many vocabulary words that I didn't understand.
	7-3	I didn't find the conversations too easy to be helpful.
Role Identification	8-1	During conversations, I felt that the chatbot recognizes that I am a Chinese college student.
	8-2	During the two conversations with the chatbot, I felt like I was talking with two different people.
Conversation Language Quality	9-1	The utterances provided by the chatbot were coherent and fluent.
	9-2	The chatbot's responses were concise and accurate.
	9-3	Unlike in real everyday conversations, the chatbot's responses were long and redundant at times.
	9-4	Interactions with the bot were similar to natural, realistic conversations and not overly formal.
Conversation Content Quality	10-1	The chatbot acknowledged what I said and provided reasonable responses.
	10-2	The chatbot provided unique and personal perspectives regarding the selected topic.
	10-3	The chatbot used personal experiences to support its opinions.
	10-4	The chatbot actively raised questions to guide the course of the conversation.
	10-5	The chatbot didn't output offensive or hurtful responses.
General Usefulness	11-1	I would find it useful to use the chatbot to review what I learned in class.
	11-2	I would recommend the chatbot to other students.
	11-3	I believe that continuing to use the chatbot will help me improve my English conversation skills.

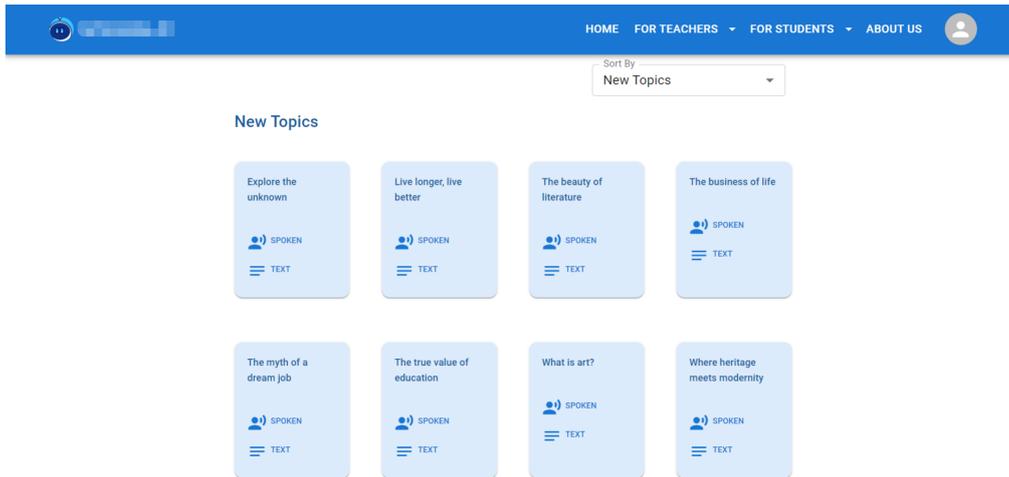


Figure 10: User Interface for Selecting a Textbook Unit as the Conversation Topic

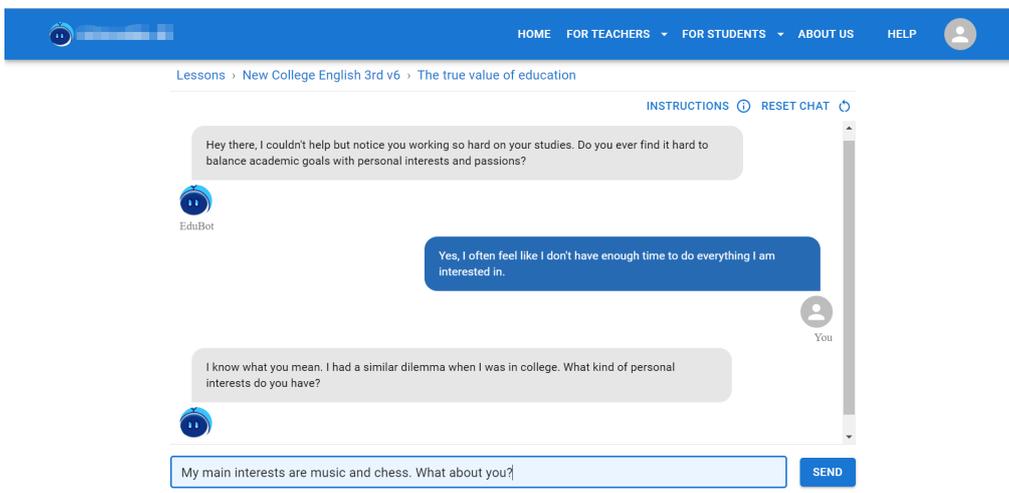


Figure 11: User Interface for Conversing with the Chatbots

### G.0.1 Participants with lower English proficiency levels found it more difficult to distinguish between the two chatbots.



Figure 12: Participants with lower English proficiency levels found it more difficult to distinguish between the two chatbots.

We observed that students in Group A were more likely to believe that the two chatbots performed the same over multiple questions. In addition, their responses were more often evenly split between the two chatbots. To verify, we calculated the following two statistics separately for each group of students: the average win rate of the “same” option over all questions and the average difference between win rates of “EduBot” and “ChatGPT” over all questions. The results are shown in Figure 12. We believe this is because it was harder for students in Group A to understand the chatbots and fully engage in the conversation.

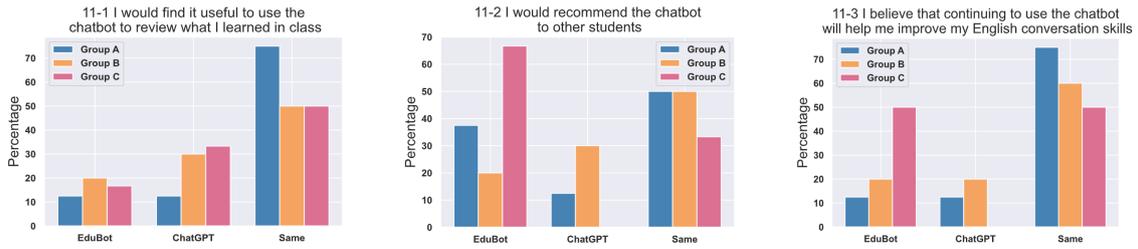


Figure 13: Participants with high English proficiency levels were more likely to prefer EduBot.

### G.0.2 Participants with high English proficiency levels were more likely to prefer EduBot.

In Figure 13, we present the three groups’ win rate results for the final section of the questionnaire. For the criteria “11-2 I would recommend the chatbot to other students” and “11-3 I believe that continuing to use the chatbot will help me improve my English conversation skills”, all participants in Group C chose either “EduBot” or “Same”. For “11-1 I would find it useful to use the chatbot to review what I learned in class”, results from Group C were in line with results from all the participants combined, with ChatGPT slightly outperforming EduBot. We believe that students in Group C more strongly preferred EduBot as a conversational training tool because they were more inclined to actively engage in conversations and provide their own thoughts instead of passively responding to the chatbot’s utterances. This caused EduBot’s advantages of providing natural responses and guiding the conversation by asking questions to be underscored in Group C’s results.

## H Analysis of User Study Conversations

We extracted all conversation histories from our user study. In the following section, we analyze the utterance lengths and coverage of target vocabulary words in the user study conversations.

### H.1 Utterance Lengths

As shown in Figure 14, we observe that in our user studies, ChatGPT produced longer outputs compared with EduBot. ChatGPT’s utterances were on average approximately 10 words longer than EduBot’s. In addition, ChatGPT occasionally produced outputs that were longer than 60 words, which rarely occurs in natural, daily conversations.

Furthermore, Figures 15 and 16 demonstrate that user study participants generally provided longer

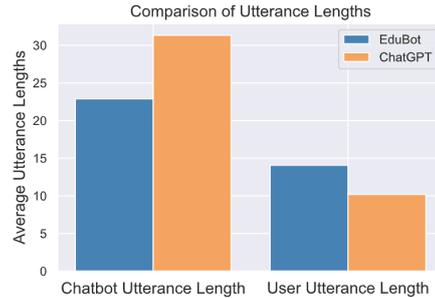


Figure 14: Comparison of utterance lengths in EduBot and ChatGPT conversations in the user study

responses when conversing with EduBot compared to ChatGPT. This indicates that EduBot’s more interactive and realistic conversation style better engages the users and guides them to practice their own conversation skills.

### H.2 Target Vocabulary Words

We also assess if EduBot can incorporate words from the target vocabulary. As shown in Figure 17, on average, conversations with EduBot included 5.55 words from the target vocabulary, while conversations with ChatGPT only included 0.62. This demonstrates that EduBot, which was further refined using curriculum-aligned data, is better suited to the user’s curriculum and English level.

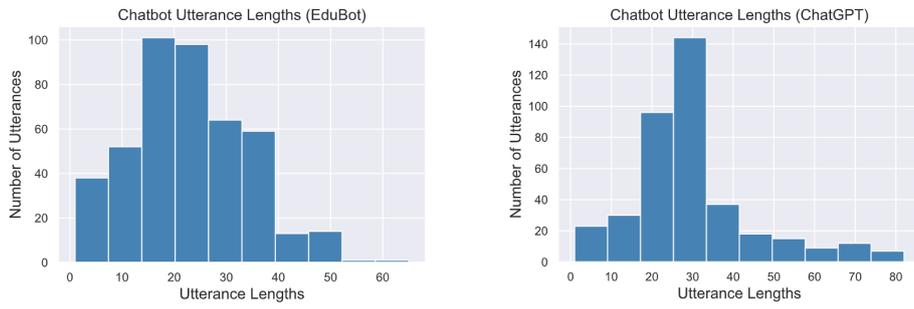


Figure 15: Lengths of chatbot utterances in the user study

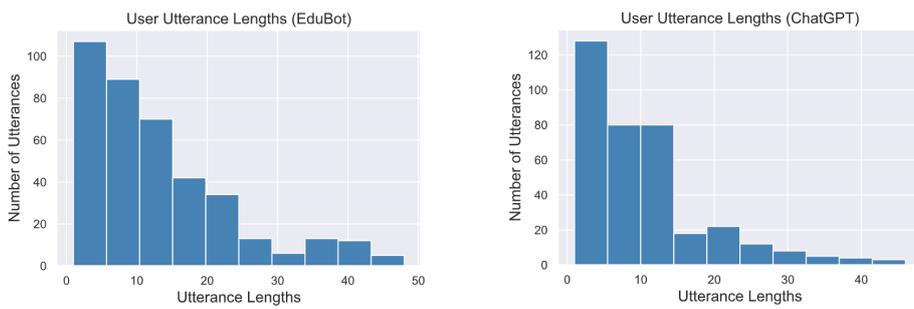


Figure 16: Lengths of user utterances in the user study

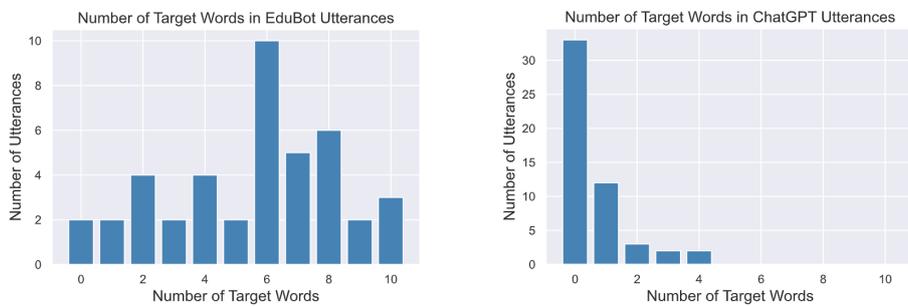


Figure 17: Coverage of target words in user study conversations

# Going beyond Imagination! Enhancing Multi-modal Dialogue Agents with Synthetic Visual Descriptions

Haolan Zhan, Sameen Maruf\*, Ingrid Zukerman and Gholamreza Haffari

Department of Data Science & AI, Monash University, Australia  
{haolan.zhan, ingrid.zukerman, gholamreza.haffari}@monash.edu

## Abstract

Building a dialogue agent that can seamlessly interact with humans, in multi-modal regimes, requires two fundamental abilities: (1) understanding emotion and dialogue acts within situated user scenarios, and (2) grounding perceived visual cues to dialogue contexts. However, recent works have uncovered shortcomings of existing dialogue agents in understanding emotions and dialogue acts, and in grounding visual cues effectively. In this work, we investigate whether additional dialogue data with only visual descriptions can help dialogue agents effectively align visual and textual features, and enhance the ability of dialogue agents to ground perceived visual cues to dialogue contexts. To this end, in the absence of a suitable dataset, we propose a synthetic visual description generation pipeline, and contribute a large-scale synthetic visual description dataset. In addition, we propose a general training procedure for effectively leveraging these synthetic data. We conduct comprehensive analyses to evaluate the impact of synthetic data on two benchmarks: MELD and IEMO-CAP. Our findings suggest that synthetic visual descriptions can serve as an effective way to enhance a dialogue agents' grounding ability, and that the training scheme affects the extent to which these descriptions improve the agent's performance.

## 1 Introduction

There have been impressive advances in large-scale vision and language models (VLMs) in performing multi-modal tasks, such as visual question answering (VQA) and image captioning (Guo et al., 2023; Chen et al., 2022; Liu and Chen, 2024). While VLMs are powerful general-purpose models for a wide range of tasks, most state-of-the-art VLMs still struggle with providing real-world, situated multi-modal assistance (Wu et al., 2023, 2024).

\*Work was done when Sameen was at Monash.



Figure 1: Visual descriptions can be an effective way to help dialogue agents interpret the visual cues from images, further enhancing the understanding ability towards human emotion and dialogue acts.

Building a situated dialogue agent that can seamlessly interact with humans in a multi-modal scenario requires two essential abilities: (1) understanding the interlocutor's emotion and dialogue acts within situated user scenarios, and (2) grounding perceived visual cues to dialogue contexts.

However, recent work (Wu et al., 2023; Liu et al., 2023; Xenos et al., 2024) has unveiled shortcomings of existing VLM-based dialogue agents with respect to these abilities. We hypothesize that current limitations can be attributed to the gap between different modalities, also known as the misalignment between visual and textual features. We argue that visual descriptions can serve as a potential way to bridge this gap by interpreting visual cues from images. To verify our hypothesis, we propose to investigate whether additional dialogue data with only visual descriptions can help dialogue agents effectively align visual and textual features, and enhance their ability to ground perceived visual cues to dialogue contexts. For instance, looking at the images in Figure 1, *visual descriptions* are capable of conveying subtle but important visual

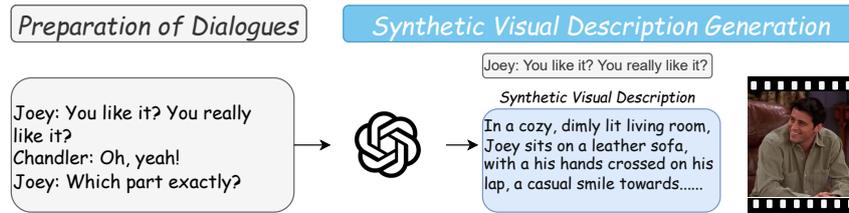


Figure 2: Synthetic data generation pipeline. Please note that the image on the right is provided for reference only, to aid in understanding the generated visual description, and is not produced by ChatGPT.

cues (e.g., *facial expression*, *human position*) about the people in these images.

Given the absence of datasets that offer annotations for visual descriptions, we devise a novel synthetic visual description generation pipeline using ChatGPT and contribute a large-scale synthetic visual description dataset by extending existing multi-dialogue corpora with additional visual descriptions. Furthermore, to effectively utilize these synthetic data, we explore several training schemes based on *knowledge distillation (KD)* (Hinton et al., 2015). Those training schemes aim to instruct dialogue agents to align the features in different modalities by distilling the ability to interpret visual cues learnt from the synthetic data.

We conduct a comprehensive analysis where we evaluate the effectiveness of synthetic data on two benchmarks: *MELD* (Poria et al., 2018) and *IEMOCAP* (Busso et al., 2008). Our results show that synthetic visual descriptions play an important role in helping dialogue agents understand and ground visual cues from images to dialogue contexts. Specifically, our method outperforms the baseline VLMs (e.g., LLaVa-1.5) by at least 6% on both emotion detection and dialogue act classification tasks. Moreover, the remarkable performance of our training framework based on knowledge distillation demonstrates that the training scheme affects the extent to which these descriptions improve a dialogue agent’s performance.

## 2 Synthetic Data Generation

To appropriately understand an interlocutor’s emotions and dialogue acts, VLM-based dialogue agents must ground perceived visual cues within dialogue contexts. We hypothesize that current VLMs are limited by a large gap between different modalities, which affects their ability to *ground* visual cues to dialogue contexts. However, the annotation that links visual cues and dialogue contexts is missing from existing widely-used datasets (e.g.,

*MELD* (Poria et al., 2018)). In this section, we investigate whether additional *visual descriptions* can help dialogue agents, and to what extent synthetic data can be used to bridge the gap between visual cues and dialogue contexts.

We can easily have access to large amounts of dialogue contexts, but it is hard to obtain the corresponding images or videos. In the absence of suitable multimodal datasets with the grounding annotations, we propose a visual description generation approach in the rest of this section. We then propose a training procedure (§ 3) for leveraging this synthetic data to improve the performance of multi-modal dialogue agents on the tasks of dialogue act and emotion prediction.

**Problem Formulation and Notation.** Given a multi-turn dialogue  $d = \{u_1, \dots, u_m\}$  consisting of a sequence of utterances  $u_i$  in plain text, our goal is to prompt ChatGPT to generate synthetic visual description  $v'_i$  for each utterance. We will get a synthetic dataset in which each of the utterances is paired with a synthetic visual description instead of a image. This synthetic dataset, augmented with the visual descriptions, will be used in training for reducing the modality gap, as explained below; see Figure 2 for an example.

**Synthetic Visual Description Generation.** Multi-modal dialogue tasks utilize plain-text dialogue and visual cues simultaneously. The motivation for the synthetic visual description generation is to explore if we can leverage it instead of real images to improve the performance of multi-modal dialogue agents. The idea is that these descriptions will stimulate a VLM-based dialogue agent to *imagine* potential visual scenes. We select three main factors that can affect visual scenes on the task of dialogue act and emotion prediction, viz (1) facial expression, (2) human action, and (3) human position; and incorporate them into synthetic visual description. We then

prompt ChatGPT via in-context learning (ICL) to generate a potential synthetic visual description for each utterance, as shown in Figure 2. We provide an example in the Appendix B to better understand the synthetic generation process.

From the MELD and IEMOCAP datasets, we have extracted and prepared 6,357 multi-turn dialogues, comprising a total of 22,126 utterances. Each utterance is associated with a synthetic visual description that depicts the potential visual scene associated with the dialogue context. The average length of each synthetic visual description is 15.6 words.

### 3 Fine-tuning a VLM-based Agent using Synthetic Data

In this section, we propose a methodology for leveraging the synthetic data produced as explained in Section 2. One intuitive way is to combine synthetic data with real data for training. However, as there is a large gap in both modalities and patterns between real images and synthetic visual descriptions, a straightforward concatenation of real and synthetic dataset would not be the best choice. We therefore propose a multi-stage training framework, which trains the dialogue agent with synthetic and real data separately, followed by a knowledge distillation training stage (Hinton et al., 2015). Specifically, we choose the state-of-the-art VLM model, LLaVa-v1.5 (Liu et al., 2023) as the backbone of our system, which integrates the visual encoder of CLIP (Radford et al., 2021) with the language decoder Vicuna (Chiang et al., 2023).

**Fine-tuning with Synthetic Data.** Suppose we have a synthetic training dataset of dialogues  $\mathcal{D}_s$ , where each dialogue  $d' = \{(u'_1, v'_1, y'_1), \dots, (u'_m, v'_m, y'_m)\}$  contains  $m$  utterances ( $u'$ ), associated synthetic visual descriptions ( $v'$ ) and output labels ( $y'$ ). We use the synthetic training dataset to fine-tune the LLaVa-v1.5 model with LoRA adapter (Hu et al., 2021), denoted by  $\theta_s$ . As the dialogues and synthetic visual descriptions are both in text, instead of feeding images to the CLIP module, we only need to use the Vicuna module as the proxy to encode the synthetic descriptions for the visual encoding.

**Fine-tuning with Real Data.** The goal of fine-tuning with real data is to adapt the dialogue model to the real multi-modal situation. We have a real dataset  $\mathcal{D}_r$  containing a set of multi-

modal dialogues, where each dialogue  $d = \{(u_1, v_1, y_1), \dots, (u_m, v_m, y_m)\}$  has  $m$  utterances ( $u$ ), corresponding images ( $v$ ) and output labels ( $y$ ). Unlike synthetic data fine-tuning, the CLIP and Vicuna module within the LLaVa-v1.5 will be used to process visual images and dialogue contexts collaboratively. This process will yield a fine-tuned adapter  $\theta_r$  for the real data.

**Knowledge Distillation.** The distillation training procedure aims to transfer the “*imagination*” ability learnt from the synthetic data to enhance dialogue agents in grounding visual cues to dialogue contexts in multi-modal settings. We conduct the knowledge distillation procedure on the fine-tuned adapters  $\theta_s$  and  $\theta_r$  by applying the KL-divergence (Kullback and Leibler, 1951) regularization in three different settings, as follows.

- *Synthetic distillation* ( $s \rightarrow r$ ): Knowledge is distilled from the synthetic adapter  $\theta_s$  to the real adapter  $\theta_r$ , based on the following training objective:

$$\max_{\theta_r} \sum_{d \in \mathcal{D}_r} \sum_{(u,v,y) \in d} \log P_{\theta_r}(y|u, v) - \gamma KL(P_{\theta_r}(\cdot|v, y) || P_{\theta_s}(\cdot|v, y))$$

where  $\log P_{\theta_r}(y|u, v)$  refers to log-likelihood probability of generated label  $y$  from the model with real adapter  $\theta_r$ . Besides, the distillation function  $KL(\cdot || \cdot)$  aims to measure and minimize the difference between  $\theta_r$  and  $\theta_s$ .  $\gamma$  is the regularisation coefficient to control the trade-off between two objectives.

- *Real distillation* ( $r \rightarrow s$ ): Knowledge is distilled from the real adapter  $\theta_r$  to the synthetic adapter  $\theta_s$ , based on the following training objective:

$$\max_{\theta_s} \sum_{d' \in \mathcal{D}_s} \sum_{(u',v',y') \in d'} \log P_{\theta_s}(y'|u', v') - \gamma KL(P_{\theta_s}(\cdot|v', y') || P_{\theta_r}(\cdot|v', y'))$$

- *Mutual distillation* ( $s \leftrightarrow r$ ): This is a mutual KD between two adapters,

$$\max_{\theta_r} \sum_{d \in \mathcal{D}_r} \sum_{(u,v,y) \in d} \log P_{\theta_r}(y|u, v) - \gamma_1 KL(P_{\theta_r}(\cdot|v, y) || P_{\theta_s}(\cdot|v, y)) - \gamma_2 KL(P_{\theta_s}(\cdot|v, y) || P_{\theta_r}(\cdot|v, y))$$

Dataset	MELD		IEMOCAP	
	Emo.	DA	Emo.	DA
UniVL	66.37	61.47	54.91	61.19
MiniGPT-4	78.00	70.33	69.00	68.93
Video-LLaMa	72.38	68.42	63.16	65.75
MultiModal-GPT	73.54	68.01	61.27	64.92
LLaVa-1.5	79.26	76.39	66.03	71.48
<i>ours</i>	<b>87.38*</b>	<b>81.03*</b>	<b>73.89*</b>	<b>77.29*</b>

Table 1: Accuracy (%) of VLM-based multi-modal dialogue agents on the emotion (**Emo.**) and dialogue act (**DA**) understanding tasks. "\*" indicates a significance of p-value < 0.05 in the Chi-Square test after Benjamini-Hochberg (BH) correction for false discovery rate (Benjamini and Hochberg, 1995).

where  $\gamma_1$  and  $\gamma_2$  are regularisation coefficients to balance the effects of two types of distillation between  $\theta_r$  and  $\theta_s$ .

## 4 Experiments

In our experiments, we aim to investigate the following two research questions: (1) How effectively do existing VLM-based agents comprehend emotions and dialogue acts from humans?, and (2) To what extent can synthetic visual descriptions enhance agents' multi-modal capabilities in understanding emotions and dialogue acts.

**Settings.** Our experiments were conducted on two datasets: *MELD* (Poria et al., 2018) and *IEMOCAP* (Busso et al., 2008), both of which are rich in annotations of emotion and dialogue acts. We evaluate the performance of each model by reporting its accuracy in predicting emotion and dialogue acts. In terms of VLMs, we select several state-of-the-art baselines including **UniVL** (Luo et al., 2021), **MiniGPT-4** (Zhu et al., 2023), **Video-LLaMa** (Zhang et al., 2023), **MultiModal-GPT** (Gong et al., 2023) and **LLaVa** (Liu et al., 2023). The details of each baseline can be found in Appendix A.

**Performance of Existing VLMs.** Table 1 presents the accuracy (%) of existing VLM-based agents on the emotion and dialogue act understanding tasks. We observe that LLaVa-1.5 outperforms other VLMs to a large extent in the MELD dataset and maintains competitive performance with MiniGPT-4 on the IEMOCAP dataset. We also note that existing VLMs mainly rely on their LLM module (e.g., Vicuna module in the LLaVa-1.5 agent), but they struggle to merge the information extracted from the CLIP (visual) module with the

Dataset	MELD		IEMOCAP	
	Emo.	DA	Emo.	DA
LLaVa-1.5 (vanilla)	79.26	76.39	66.03	71.48
(1) <i>synthetic data (s)</i>	75.84	68.16	59.98	66.92
(2) <i>real data (r)</i>	82.67	78.25	69.72	72.94
(3) <i>mixture (s then r)</i>	84.01	79.86	71.35	74.16
(4) <i>mixture (r then s)</i>	76.15	71.55	62.04	68.28
(5) <i>distillation (s → r)</i>	<b>87.38*</b>	<b>81.03*</b>	73.89	<b>77.29</b>
(6) <i>distillation (r → s)</i>	80.19	77.91	65.63	71.74
(7) <i>distillation (r ↔ s)</i>	85.43	79.59	<b>74.52</b>	76.13

Table 2: Ablation studies of different types of training data and distillation settings. "\*" indicates a significance of p-value < 0.05 in the Chi-Square test with BH correction.

LLM (textual) module, mainly due to the modality misalignment. The results support our hypothesis that visual descriptions can help bridge the gap by interpreting visual cues from images. We further provide an in-depth analysis of the impact of visual descriptions.

**The Effectiveness of Synthetic Data.** We conducted comprehensive ablation studies to investigate the effectiveness of using synthetic data to enhance the performance of our agent. We selected the top-performing VLM model, LLaVa-1.5, from Table 1 as our baseline. The results are presented in Table 2, which outlines seven different data configurations, including: (1) training only on *synthetic* data, (2) training only on *real* data, (3) mixed training involving initial training on *synthetic* data followed by *real* data, (4) mixed training involving the reverse sequence, and employing distillation techniques as discussed in Section 3, viz (5) synthetic distillation (synthetic→real), (6) real distillation (real→synthetic) and (7) mutual distillation (real↔synthetic).

The findings in Table 2 indicate that incorporating knowledge distillation into the training process enables LLaVa-1.5 to surpass the performance achieved through either naive mixed training or training solely on synthetic data or on real data. Notably, among the three distillation approaches ((5)-(7)), the strategy of distilling knowledge from synthetic to real data (*distillation (s → r)*) yielded the best results overall. In contrast, the performance of distillation from real to synthetic data was largely equivalent to that of LLaVa-1.5. This suggests that synthetic data must be utilized judiciously, as a significant discrepancy between real and synthetic data can adversely affect performance.

## 5 Related Work

**Visual Dialogue.** The visual dialogue task was proposed by Das et al. (2017). It requires an agent to answer multi-round questions about a given image, similarly to Visual Question Answering (VQA) (Das et al., 2017; Jiang et al., 2020b; Huber et al., 2018). Previous work (Wu et al., 2018; Kottur et al., 2018; Yang et al., 2019; Guo et al., 2019; Niu et al., 2019; Kang et al., 2019; Jiang et al., 2020a; Yang et al., 2021) focused on developing different attention mechanisms to model the interactions among image, question and dialogue history (Wang et al., 2020). With the rapid development of large-scale vision-language models (VLMs) (Chen et al., 2022; Dai et al., 2022; Wu et al., 2023; Zhu et al., 2023; Zhang et al., 2023), recent work focuses on building unified models that can handle multiple tasks. However, most models are still unable to support situated interaction with humans in real scenarios, especially capturing human emotions and dialogue acts, and grounding to their dialogue contexts.

**Learning from synthetic data.** There has been some work on learning from synthetic data for dialogue systems (Dai et al., 2022; Kim et al., 2022; Semnani et al., 2023; Bao et al., 2023; Abdullin et al., 2024; Zhan et al., 2024). Synthetic data are easy to generate, and are particularly useful for providing detailed labelling to reduce human labor, such as dialogue acts (Zhan et al., 2023), knowledge injection (Bao et al., 2023) or simulating dialogues in new scenarios, such as the rapid generation of a sequence of QA from documents (Dai et al., 2022). However, these works mainly focus on plain text dialogues, rather than multi-modal dialogues. We propose a novel framework to utilize synthetic data to address this gap, and thereby enhance the abilities of multi-modal dialogue agents on the task of emotion and dialogue act classification.

## 6 Conclusion

Our work demonstrates the potential of synthetic visual descriptions to improve the performance of dialogue agents, particularly in understanding emotions, dialogue acts and grounding visual cues to dialogue contexts. By introducing a novel synthetic visual description generation pipeline and a large-scale dataset, along with an effective training procedure, we have taken a crucial step towards overcoming the limitations of multi-modal dialogue agents.

The positive outcomes observed in our experiments highlight the importance of appropriate training schemes to fully leverage synthetic data, pointing towards a promising direction for future research.

## Limitations

As our work provides an initial step into incorporating synthetic visual descriptions into multimodal dialogue agents, we do not offer an exhaustive analysis of the synthetic data, nor do we identify the most suitable use cases for evaluating the effectiveness of synthetic data in such scenarios. Besides, we did not analyse why certain distillation schemes do better than others. Additionally, it is promising to conduct further evaluation to determine whether enhancing the agents’ grounding capabilities could also improve their response abilities.

## Acknowledgement

This material is based on research sponsored by DARPA under agreement number HR001122C0029. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright notation thereon.

## References

- Yelaman Abdullin, Diego Molla-Aliod, Bahadorreza Ofoghi, John Yearwood, and Qingyang Li. 2024. Synthetic dialogue dataset generation using llm agents. *arXiv preprint arXiv:2401.17461*.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736.
- Jianzhu Bao, Rui Wang, Yasheng Wang, Aixin Sun, Yitong Li, Fei Mi, and Ruifeng Xu. 2023. [A synthetic data generation framework for grounded dialogues](#). In *ACL*, pages 10866–10882, Toronto, Canada. ACL.
- Yoav Benjamini and Yosef Hochberg. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300.
- Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42:335–359.

- Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed Elhoseiny. 2022. Visualgpt: Data-efficient adaptation of pretrained language models for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18030–18040.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023).
- Zhuyun Dai, Arun Tejasvi Chaganty, Vincent Y Zhao, Aida Amini, Qazi Mamunur Rashid, Mike Green, and Kelvin Guu. 2022. Dialog inpainting: Turning documents into dialogs. In *International Conference on Machine Learning*, pages 4558–4586. PMLR.
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. 2017. Visual dialog. In *CVPR*.
- Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv preprint arXiv:2305.04790*.
- Dalu Guo, Chang Xu, and Dacheng Tao. 2019. Image-question-answer synergistic network for visual dialog. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10434–10443.
- Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and Steven Hoi. 2023. From images to textual prompts: Zero-shot visual question answering with frozen large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10867–10877.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- 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*.
- Bernd Huber, Daniel McDuff, Chris Brockett, Michel Galley, and Bill Dolan. 2018. Emotional dialogue generation using image-grounded language models. In *CHI*.
- Xiaoze Jiang, Jing Yu, Zengchang Qin, Yingying Zhuang, Xingxing Zhang, Yue Hu, and Qi Wu. 2020a. Dualvd: An adaptive dual encoding model for deep visual understanding in visual dialogue. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11125–11132.
- Xiaoze Jiang, Jing Yu, Yajing Sun, Zengchang Qin, Zihao Zhu, Yue Hu, and Qi Wu. 2020b. Dam: De-liberation, abandon and memory networks for generating detailed and non-repetitive responses in visual dialogue. In *IJCAI*.
- Gi-Cheon Kang, Jaeseo Lim, and Byoung-Tak Zhang. 2019. Dual attention networks for visual reference resolution in visual dialog. 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 2024–2033.
- Gangwo Kim, Sungdong Kim, Kang Min Yoo, and Jaewoo Kang. 2022. Generating information-seeking conversations from unlabeled documents. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2362–2378.
- Satwik Kottur, José MF Moura, Devi Parikh, Dhruv Batra, and Marcus Rohrbach. 2018. Visual coreference resolution in visual dialog using neural module networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 153–169.
- Solomon Kullback and Richard A Leibler. 1951. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Mengchen Liu and Chongyan Chen. 2024. An evaluation of gpt-4v and gemini in online vqa. *arXiv preprint arXiv:2312.10637*.
- Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon Bharti, and Ming Zhou. 2021. Univl: A unified video and language pre-training model for multimodal understanding and generation. *arXiv preprint arXiv:2002.06353*.
- Zheng-Yu Niu, Hua Wu, Haifeng Wang, et al. 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.
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2018. Meld: A multimodal multi-party dataset for emotion recognition in conversations. *arXiv preprint arXiv:1810.02508*.
- 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.

- Sina Semnani, Violet Yao, Heidi Zhang, and Monica Lam. 2023. [WikiChat: Stopping the hallucination of large language model chatbots by few-shot grounding on Wikipedia](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 2387–2413, Singapore. Association for Computational Linguistics.
- Yue Wang, Shafiq Joty, Michael R Lyu, Irwin King, Caiming Xiong, and Steven CH Hoi. 2020. Vd-bert: A unified vision and dialog transformer with bert. *EMNLP*.
- Qiucheng Wu, Handong Zhao, Michael Saxon, Trung Bui, William Yang Wang, Yang Zhang, and Shiyu Chang. 2024. Vsp: Assessing the dual challenges of perception and reasoning in spatial planning tasks for vlms. *arXiv preprint arXiv:2407.01863*.
- Te-Lin Wu, Satwik Kottur, Andrea Madotto, Mahmoud Azab, Pedro Rodriguez, Babak Damavandi, Nanyun Peng, and Seungwhan Moon. 2023. Simmc-vr: A task-oriented multimodal dialog dataset with situated and immersive vr streams. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6273–6291.
- Wei Wu, Xu Sun, and Houfeng Wang. 2018. Question condensing networks for answer selection in community question answering. In *ACL*.
- Alexandros Xenos, Niki Maria Foteinopoulou, Ioanna Ntinou, Ioannis Patras, and Georgios Tzimiropoulos. 2024. Vllms provide better context for emotion understanding through common sense reasoning. *arXiv preprint arXiv:2404.07078*.
- Tianhao Yang, Zheng-Jun Zha, and Hanwang Zhang. 2019. Making history matter: History-advantage sequence training for visual dialog. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2561–2569.
- Ze Yang, Wei Wu, Huang Hu, Can Xu, Wei Wang, and Zhoujun Li. 2021. Open domain dialogue generation with latent images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14239–14247.
- Haolan Zhan, Zhuang Li, Xiaoxi Kang, Tao Feng, Yuncheng Hua, Lizhen Qu, Yi Ying, Mei Rianto Chandra, Kelly Rosalin, Jureynolds Jureynolds, et al. 2024. Renovi: A benchmark towards remediating norm violations in socio-cultural conversations. *arXiv preprint arXiv:2402.11178*.
- Haolan Zhan, Sameen Maruf, Lizhen Qu, Ingrid Zukerman, and Gholamreza Haffari. 2023. Turning flowchart into dialog: Plan-based data augmentation for low-resource flowchart-grounded troubleshooting dialogs. *arXiv preprint arXiv:2305.01323*.
- Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

## A Baseline Models

**UniVL** (Luo et al., 2021) is a unified video and language pre-training model for multi-modal understanding and generation. UniVL model adapts Transformer as the backbone and has individual language and video encoder, following with a cross-encoder and decoder module.

**MiniGPT-4** (Zhu et al., 2023) contains a vision encoder with a pre-trained ViT and Q-Former model, a single linear projection layer, and an advanced Vicuna large language model (LLM). MiniGPT-4 freezes the vision part and only requires training the linear projection layer to align the visual features with the Vicuna.

**Video-LLaMa** (Zhang et al., 2023) maintains a similar architecture with the MiniGPT-4, including the ViT and Q-Former for the visual and audio encoder. On the top of the architecture, a LLM (LLaMa or Vicuna) is followed to align multi-modal features with contextual features.

**MultiModal-GPT** (Gong et al., 2023) is based on the open-flamingo (Alayrac et al., 2022) model. MultiModal-GPT consists of a vision encoder from CLIP, a sampler to receive the spatial features from the vision encoder, and a language decoder LLaMA.

**LLaVa-1.5** (Liu et al., 2023) considers the pre-trained CLIP visual encoder ViT-L/14, then the visual feature will be sent into a simple linear layer to connect image features into the word embedding space. Finally, LLaVa-1.5 choose the Vicuna as the LLM backbone.

## B An Example of Synthetic Data Generation

We provide an example for better understanding the synthetic data generation process. Suppose we have prepared a dialogue shown in the Figure 2. Then, we prompt ChatGPT to generate a synthetic visual description. We present an example in Figure 3 to illustrate this process.

**Synthetic Data Generation Prompt**

Can you imagine a potential scene for the context of the following dialogue? Try to incorporate several elements: Facial expression, Human Action and Human Position.

Here is an example:  
<Dialogue Context>  
Ross: Or! Or, we could go to the bank, close our accounts and cut them off at the source.

<Visual Description>  
Ross raises his eyebrows, leans forward, and speaks with a hint of frustration and determination.

Now, please generate a reasonable visual description for the last utterance of the following dialogue:

<Dialogue Context>  
Joey: You like it? You really like it?  
Chandler: Oh, yeah!  
Joey: Which part exactly?

Please try to generate a concise description within 30 words.

---

<Output - Synthetic Visual Description>  
In a cozy, dimly lit living room, Joey sits on a leather sofa, with his hands crossed on his lap, leans forward, a casual smile towards Chandler.

Figure 3: An example of synthetic data generation.

# User Review Writing via Interview with Dialogue Systems

Yoshiki Tanaka and Michimasa Inaba  
The University of Electro-Communications  
{y-tanaka, m-inaba}@uec.ac.jp

## Abstract

User reviews on e-commerce and review sites are crucial for making purchase decisions, although creating detailed reviews is time-consuming and labor-intensive. In this study, we propose a novel use of dialogue systems to facilitate user review creation by generating reviews from information gathered during interview dialogues with users. To validate our approach, we implemented our system using GPT-4 and conducted comparative experiments from the perspectives of system users and review readers. The results indicate that participants who used our system rated their interactions positively. Additionally, reviews generated by our system required less editing to achieve user satisfaction compared to those by the baseline. We also evaluated the reviews from the readers' perspective and found that our system-generated reviews are more helpful than those written by humans. Despite challenges with the fluency of the generated reviews, our method offers a promising new approach to review writing.

## 1 Introduction

User reviews posted on e-commerce and review sites benefit both buyers and sellers. For users considering a purchase, these reviews provide product details and insights, thereby supporting the decision-making process of purchase. Additionally, they act as feedback for sellers, enabling product quality improvement. Reviews that contain specific product details or are based on personal user experiences are especially valuable. However, writing these insightful reviews can be time-consuming and labor-intensive, as the task requires careful language and proofreading for potential buyers.

To reduce users' burden of writing reviews, researchers (e.g., (Dong et al., 2017; Li and Tuzhilin, 2019)) have focused on automated review genera-

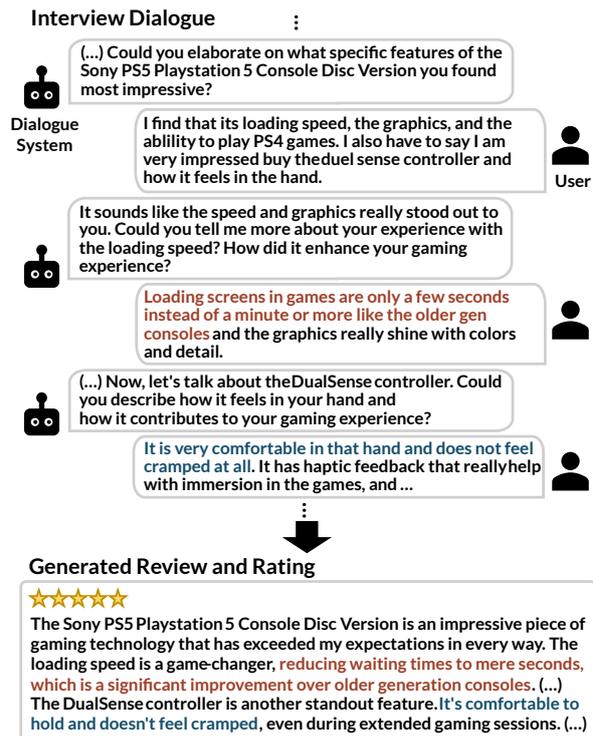


Figure 1: Example of review creation supported by the proposed method. In the dialogue and review, the contents of the same-colored text correspond to each other.

tion, thereby making the review creation process more efficient.

Despite these efforts, the ability of existing review generation systems to personalize for users and products remains limited. Previous studies have explored utilizing a wide range of information about products and users for generating reviews. However, available subjective information is limited. Specifically, without direct input from the user, incorporating the user's actual experiences with the product into the generated review is challenging. This constraint significantly limits the system's ability to personalize for the user. To overcome this problem, we focused on supporting the creation of reviews by directly eliciting infor-

mation about products from users.

In this study, we propose the novel utilization of dialogue systems for creating user reviews. Figure 1 shows an example of the review creation process supported by the dialogue system according to our proposed method. First, the dialogue system acts as an interviewer, eliciting user opinions on products through interview dialogues. Second, the review text generator generates review text based on the dialogue history. Finally, the rating predictor predicts a rating consistent with the generated review text. Our method allows users to easily create reviews by simply interacting with the system, thus reducing the effort involved in review creation.

To evaluate our method, we implemented a system incorporating our approach using GPT-4. Subsequently, we conducted experiments using our system, collecting data on dialogues between the system and users, the generated reviews, predicted ratings, and participants' feedback on our system. We discuss the effectiveness of our method after analyzing the collected data. In summary, our main contributions are as follows:

1. As a novel application of dialogue systems, we propose a method for supporting user review creation. Furthermore, we developed a system incorporating our approach using GPT-4.
2. We conducted a comprehensive survey from the perspectives of system users and review readers, showing that our method can provide high-quality and helpful reviews for both parties.

## 2 Related Work

### 2.1 Interview Dialogue System and Dataset

The interview dialogues are aimed at eliciting information from the interviewees. Prior research suggests that surveys conducted on chatbot platforms yield higher-quality responses than web survey platforms (Kim et al., 2019). This finding indicates that employing dialogue systems to collect user opinions and impressions is a promising approach.

Researchers have collected interview dialogue data on various topics, including radio (Majumder et al., 2020), news (Zhu et al., 2021), sports (Sun et al., 2022), and cooking (Okahisa et al., 2022).

The objectives of these collections vary from analyzing dialogue patterns (Majumder et al., 2020; Okahisa et al., 2022) to dialogue summarization (Zhu et al., 2021). Here, we utilize the interview dialogue system to support the creation of helpful reviews.

### 2.2 Review Generation

User reviews reflect user's opinions and requests regarding a product. These insights benefit buyers and sellers. Additionally, user reviews have a wide range of applications. Previous research has applied reviews to natural language processing tasks such as recommendations (Qiu et al., 2021), opinion summarization (Bražinskas et al., 2020), and task-oriented dialogue (Zhao et al., 2023).

User reviews that include detailed information about the product and user experiences are useful. However, writing these reviews is a labor-intensive task for humans. To increase the efficiency of this process, researchers have proposed automated review generation models, enhancing their review generation capabilities by utilizing information such as ratings (Dong et al., 2017; Sharma et al., 2018; Li et al., 2019; Kim et al., 2020), images (Truong and Lauw, 2019; Vu et al., 2020), past reviews written by the user (Li and Tuzhilin, 2019), and aspect-oriented features (Li and Tuzhilin, 2019). Unlike these studies, we focus on the collaborative writing of user reviews with the support of the dialogue system.

Some researchers have focused on supporting users in creating reviews, similar to our approach (Ni and McAuley, 2018; Bhat et al., 2023). For example, Ni and McAuley proposed utilizing short phrases related to products that are provided by customers, such as review summaries and product titles, as auxiliary data for generating reviews (Ni and McAuley, 2018). In their system, the user provides information in a unidirectional manner. In contrast, we utilize an interview-specific dialogue system to collect information from the user through interactive interaction. The dialogue system can ask follow-up questions to obtain additional details regarding a product although this information may be ambiguous. This capability supports the creation of detailed reviews.

### 2.3 Dialogue Summarization

In our method, we proposed to convert conversational data (i.e., interview dialogue history) into non-conversational data (i.e., review texts).

Therefore, our work is closely related to dialogue summarization research. To build an effective model for dialogue summarization, researchers have proposed diverse approaches to learning methods (Zou et al., 2021; Li et al., 2023; Zhong et al., 2022; Zhang et al., 2021). Additionally, researchers have built dialogue summarization datasets that can be used for training models; these datasets cover daily life conversations (Gliwa et al., 2019; Chen et al., 2021), meetings (Carletta et al., 2006; Zhong et al., 2021), TV series (Chen et al., 2022), and media dialogue (Zhu et al., 2021). While these studies aim to condense dialogue histories into brief texts, our approach takes a different direction. We focus on extracting useful product information for readers from interview dialogues and organizing it into a non-conversational data format, rather than compressing it into shorter text.

### 3 Methodology

To create useful reviews, reviewers must provide detailed product information. Interview dialogue systems are employed to effectively elicit this information. To enhance readability, we propose organizing the dialogue history into a non-dialogue format. Our method comprises three processes: interview dialogue, review text generation, and rating prediction. In this paper, the systems that perform these processes are referred to as the “interview dialogue system,” the “review text generator,” and the “rating predictor,” respectively. Our system utilizes these components in sequence to generate reviews as the output. An overview of our system is shown in Figure 2. We use the gpt-4-0613 model to implement our system.

#### 3.1 Interview Dialogue System

To assist potential buyers in making purchase decisions, guiding users to create helpful reviews is crucial. In our approach, therefore, our system should be designed to effectively collect information from the user. To achieve this, we propose utilizing an interview dialogue system. For the interview dialogue system, it is desirable to elicit both the pros and cons of a product in a balanced and detailed manner. Specifically, the system should be capable of asking follow-up questions about the content mentioned by the user or changing the topic to inquire about different aspects of the product.

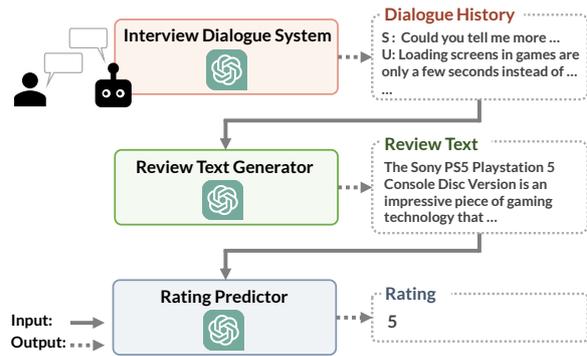


Figure 2: Overview of our system. First, the interview dialogue system interviews the user to elicit their impressions and requests about the product they used. Next, the review text generator uses the dialogue history as input to generate a review text. Finally, the rating predictor predicts a rating consistent with the sentiment of the generated review text.

We designed a prompt that incorporated instructions for the system to perform these behaviors. Moreover, aiming to both collect sufficient information for creating reviews and ensure users don’t become bored, we added constraints regarding the number of turns to the prompt. In our experiments, we adopted instructions to ask at least 8 questions and conclude the interview within 15 turns. Additionally, to ensure the interview does not continue indefinitely, we externally implemented a setting in the interview dialogue system to end the dialogue after 15 turns. The prompt template for the interview dialogue is shown in Appendix A.1.

#### 3.2 Review Text Generator

Although the dialogue history between interview dialogue systems and users offers useful and detailed product information, it often contains redundancies. Consequently, it is not appropriate to post it directly as a user review. Therefore, we propose transforming the dialogue history into a format suitable for reviews. Our review text generator aims to capture the essence of the interview dialogue history while generating review texts from the perspective of the user. To generate reviews that align with the user’s feedback, the system must faithfully reflect the content of the dialogue history in the review text. Our prompts include instructions to concisely summarize important information mentioned during the interview and generate the main body of the user review. The prompt template for the review generation is shown in Appendix A.2.

### 3.3 Rating Predictor

In e-commerce and review platforms, customer ratings are aggregated into a single score, providing other users with an initial impression of the product. For an aggregated score to be reliable, reviewers must assign ratings that accurately reflect the content of their review text. While the ratings impact the reputation widely, considering the potential for human error in assigning ratings, automating the task might be an effective solution. Our rating predictor automatically outputs a rating consistent with the sentiment of the input review text, ranging from 1 to 5 as an integer. Ratings consistent with the content of the review texts could reduce exaggerated scoring caused by user subjectivity. As a result, this can improve the reliability of the ratings.

We utilized GPT-4 to implement a rating predictor. To enhance predictive performance, we designed prompts that apply chain-of-thought prompting (Wei et al., 2022; Wang et al., 2023; Kojima et al., 2022), that feeds large language models not only examples of question-and-answer pairs but also examples of the thought processes leading to those answers. In this study, we collected five sets of product titles, review texts, and ratings from Amazon.com to create output exemplars, each corresponding to ratings from 1 to 5. Subsequently, for each set, we wrote descriptions of the reasoning paths leading to the prediction of the rating from the product title and review text. We used these as few-shot exemplars within the context. Similarly, for target reviews, GPT-4 is encouraged to generate a reasoning path and an answer.

## 4 Experiments

We aim to facilitate the review-writing process for reviewers and provide helpful reviews to readers. To investigate the practicality of our method, we conducted evaluations from the perspectives of system users (Section 4.1) and review readers (Section 4.2).

### 4.1 Participant Evaluation

To evaluate our method, we collected feedback through interviews, generated reviews, and questionnaires. Data collection was conducted through Amazon Mechanical Turk (MTurk)<sup>1</sup>.

<sup>1</sup><https://www.mturk.com>

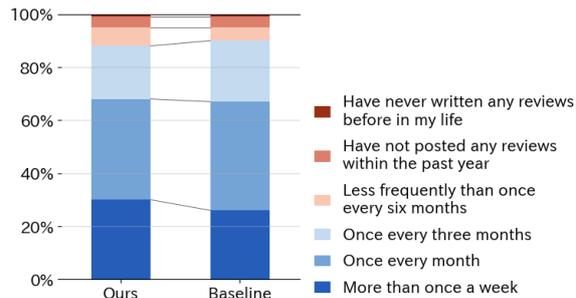


Figure 3: Participant responses to “In the past year, how often have you posted reviews?”

### 4.1.1 Experimental Setup

We tuned the temperature for each system. For the interview dialogue system, the temperature was set to 0.2. The review text generator and rating predictor generate outputs that are faithful to the input. Therefore, we set the temperature to 0 for these systems to suppress the diversity of the generated text.

### 4.1.2 Baseline System

To demonstrate the effectiveness of using interview dialogue systems that adapt questions based on the context, we constructed a baseline system. The baseline system replaces the interview dialogue system with one that asks manually created questions in a fixed order. To construct the baseline system, we manually created nine questions on topics such as the reason for purchasing the product and the evaluation of the product in comparison with other products. All questions asked by the baseline system are listed in Appendix B. We collected data using this system in the same manner as with our proposed system.

### 4.1.3 Evaluation Procedure

Initially, participants conducted an interview dialogue with our interview dialogue system. After the interview, they were presented with the generated reviews and ratings. Participants then completed a post-interview survey comprising multiple-choice and open-ended questions. For each setting, we recruited 100 participants located in AU, CA, NZ, GB, or the US and had a 95% approval rate with at least 500 previously approved HITs.

### 4.1.4 Post-Interview Survey

After the interview, participants responded to a post-interview survey. Several questions in this

Table 1: Likert Items in Post-interview Survey

Dimension	Labels in Figure 4	Statements
Interview	Enjoyable	<i>How fun was your interaction with the chatbot interviewer?</i>
	Skillful	<i>The interviewer skillfully elicited your impressions or opinions.</i>
	In-depth	<i>The chatbot interviewer attempted to elicit your impressions or opinions in depth.</i>
Review	Faithful	<i>The system-generated review faithfully reflects what you said during your interviews.</i>
	Concise	<i>The system-generated review offers a concise summary of the points you mentioned during the interview.</i>
System	Quality	<i>Please rate the overall quality of the system.</i>
	Burdened(I)	<i>I felt burdened to have an interview chat about the product.</i>
	Burdened(R)	<i>Writing a review with the support of the system is more burdensome than writing a review yourself.</i>

Participants' Responses: Ours (Upper) vs. Baseline System (Lower)

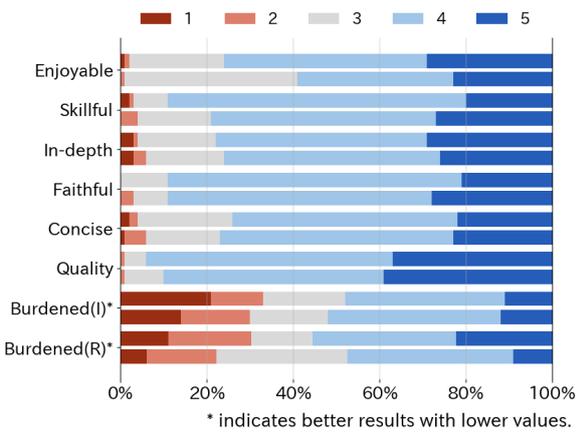


Figure 4: Participant responses to questions on a Likert scale from 1 (Strongly disagree) to 5 (Strongly agree) in a post-interview survey. For each question, the upper bar shows the results from our system and the lower bar shows the baseline results.

survey were answered using a 5-point Likert scale. These questions are related to the interview dialogue, the generated reviews, and the overall system (See Table 1).

We also asked participants how frequently they post reviews to compare with their usual review-writing experiences. As shown in Figure 3, 95% of the participants posted at least one user review in the past year. Additionally, participants were asked to rate the product they selected by responding to the question: "If you were to rate this product again, what rating would you give it?" and provided a rating from 1 to 5.

#### 4.1.5 Participant Feedback Analysis

Figure 4 illustrates the distribution of responses to eight questions<sup>2</sup>. Regarding the dimensions of

<sup>2</sup>For Burdened(R), we excluded responses from participants who selected the "Have never written any reviews before in my life" option to the question in Figure 3.

the interview and review, most participants evaluated two components positively: our interview dialogue system and our review text generator. Participants showed a similar positive trend across two settings for the four items: In-depth, Faithful, Concise, and Quality. Notably, for Quality, 90% or more of the participants rated the overall quality of the systems positively.

Our system provided users with more enjoyable interviews and higher satisfaction regarding the generated reviews compared with the baseline system. As shown in Figure 4, when using our system based on GPT-4, more participants agreed that interacting with the system was fun. Moreover, the difference in the methods used to elicit information—our interview dialogue system and the baseline—impacts users' enjoyment, with statistically significant differences (Mann-Whitney U test,  $p < 0.05$ ). Participants also responded to the multiple-choice question, "If you had to edit and post a system-generated review to your satisfaction, how much of it would you need to rewrite?". Figure 5 shows that different types of systems resulted in varied response distributions. In particular, 38% and 27% of participants using the baseline system and our system, respectively, responded that they needed to rewrite more than 50% of the review. These results indicate that our system can provide reviews with higher satisfaction than the baseline.

Our system imposed a greater burden on participants. Figure 4 shows that a higher percentage of participants agreed that *writing a review with the support of our system is more burdensome than writing alone*, compared to the baseline. We argue that the response time of the system is one of the reasons for this difference. Our GPT-4-based system, which generates responses based on users' utterances, takes a longer time to gener-

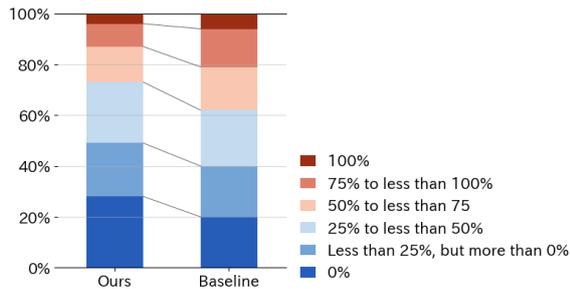


Figure 5: Participant responses to “If you had to edit and post a system-generated review to your satisfaction, how much of it would you need to rewrite?”

ate responses than a baseline that asks predefined questions. Notably, several participants suggested that the response speed of our system should be improved. In response to the free-form question “What is one enhancement that can be made to improve this system?”, we received answers such as “more fast replies” and “need quick reply.” In our experiments, unlike the ChatGPT interface<sup>3</sup>, we did not employ real-time response generation using streaming functionality. Adding this feature would be an effective modification to enhance our system’s response speed, which is expected to significantly improve user experience.

#### 4.1.6 Case Study

Our interview dialogue system can generate follow-up questions that explore the content of users’ ambiguous responses in depth. Table 2 shows an example of the data collected, comprising the dialogue history regarding an electric shaver and the corresponding review text generated. During the interview, our system initially asked about the participant’s overall satisfaction with the product, to which the participant replied, “... well satisfied but with few minor issues.” Based on this response, our system posed follow-up questions to clarify the aspects that the participant was satisfied with and the issues they mentioned. As demonstrated in this example, our system can elicit deeper information about products from users compared with the baseline system.

Additionally, our interview dialogue system can change topics during an interview to collect information on various aspects of a product. For the interview dialogues shown in Table 2, the first three turns focus on the participant’s satisfaction with the product. Subsequently, the system changes the

topic by saying “Now, let’s go back to the issue” and thus shifting focus to the issues that the user mentioned in the first turn. In the next turn, our system shifts the topic again to highlight the product’s impressive features. By switching topics in this manner, our system can acquire information on various aspects of the product.

Our system can also generate review texts and ratings that faithfully reflect the content mentioned in the interview. As shown in the lower part of Table 2, the generated review text closely aligns with user comments such as “i would say well satisfied ...” and “some times small hair from the beard gets stuck.” Moreover, our system accounts for issues mentioned by the user and predicts a rating of 4, one less than the maximum.

## 4.2 Third-Party Evaluation

Helpful user reviews offer effective decision-making support for potential buyers and sellers. To assess the quality of the reviews generated by our system from the reader’s perspective, we conducted multidimensional evaluations.

### 4.2.1 Experimental Procedure

This evaluation was conducted in a pairwise setting. In this experiment, annotators were presented with two reviews: one generated by our system and another for comparison. We compared our system-generated reviews with Amazon reviews to evaluate the quality difference from human-written reviews. Additionally, to assess the differences from the perspective of the readers, we compared our system-generated reviews with those generated by a baseline system, which were collected in Section 4.1.

Annotators chose either the superior option or a tie for each of the eight dimensions, as shown in Table 3. We recruited the annotators through MTurk.

### 4.2.2 Selection Criteria for Comparative Reviews

We collected human-written reviews from the Amazon Review Dataset (Ni et al., 2019). To facilitate comparison for the annotators, we initially removed any human-written reviews where the category or rating did not match those generated by our system. After this filtering, to avoid selecting poor-quality reviews, we retained reviews in the top 5% of human-written reviews based on helpfulness votes, using these as compar-

<sup>3</sup><https://chatgpt.com/>

Table 2: Example of a dialogue between our system and the user (top, S: system, U: user) and generated review text (bottom). In this example, our rating predictor predicted a rating of 4.

---

S (...) To start off, could you tell me about your overall satisfaction with this product?

U **i would say well satisfied** but with few minor issues.

S I see, you're generally satisfied but there are a few minor issues. Could you elaborate on what you particularly like about the Braun Series 9 9370cc Electric Shaver?

U **some times small hair from the beard gets stuck**s in particular spot on the blade area in the machine.

S I understand that small hairs getting stuck in the blade area is one of the issues you've encountered. But before we delve into that, could you first share what you find satisfying or impressive about the product? What are its positive points in your opinion?

U the postivie sides is smooth shaving experience and battery capacity which is good

S It's great to hear that you appreciate the smooth shaving experience and the battery capacity. Now, let's go back to the issue you mentioned earlier about small hairs getting stuck. Could you tell me more about this problem? How often does it happen and how does it affect your shaving experience?

U i wouls say not often, ifaced this issue only once.

S I see, (...) That's good to know. Now, aside from the smooth shaving experience and the battery capacity, are there any other features of the Braun Series 9 9370cc Electric Shaver that you find beneficial or impressive?

...

---

**I am generally satisfied** with the Braun Series 9 9370cc Rechargeable Wet & Dry Men's Electric Shaver with Clean & Charge Station. The product offers a smooth shaving experience, which I find impressive. The battery capacity is also commendable, lasting for a good amount of time. I appreciate the versatility of the product, as it can be used for both wet and dry shaving. The Clean & Charge Station that comes with the shaver is quite efficient and useful. It effectively cleans and charges the shaver, making it more convenient compared to other products I've used before. **However, I did encounter a minor issue where small hairs from my beard got stuck** in a particular spot on the blade area. This happened only once and did not significantly affect my overall shaving experience. Considering all the features and my experience with the product, I believe it offers good value for its price. I would recommend the Braun Series 9 9370cc Electric Shaver to others, mainly for its battery performance and its variety of features, including its wet and dry usage.

---

ison candidates. Finally, among these candidates, we selected the review of the product whose title achieved the highest ROUGE-L score (Lin, 2004).

We applied the same process when matching reviews generated by the baseline system; however, for some reviews collected using our system, there were no corresponding reviews in the baseline data that matched both category and rating. Therefore, we excluded these unmatched data. Consequently, while we conducted comparisons for all 100 reviews written by humans, only 96 baseline-generated reviews met the criteria.

### 4.2.3 Results and Discussions

The overall results are presented in Table 4. The annotators prefer the reviews generated by our system to those written by humans or generated by the baseline system. Notably, the reviews generated by our system are helpful, provide a balanced view of pros and cons, and offer comprehensive information. These findings indicate that our interview dialogue system is capable of eliciting a wide range of information about products from users through topic transitions.

The reviews generated by our system lack the fluency of human-written reviews. For instance, our review text generator tends to use the formal product title when referring to the product. Addi-

tionally, human-written reviews contain more individual experiences compared with those generated by our system. Despite these limitations, our system has high scalability, offering the potential for improvement. Specifically, our system's output could be enhanced by refining the prompts to generate texts that are more human-like and elicit detailed usage experiences from users.

By replacing the baseline system, which uses fixed questions, with our interview dialogue system, we observe improvements across all metrics. Notably, our system can generate reviews that are rich in experience-based information, contain more detailed information, and cover a broader range of topics. This demonstrates that our system can elicit more detailed and extensive information from users through follow-up questions and topic transitions.

### 4.3 Discussion on Predicted Ratings

To further explore the characteristics of the reviews and ratings generated by our system, we analyze them along two axes: the difference based on the source of the ratings (comparing ratings assigned by humans to those predicted by our system) and the difference based on the annotators (comparing the ratings given by system users to those assigned by third parties). To obtain ratings

Table 3: Questions in comparative evaluation

Labels in Table 4	Questions
Helpfulness	<i>Which review would be more helpful for making a purchase decision?</i>
Fluency	<i>Which review exhibits a more fluent and human-like writing style?</i>
Conciseness	<i>Which review is more concise and to the point?</i>
Experience	<i>Which review provides more information based on the actual usage experience of the product?</i>
Balance	<i>Which review presents a more balanced view of the product’s pros and cons?</i>
Depth	<i>Which review provides more in-depth information about any specific aspect of the product?</i>
Coverage	<i>Which review mentions a more comprehensive range of product aspects?</i>
Overall	<i>Which review is overall more preferable?</i>

Table 4: Results of third-party evaluation. The values represent the percentage of votes each received.

Reviews	Helpfulness	Fluency	Conciseness	Experience	Balance	Depth	Coverage	Overall
Human	38.0	<b>47.0</b>	37.0	<b>57.0</b>	37.0	43.0	40.0	41.0
Tie	6.0	15.0	6.0	9.0	15.0	10.0	5.0	7.0
Ours	<b>56.0</b>	38.0	<b>57.0</b>	34.0	<b>48.0</b>	<b>47.0</b>	<b>55.0</b>	<b>52.0</b>
Baseline	38.5	28.1	45.8	21.9	34.4	35.4	35.4	37.5
Tie	12.5	33.3	6.2	16.7	15.6	12.5	10.4	17.7
Ours	<b>49.0</b>	<b>38.5</b>	<b>47.9</b>	<b>61.5</b>	<b>50.0</b>	<b>52.1</b>	<b>54.2</b>	<b>44.8</b>

Table 5: Average absolute difference in ratings between Amazon customers and Turkers (top-left), between system-predicted ratings and Turkers’ ratings for system-generated reviews (top-right), and between system-predicted ratings and participants’ ratings (bottom-right, see Section 4.1.4).

Annotator/Source	Human-written	System-generated
Turkers	0.59	0.12
Participants in Section 4.1	-	0.57

assigned by third parties, we newly recruited annotators from MTurk and asked them to assign ratings to both the human-written reviews (left column)<sup>4</sup> and those generated by our system (right column). We also collected ratings assigned by participants from the experiments in Section 4.1. Note that these participants, unlike the Turkers, had seen the ratings predicted by our system.

The results in the top row of Table 5 demonstrate that the difference between the ratings predicted by our system and those assigned by third parties is remarkably smaller than the difference found in human-written reviews. This finding indicates that the sentiment of the reviews generated by our system is easily comprehensible to readers.

The ratings predicted by our system, as shown in the right column of Table 5, align more closely with those assigned by third-party annotators than with those of system users. This finding indicates that our system emphasizes objectivity over sub-

jectivity in its ratings.

The aforementioned observations indicate that our system generates review texts that are easy for humans to understand and provide more objective ratings. This finding suggests that our interview dialogue system and review text generator can generate reviews that accurately capture reviewers’ sentiments, thereby supporting informed purchasing decisions, while the rating predictor also provides highly objective and reliable ratings.

## 5 Conclusion

In this study, we present a novel method for utilizing dialogue systems to facilitate user review creation. Our approach involves three processes: interview dialogue, review text generation, and rating prediction. Although ensuring the fluency of the system-generated reviews remains a challenge, our method provides high-quality and helpful reviews for both reviewers and their readers.

Our method possesses high scalability. For instance, feeding product descriptions into our interview dialogue system could lead to deeper interview dialogues about more detailed information. However, our experiments have shown that even without such extensions, our system is capable of providing reviews that are more helpful than human-written ones. Furthermore, adapting our dialogue system’s strategies to user preferences during review writing could improve user experience. Further research can accomplish this objective by conducting a more detailed analysis of user preferences.

<sup>4</sup>For the annotations, we used 100 human-written reviews selected in Section 4.1.

## References

- Advait Bhat, Saaket Agashe, Parth Oberoi, Niharika Mohile, Ravi Jangir, and Anirudha Joshi. 2023. [Interacting with next-phrase suggestions: How suggestion systems aid and influence the cognitive processes of writing](#). In *Proceedings of the 28th International Conference on Intelligent User Interfaces*, IUI '23, page 436–452, New York, NY, USA. Association for Computing Machinery.
- Arthur Bražinskas, Mirella Lapata, and Ivan Titov. 2020. [Unsupervised opinion summarization as copycat-review generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5151–5169, Online. Association for Computational Linguistics.
- Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, Guillaume Lathoud, Mike Lincoln, Agnes Lisowska, Iain McCowan, Wilfried Post, Dennis Reidsma, and Pierre Wellner. 2006. The ami meeting corpus: A pre-announcement. In *Machine Learning for Multimodal Interaction*, pages 28–39, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Mingda Chen, Zewei Chu, Sam Wiseman, and Kevin Gimpel. 2022. [SummScreen: A dataset for abstractive screenplay summarization](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8602–8615, Dublin, Ireland. Association for Computational Linguistics.
- Yulong Chen, Yang Liu, and Yue Zhang. 2021. [DialogSum challenge: Summarizing real-life scenario dialogues](#). In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 308–313, Aberdeen, Scotland, UK. Association for Computational Linguistics.
- Li Dong, Shaohan Huang, Furu Wei, Mirella Lapata, Ming Zhou, and Ke Xu. 2017. [Learning to generate product reviews from attributes](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 623–632, Valencia, Spain. Association for Computational Linguistics.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. [SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization](#). In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79, Hong Kong, China. Association for Computational Linguistics.
- Jihyeok Kim, Seungtaek Choi, Reinald Kim Amplayo, and Seung-won Hwang. 2020. [Retrieval-augmented controllable review generation](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2284–2295, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Soomin Kim, Joonhwan Lee, and Gahgene Gweon. 2019. [Comparing data from chatbot and web surveys: Effects of platform and conversational style on survey response quality](#). In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI '19, page 1–12, New York, NY, USA. Association for Computing Machinery.
- Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 22199–22213. Curran Associates, Inc.
- Junyi Li, Wayne Xin Zhao, Ji-Rong Wen, and Yang Song. 2019. [Generating long and informative reviews with aspect-aware coarse-to-fine decoding](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1969–1979, Florence, Italy. Association for Computational Linguistics.
- Pan Li and Alexander Tuzhilin. 2019. [Towards controllable and personalized review generation](#). 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 3237–3245, Hong Kong, China. Association for Computational Linguistics.
- Yu Li, Baolin Peng, Pengcheng He, Michel Galley, Zhou Yu, and Jianfeng Gao. 2023. [DIONYSUS: A pre-trained model for low-resource dialogue summarization](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1368–1386, Toronto, Canada. Association for Computational Linguistics.
- 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.
- Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. 2020. [Interview: Large-scale modeling of media dialog with discourse patterns and knowledge grounding](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8129–8141, Online. Association for Computational Linguistics.
- Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. [Justifying recommendations using distantly-labeled reviews and fine-grained aspects](#). 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 188–197, Hong Kong, China. Association for Computational Linguistics.

- Jianmo Ni and Julian McAuley. 2018. [Personalized review generation by expanding phrases and attending on aspect-aware representations](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 706–711, Melbourne, Australia. Association for Computational Linguistics.
- Taro Okahisa, Ribeka Tanaka, Takashi Kodama, Yin Jou Huang, and Sadao Kurohashi. 2022. [Constructing a culinary interview dialogue corpus with video conferencing tool](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3131–3139, Marseille, France. European Language Resources Association.
- Zhaopeng Qiu, Xian Wu, Jingyue Gao, and Wei Fan. 2021. [U-bert: Pre-training user representations for improved recommendation](#). volume 35, pages 4320–4327.
- Vasu Sharma, Harsh Sharma, Ankita Bishnu, and Labhesh Patel. 2018. [Cyclegen: Cyclic consistency based product review generator from attributes](#). In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 426–430, Tilburg University, The Netherlands. Association for Computational Linguistics.
- Hanfei Sun, Ziyuan Cao, and Diyi Yang. 2022. [SPORTSINTERVIEW: A large-scale sports interview benchmark for entity-centric dialogues](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 5821–5828, Marseille, France. European Language Resources Association.
- Quoc-Tuan Truong and Hady Lauw. 2019. [Multimodal review generation for recommender systems](#). In *The World Wide Web Conference, WWW '19*, page 1864–1874, New York, NY, USA. Association for Computing Machinery.
- Xuan-Son Vu, Thanh-Son Nguyen, Duc-Trong Le, and Lili Jiang. 2020. [Multimodal review generation with privacy and fairness awareness](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 414–425, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Xinyuan Zhang, Ruiyi Zhang, Manzil Zaheer, and Amr Ahmed. 2021. [Unsupervised abstractive dialogue summarization for tete-a-tetes](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14489–14497.
- Chao Zhao, Spandana Gella, Seokhwan Kim, Di Jin, Devamanyu Hazarika, Alexandros Papangelis, Behnam Hedayatnia, Mahdi Namazifar, Yang Liu, and Dilek Hakkani-Tur. 2023. [“what do others think?”: Task-oriented conversational modeling with subjective knowledge](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 309–323, Prague, Czechia. Association for Computational Linguistics.
- Ming Zhong, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2022. [Dialoglm: Pre-trained model for long dialogue understanding and summarization](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11765–11773.
- Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir Radev. 2021. [QMSum: A new benchmark for query-based multi-domain meeting summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5905–5921, Online. Association for Computational Linguistics.
- Chenguang Zhu, Yang Liu, Jie Mei, and Michael Zeng. 2021. [MediaSum: A large-scale media interview dataset for dialogue summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5927–5934, Online. Association for Computational Linguistics.
- Yicheng Zou, Bolin Zhu, Xingwu Hu, Tao Gui, and Qi Zhang. 2021. [Low-resource dialogue summarization with domain-agnostic multi-source pretraining](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 80–91, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

## A Prompt Template

### A.1 Prompt for Interview Dialogue

Table 6 shows a prompt template for interview dialogues. [PRODUCT\_NAME] is a placeholder for the product title, which will be replaced with the product title selected by the participant. [MAX\_QUESTION] and [MIN\_QUESTION] are placeholders for the maximum and minimum number of dialogue turns. In our experiments, we used 15 and 8, respectively.

## A.2 Prompt for Review Generation

Table 7 shows a prompt template for review generation. Similar to that for interviewing, [PRODUCT\_NAME] is a placeholder for the product title. [DIALOGUE] is a placeholder for the dialogue history, into which the interview dialogue history between our system and the participants is inserted.

## B Baseline Details

Table 8 shows a prompt template for review generation. Similar to that for interviewing, [PRODUCT\_NAME] is a placeholder for the product title. [DIALOGUE] is a placeholder for the dialogue history, into which the interview dialogue history between our system and the participants is inserted.

Table 6: Prompt template for interviewing.

---

Your role is “interviewer” and my role is “interviewee”.  
 About the product I am going to present, please elicit my impressions and opinions from me when I have touched it.

Note the following statements.

1. The interviewer elicits the interviewee’s satisfaction and dissatisfaction (the positive and negative points) with the product in a well-balanced and detailed.
2. In response to the interviewee’s response, the interviewer asks more in-depth questions about the aspect or elicits feedback about other aspects of the product.
3. Be sure to attach the name of your role at the beginning of your utterance. Since your role is “interviewer”, your generation should begin with “Interviewer:”.
4. Don’t generate interviewee’s utterances.
5. Add “[Wait\_for\_Response]” at the end of your utterance and wait for my response.
6. You must ask at least [MIN\_QUESTION] questions. In other words, the dialogue must continue for [MIN\_QUESTION] or more turns.
7. Having fulfilled the 6th statement, you can terminate the interview at your discretion. However, the interview must be completed within [MAX\_QUESTION] turns.
8. When you terminate the interview, add “[End\_of\_Interview]” at the end of your utterance.

Now, please elicit my impressions and opinions about the following product from me.  
 [PRODUCT\_NAME]

---

Table 7: Prompt template for review generation.

---

[DIALOGUE]

The above is a dialogue about “[PRODUCT\_NAME]” between the interviewer and the interviewee who has touched on this product.

Write a customer review about the product as if written by the interviewee, by briefly summarizing the important information mentioned in the above interview, such as the good and bad points of the product and the interviewee’s experience with it.

Do not output the review’s title.

The following is a body of the product review of the product written by the interviewee:

---

Table 8: Questions asked by the baseline system

---

Q-1	First, could you tell me about the features and functions of this product? What kind of product is this?
Q-2	What made you decide to purchase this product?
Q-3	If you have any points that you like or are satisfied with this product, please tell me in detail.
Q-4	What are the advantages of this product compared to other products?
Q-5	If you have any dissatisfaction with this product or areas for improvement for this product, please tell me in detail.
Q-6	What are the disadvantages of this product compared to other products?
Q-7	Who would this product be suitable for?
Q-8	Is this product worth the price? Also, why do you think so?
Q-9	Finally, do you have any requests or impressions about the product?

---

# Conversational Feedback in Scripted versus Spontaneous Dialogues: A Comparative Analysis

Ildikó Pilán<sup>1</sup>, Laurent Prévot<sup>2,3</sup>, Hendrik Buschmeier<sup>4</sup>, and Pierre Lison<sup>1</sup>

<sup>1</sup> Norwegian Computing Center, Oslo, Norway

<sup>2</sup> CEFC, CNRS, MEAE, Taipei, Taiwan

<sup>3</sup> Aix Marseille Université & CNRS, LPL, Aix-en-Provence, France

<sup>4</sup> Faculty of Linguistics and Literary Studies, Bielefeld University, Bielefeld, Germany

{pilan,plison}@nr.no

laurent.prevot@univ-amu.fr

hbuschme@uni-bielefeld.de

## Abstract

Scripted dialogues such as movie and TV subtitles constitute a widespread source of training data for conversational NLP models. However, there are notable linguistic differences between these dialogues and spontaneous interactions, especially regarding the occurrence of *communicative feedback* such as backchannels, acknowledgments, or clarification requests. This paper presents a quantitative analysis of such feedback phenomena in both subtitles and spontaneous conversations. Based on conversational data spanning eight languages and multiple genres, we extract lexical statistics, classifications from a dialogue act tagger, expert annotations and labels derived from a fine-tuned Large Language Model (LLM). Our main empirical findings are that (1) communicative feedback is markedly less frequent in subtitles than in spontaneous dialogues and (2) subtitles contain a higher proportion of negative feedback. We also show that dialogues generated by standard LLMs lie much closer to scripted dialogues than spontaneous interactions in terms of communicative feedback.

## 1 Introduction

While the amount of text data available for training or fine-tuning LLMs is large and growing steadily, spoken conversational data remains relatively scarce. Although corpora of spontaneous spoken interactions have been collected for various languages (Dingemans and Liesenfeld, 2022), those are generally of a modest size and limited to specific topics or tasks. Due to this scarcity of available data, a common approach for the development of conversational models is to rely on corpora of authored dialogues extracted from movie scripts (Danescu-Niculescu-Mizil and Lee, 2011) or movie and TV subtitles (Lison et al., 2018; Davies, 2021).

However, those dialogues are markedly different from spontaneous interactions. Most importantly, movie scripts and subtitles are explicitly written

with the aim of *narrating a story*. Subtitles must also abide to strict length constraints, and thus tend to only transcribe the most salient part of each turn. As a consequence, many conversational phenomena such as disfluencies (Shriberg, 1996), overlapping talk (Schegloff, 2000), and backchannels (Yngve, 1970) are either absent or uncommon in those dialogues, unless their presence happens to contribute to the storyline (Berliner, 1999; Chepinchik and Thompson, 2016).

This paper provides a quantitative analysis of how subtitles differ from spontaneous dialogues, focusing more specifically on *conversational feedback* (Allwood et al., 1992) and *grounding* (Clark and Schaefer, 1989) phenomena. To highlight differences in linguistic properties between subtitles and spontaneous conversation corpora, we first compile a range of lexical statistics and use a dialogue act tagger to estimate the relative frequencies of various feedback signals. To obtain more fine-grained estimates on three core feedback categories, respectively *Agreement / Acceptance*, *Acknowledgement / Backchannel* and *Negative Feedback*, we collect manual annotations on multiple dialogue samples and fine-tune a LLM on those annotations to automatically detect the presence of those feedback in our corpora. Finally, we apply the fine-tuned LLM on synthetic dialogues generated with standard autoregressive LLMs, and show that those dialogues are comparatively much closer to scripted dialogues than to spontaneous interactions when it comes to the frequency and type of conversational feedback. Those experiments are conducted for eight languages (English, Chinese, French, German, Hungarian, Italian, Japanese and Norwegian) for which corpora of spontaneous dialogues are readily available.

The paper is structured as follows. Section 2 reviews related work, and Section 3 presents the corpora employed in our experiments. Section 4 describes the observed lexical distributions of feed-

back phenomena and Section 5 compares them to estimates derived with a dialogue act tagger. In Section 6, we describe the manual annotation of dialogue samples and the fine-tuning of an LLM to automate this process. Finally, Section 7 describes the results of applying this LLM-based method to synthetic dialogues, and Section 8 concludes.

## 2 Related Work

### 2.1 Conversational Feedback and Grounding

A key aspect of any communicative activity is the management of the common ground, a process often called *conversational grounding* (Clark and Schaefer, 1989). The study of grounding and related phenomena, such as conversational feedback (Allwood et al., 1992), has been instrumental to cognitive approaches to communication (Clark, 1996), and to dialogue system development (Traum, 1994; Paek and Horvitz, 2000; Yaghoubzadeh et al., 2015).

Feedback and grounding can happen at any of the *levels of communication* that includes simple contact, perception, understanding and higher-level evaluation of what had been said (Allwood et al., 1992; Clark, 1996). Conversational feedback may appear at different positions in a dialogue. However, a number of corpus studies found that they have a tendency to occur at specific places, mostly where they cause little interference (Kjellmer, 2009). These places of occurrence have also been referred to as *Feedback Relevant Spaces* (Heldner et al., 2013; Howes and Eshghi, 2021). Although, arguably, any utterance relates directly or indirectly to grounding (through implicit and high level pragmatic inference, Clark and Schaefer 1989), *acknowledgments* and other positive feedback signals (see Ex. (1)), along with *repair* (see Ex. (2)), have been identified as the most prominent grounding mechanisms (Jefferson, 1972; Bunt, 1994). Their frequency in human-human dialogue is known to be very high (e.g., Stolcke et al., 2000a) and universal across languages (Liesenfeld and Dingemanse, 2022; Dingemanse et al., 2015). These conversational signals, while they do not cover all grounding phenomena, can therefore be seen as a useful proxy to quantify feedback in a dialogue.

- (1) A: and uh it really does irk me to see those guys out there uh you know making that ///much money///

B: ///yeah///<sup>1</sup>

Recent works have emphasized the role of feedback and grounding signals in their study of human-human conversations (Fusaroli et al., 2017; Dideriksen et al., 2022; Dingemanse and Liesenfeld, 2022) as well as human-agent interaction (Visser et al., 2014; Hough and Schlangen, 2016; Buschmeier and Kopp, 2018; Axelsson et al., 2022).

The literature tends to merge the two closely related concepts of *backchannels* and *acknowledgments*. Backchannels (Yngve, 1970), or *continuers* (Schegloff, 1982), are not positioned on the main channel, but uttered by the “listener”, often as low intensity unobtrusive overlapping speech (Heldner et al., 2010) or non-verbally (Allwood et al., 2007; Truong et al., 2011). Acknowledgments, on the other hand, have a slightly broader, functional definition of minimal positive feedback (Jefferson, 1984; Allwood et al., 1992).

There is a large body of work on lexical markers, also called *cue phrases* or *discourse markers* (Jefferson, 1984; Allwood et al., 1992; Muller and Prévot, 2003), since they present interesting linguistic features and constitute convenient explicit cues for detecting feedback utterances automatically (Jurafsky et al., 1998; Gravano et al., 2012; Prévot et al., 2015). Gravano et al. (2012) developed a list of affirmative cue words made of *alright*, *mm-hm*, *okay*, *right*, *uh-huh*, *yeah*. Form-Function studies of similar lists have been made at least for Swedish (Allwood, 1988), U.S. English (Ward, 2006), and French (Prévot et al., 2015).

Few studies have, however, concentrated on direct negative feedback associated with rejection and corrective dialogue acts. Although Allwood et al. (1992) suggests a polarity dimension for characterizing feedback, most recent studies have focused on positive feedback. Indeed, in collaborative dialogue and everyday conversations, which are the two genres dominating available datasets, positive feedback constitutes the large majority of explicit feedback (e.g., Malisz et al., 2016). Negative feedback is instead often expressed constructively, using repair mechanisms, specifically *clarification requests* (Purver, 2004). These may rely on simple lexical cues (e.g., for English, *pardon?*, *huh?*), sluices (such as *what?*, *who?*), or on clarification ellipsis, as in the following example (Fernández et al., 2007):

<sup>1</sup>Notation: ///text/// produced in overlap with the speech of the other speaker. From Switchboard (Godfrey et al., 1992)

- (2) **A:** and then we're going to turn east  
**B:** mmhmm  
**A:** not straight east slightly sort of northeast  
**B:** slightly northeast?<sup>2</sup>

The occurrence of feedback signals in dialogue transcriptions can be detected using various types of sequence labeling models from classical hidden Markov models (Stolcke et al., 2000b) to modern neural architectures and large language models (Liu et al., 2017; Noble and Maraev, 2021).

## 2.2 Analysis of Subtitles

Subtitles are typically short written text snippets and they accompany audiovisual content on the screen. They are often subject to condensation and normalization, where non-standard verbal elements (repetitions, signs of hesitation etc.) are omitted or replaced by more standard alternatives (Gottlieb, 2012) due to constraints on the length, readability and writing conventions. As subtitles are displayed alongside audiovisual content, viewers can typically recover omitted dialogue-relevant cues from the accompanying images and sounds. *Interlingual subtitling* – where the original language of the audio is different from the subtitling language – differs somewhat from *intralingual subtitling*, which is meant for same-language audio and subtitles which also records non-verbal elements writing for the benefit of hearing impaired audiences or non-native speakers (Gottlieb, 2012).

Rühlemann (2020) compared real conversations and scripted ones and observed that continuers were absent from the latter. Prevot et al. (2019) compared data from the *Open Subtitles* corpus (Lison and Tiedemann, 2016; Lison et al., 2018) in English, French and Mandarin with both written and conversational corpora and found that OpenSubtitles occupied an intermediate position between written and conversational data in terms of lexical and syntactic features. This paper builds upon those earlier works but focuses specifically on communicative feedback, using a combination of lexical statistics, manual and automate annotations to quantify its frequency in various corpora.

## 3 Corpora

We rely on data from both OpenSubtitles and existing, publicly available corpora of real conversations covering eight different languages (see Table 1).

<sup>2</sup>From HCRC Map Task (Anderson et al., 1991).

## 3.1 Spontaneous Dialogues

**German (de)** We use the Hamburg MapTask corpus (HZSK, 2010), in which twelve dyads of (L2) speakers of German engage in dyadic task-oriented short dialogues.

**English (en)** For English, we use Switchboard (SWBD), consisting of dyadic topic oriented phone conversation (Godfrey et al., 1992) as well as Fisher (Cieri et al., 2004) for some experiments; AMI, with multi-party multimodal task-oriented dialogues (Carletta, 2007); HCRC MapTask (MT) comprising dyadic task-oriented short dialogues (Anderson et al., 1991); and STAC, a multi-party negotiation chat corpus (Asher et al., 2016).

**French (fr)** We include CID, consisting of dyadic, 1-hour long, loosely topic-oriented face-to-face conversations with 16 participants (Blache et al., 2017); French MapTask with 16 participants (Gorisch et al., 2014); and Aix-DVD, dyadic face-to-face conversations about movie preferences of 16 participants (Prévot et al., 2016).

**Hungarian (hu)** We employ BUSZI-2 corpus (Budapest Sociolinguistic Interview, Váradi, 2003), where 50 participants with different educational levels participated in a 30-minute directed conversation and then performed language tasks (e.g. grammaticality judgments).

**Italian (it)** We use the CLIPS corpus (Savy and Cutugno, 2009), consisting of both a map task and a difference spotting task between images. We exclude dialogues with a high proportion (> 10%) of utterances with dialectal words.

**Japanese (ja)** This language is represented by the transcripts of the CallHome Japanese corpus (Den and Fry, 2000) consisting of 120 unscripted telephone conversations between native speakers, mostly family members or close friends.

**Norwegian (no)** We use the NoTa-Oslo corpus (Johannessen et al., 2007), containing interviews and conversations from 2004–2006 with 166 informants from the Oslo area. The dialogues consist of 10-minute semi-formal interviews and 30-min informal dialogues with other informants.

**Mandarin Chinese (zh)** The source of our Mandarin Chinese data was CALLHOME (Wheatley, 1996) consisting of unscripted telephone conversations between native speakers.

Language	de	en	fr	hu	it	ja	no	zh	total
# Spontaneous dialogues	24	2766	48	50	88	120	259	120	<b>3475</b>
# Utterances	4K	373K	27K	31K	24K	39K	86K	18K	<b>602K</b>
# Subtitles	98	100	100	68	95	74	87	93	<b>715</b>
# Utterances	131K	140K	126K	93K	138K	106K	98K	114K	<b>946K</b>

Table 1: Overview of dialogue data sources for both spontaneous conversations and subtitles employed in this paper.

### 3.2 Subtitles

The scripted dialogues are extracted from [Open-Subtitles 2018](#) (Lison et al., 2018), a large collection of over 3.7 million subtitles (amounting to  $\approx$  22.1 billion words) extracted from the Open-Subtitles.org database and covering 60 languages. We include both (1) subtitles for the hearing impaired, where the subtitle language and the original audio language are identical and (2) subtitles for foreign audiences. The subtitles are then filtered according to several criteria. Only recent movies (year  $\geq$  1990) are included to reflect contemporary language use, as is the case for the corpora of spontaneous conversations. We also omit subtitles with less than 100 utterances and exclude genres that are less relevant for this study (Documentary, Reality-TV, Biography, Sport, Musical, Music, Adult, Animation, Short and Game-Show).

We sample up to ten movies per audience type (hearing impaired vs. foreign audience) from the five largest genres, namely drama, comedy, crime, action, and romance. Table 1 shows the number of movies and utterances per language for the selected subtitles. Note that subtitles are typically segmented by dialogue turns or sentences instead of utterances. The term “utterance” should therefore be understood broadly in this paper.

This paper focuses on the textual aspects of grounding phenomena. While speech and non-linguistic aspects of communicative feedback (such as timing, intonation, gestures or gaze) are both important and well-studied, in particular for acknowledgements and backchannels, those information are not available in subtitles corpora, which are intrinsically limited to text transcriptions.

## 4 Lexical Analysis

Lexical statistics of acknowledgment cues gives us a first picture of the feedback frequency. Acknowledgments tend to be produced by the addressee (not the main speaker) and are therefore often short productions uttered in overlap and potentially with a lower voice. Out of those three properties (brevity,

overlap, lower volume), only the first is practically measurable in our experiments, as the subtitles are by construction text-based.

Given their relation to acknowledgments, we first analyse “very short utterances” (Edlund et al., 2009), defined here as three tokens or less. Feedback is also very well represented at initial positions of longer turns/contributions. We therefore targeted two locations: *very short utterances* (all tokens) and *initial positions* (one token) of all other utterances. Comparing term frequencies between these locations and the overall corpus allowed us to compile language-specific lists of *cue words*. Those lists of cue words (presented in Table 3 in the Appendix) are divided into four core classes of feedback:

- positive feedback/acknowledgment (+)
- neutral/continuer (=)
- negative feedback (-)
- clarification request (?).

We plot in Figure 1 the frequencies of those feedback classes in each corpus, either in terms of absolute frequency (left side) or by looking at the relative proportions of the feedback classes (right side). Figure 2 shows the lexical distribution of the most frequent lexical items observed in the utterances of plot (b) for English.

We observe that the statistics based on cue words differ substantially between subtitles and spontaneous dialogues. This difference is observed across all languages and sub-genres, (see Appendix A for other languages). We sought to identify and reduce other sources of variation between corpora. STAC, as a chat corpus, exhibits different patterns than other dialogue corpora, notably due to the presence of emojis. Similarly, for English and French, we explored the impact of politeness expression (highly frequent in OpenSubtitles). Those peculiarities did not, however, change the overall picture of our analysis (see Figure 13 in Appendix A).

One key difference between real dialogues and subtitles relates to the overall frequency of feedback cues, which is much higher in spontaneous

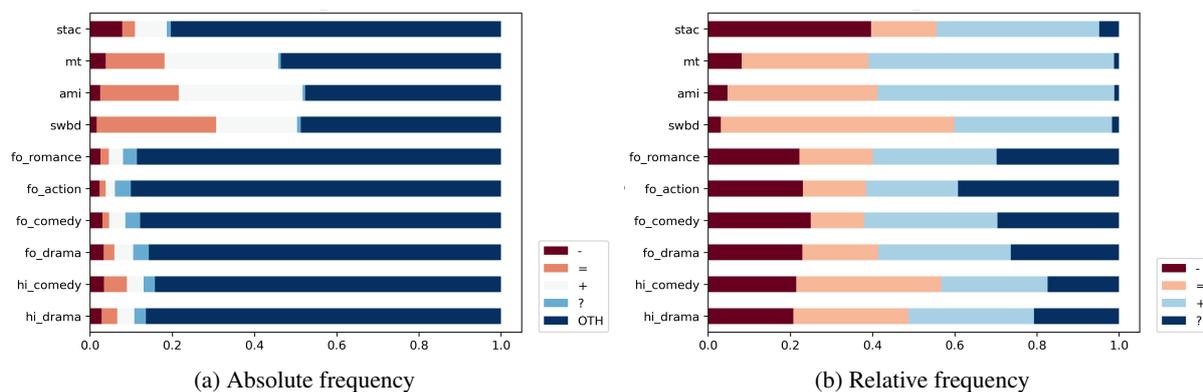


Figure 1: Frequency of conversational feedback of various types among utterances in the English corpora (both spontaneous and subtitles) based on manually curated lists of cue words to detect. Fig. (a) shows the absolute frequency while Fig. (b) zooms in on utterances labelled with at least one feedback. + denotes positive feedback/acknowledgement, = neutral/continuer feedback, - negative feedback, ? clarification requests and 'OTH' is for other utterances. fo and hi respectively stand for 'foreign audience' and 'hearing-impaired' subtitles. Corpora without these prefixes are spontaneous dialogues.

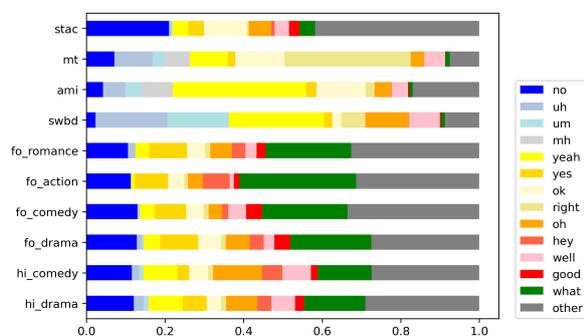


Figure 2: Most common lexical items associated with communicative feedback, as detected through manually curated lists of cue words in English, factored by corpus.

dialogues (40–50%) than in subtitles (10–20%), as observed in figure 1(a). Furthermore, as shown in Figure 1(b), feedback in spontaneous dialogues consists mostly in positive or neutral (*continuers*) feedback, while subtitles have few neutral signals but seem to exhibit a much higher proportion of negative feedback and clarification requests.

We compared our English cue word lists against the annotations in Switchboard. After grouping feedback-related labels into a single *Feedback* category, we find that the cue word lists yield an  $F_1$  score of 0.76.

## 5 Dialogue Act Tagging

Although lexical statistics do highlight substantial differences in subtitles and spontaneous dialogues, they remain imprecise estimates, as many cue words related to feedback tend to be ambiguous. In this section, we refine our analysis using a dialogue act tagging model trained on the DAMSL-

Switchboard corpus.

### 5.1 Data

We map the original set of Switchboard (SWBD) tags, and their clustered DAMSL-SWBD equivalents, into five coarse dialogue act (DA) classes: *Forward looking*, *Yes/no answers*, *Assessment*, *Backchannel* and *Other*. The two classes most directly relevant for feedback, namely *Backchannel* and *Assessment*, are inspired, in part, by Mezza et al. (2018). Distinguishing between these two feedback-related classes is also motivated by Goodwin (1986), who outline a number of positional and functional differences between these. The *Backchannel* category consists of the SWBD-DAMSL labels<sup>3</sup> *Acknowledge (Backchannel)*, (SWBD tag b), *Backchannel in question form (bh)*, *Response Acknowledgment (bk)*, *Summarize/reformulate (bf)* and *Signal-non-understanding (br)*. As this latter tag suggests, negative feedback signals are also part of the *Backchannel* category, since they are too few to reliably learn a separate class from. The *Assessment* category comprises not only the labels *Agree/Accept (aa)*, but also *Appreciation (ba)* and *Exclamation (fe)*. The forward looking category contains utterances expressing explanations, instructions and suggestions as well as questions. Table 4 in Appendix B shows the distribution of instances per label and their SWBD tag.

<sup>3</sup>[web.stanford.edu/jurafsky/ws97/manual.august1.html](http://web.stanford.edu/jurafsky/ws97/manual.august1.html)

DA group	Data	de	en	fr	hu	it	ja	no	zh
<b>Assessment</b>	<b>SPCONV</b>	16.50	9.11	4.62	15.49	12.64	15.74	17.05	6.96
	<b>SUBS</b>	9.08	7.07	7.72	9.29	8.34	6.48	6.53	5.00
<b>Backchannel</b>	<b>SPCONV</b>	11.57	10.79	11.96	4.28	5.73	18.96	2.67	5.65
	<b>SUBS</b>	3.49	3.72	3.44	3.48	3.45	3.74	3.47	3.00
<b>Yes/no answer</b>	<b>SPCONV</b>	2.22	1.15	1.24	4.00	6.55	2.84	5.09	1.00
	<b>SUBS</b>	1.97	1.37	1.68	1.47	1.38	1.15	2.32	0.76

Table 2: Proportions (%) of the relevant dialogue act groups detected by the BERT-based dialogue act tagger in the spontaneous conversation (SPCONV) and in the subtitle (SUBS) corpora.

## 5.2 Model Training

We fine-tune the monolingual bert-base-cased pre-trained model (Devlin et al., 2019) using 80% of the Switchboard data as training and 20% for development and testing. We set up the task as a sequence classification problem, including the preceding utterance as context. We train the model with a batch size of 8, a learning rate of  $4E-5$  and default values for the other parameters. We run and compare three different random seeds, yielding similar performance. To improve recall, we also adjust the probability thresholds for the feedback classes.

The model performs relatively well on the Switchboard test set, yielding an accuracy of 0.81. The  $F_1$  scores for the *Assessment* and *Backchannel* classes are respectively 0.59 and 0.83. This score difference may be due to *Backchannel* instances being better represented in the training data, as well as some label confusion between the *Assessment* and the *Yes/No question* categories.

## 5.3 Empirical Results

We then use the trained dialogue act tagger to detect conversational feedback signals in both the spontaneous dialogue and subtitles. For non-English corpora, we machine translate the data using the Google Translate API. Feedback-annotated conversational corpora is non-existent for most languages and the quality of current MT systems is generally considered high enough to serve as a viable alternative (Isbister et al., 2021).

Table 2 presents the empirical results obtained with our dialogue act tagger on both spontaneous dialogues and subtitle corpora. We observe that backchannels are considerably more frequent (by a factor three) in spontaneous dialogues than in subtitles for half of the languages – which is in line with the results of our lexical analysis in Section 4. The number of utterances labeled as *Assessment*

differs less, but subtitles still seem to contain less of this feedback type in almost all genres and languages except French (see Appendix B for details). Given that the tagger is only trained on a single corpus, some of the differences found may also be attributed to the generalization ability of the tagger to certain domains. We therefore also conduct some manual error analysis.

## 5.4 Error Analysis

In general, the proportion of the *Backchannel* category for the spontaneous conversations is lower for Hungarian, Italian, Norwegian and Mandarin than for the other languages. This is likely due to the use of infrequent spelling variants of backchannel signals such as *hmm*, *mh*. We have also found that the tagger has difficulties detecting feedback when they are part of longer utterances, whether they appear in an utterance-initial position or not. We also observe a general tendency to associate sentence-final question marks to feedback cues. When inspecting the most frequent utterances tagged as feedback, we also notice that short utterances pose some challenges for machine translation due to polysemy, e.g., *Cosa?* “Thing?”, also translatable as “What?”, in Italian.

## 6 Further Annotations

The results from the dialogue tagger do show some clear trends regarding the extent to which communicative feedback is expressed in subtitles compared to spontaneous interactions. However, the use of DAMSL-Switchboard as sole source of training data is a limiting factor in our analysis, in particular when it comes to non-English dialogues, which must be machine-translated prior to labelling. Furthermore, the tagger does not provide information about the frequency of negative feedback, although the lexical analysis from Section 4 does seem to

point towards a higher frequency of those communicative signals in subtitles.

We therefore complement the analyses of the two previous sections with a manual annotation effort. To this end, we sample from each corpus a set of 300 utterances to annotate. However, as evidenced by the results of the previous sections, many utterances of our corpora do not seem to contain any communicative feedback. To ensure the annotation process can cover a sufficiently broad variety of feedback signals despite this class imbalance, we do not select the utterances purely at random, but select half among those marked as feedback-relevant by the cue words of Section 4, and the other half among those that do not.

### Annotation Process

We recruited 6 annotators with prior expertise in linguistic annotation and proficient in the language corresponding to the corpus to annotate. Those annotators were provided each utterance in its context, and were tasked to decide whether the utterance in question contains one of the following three categories of communicative feedback: defined in the annotation guidelines as such:

**AGREE\_ACCEPT** : indicates that the speaker agrees or accepts what has been said.

**ACK\_BACK** : indicates that the speaker is listening to her interlocutor, or at least heard what has been said, without necessarily agreeing with it or committing to its content.

**NEGATIVE\_FEEDBACK** : indicates that the speaker could not hear or understand her interlocutor, or even rejects or disagrees with what the other person has said.

Answers to explicit questions should not be considered as feedback. Each utterance can be tagged with zero, one, or multiple feedback labels. These categories specifically target and distinguish between different conversational feedback phenomena and are therefore somewhat more comprehensive than the categories employed by the tagger of the previous section. There, similar categories were derived by merging the available feedback-relevant dialogue act labels from the SWBD annotations.

A total of 24 corpus samples, each comprising 300 utterances, were annotated<sup>4</sup>. Three corpus samples (respectively for English, French and

<sup>4</sup>The full set of annotated dialogue samples is available at [https://github.com/NorskRegnesentral/conv\\_feedback](https://github.com/NorskRegnesentral/conv_feedback).

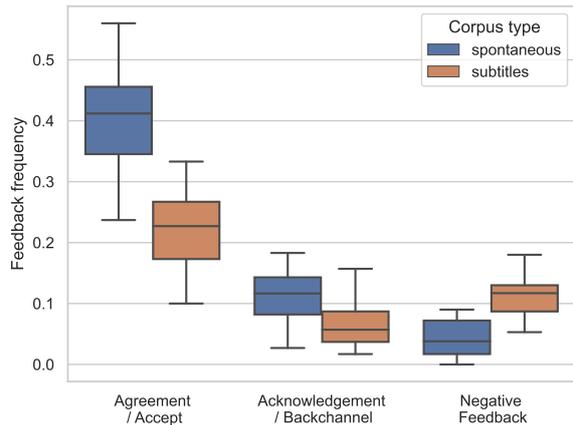


Figure 3: Frequency of communicative feedback depending on the source of the dialogue sample (spontaneous interactions or subtitles) and the category of feedback, based on annotations from human experts.

Chinese) were doubly annotated, and the Kappa’s score of their agreement was found to be 0.59 for *AGREE\_ACCEPT*, 0.42 for *ACK\_BACK* and 0.54 for *NEGATIVE\_FEEDBACK* across the 3 samples. This relatively low inter-annotator agreement illustrates the challenging nature of the annotation task, in particular due to the lack of explicit turn boundaries in subtitles, making it at times difficult to determine the context behind each utterance.

### 6.1 Annotation Results

Figure 3 illustrates the frequencies of the three feedback categories across the 24 annotated samples. We observe again a lower proportion of both *Agree / Accept* and *Acknowledgement / Backchannel* feedbacks in the subtitles compared to real interactions. The proportion of *Negative feedback* is, however, higher for the subtitles. We hypothesise that this may stem from the fact that disagreements between interlocutors are more interesting from the storytelling perspective, and are therefore more common in subtitles than in real interactions.

We investigated whether subtitles for foreign audiences differed from subtitles written for the hearing impaired (as those often need to adhere more closely to the original on-screen conversation), but did not find any substantial disparity.

### 6.2 LLM-based Annotation

The frequencies of Figure 3 are obtained using the manually annotated dialogue samples. However, those samples only cover a small fraction of available corpora. Furthermore, as the sampling procedure relied on the use of cue-words to cover a

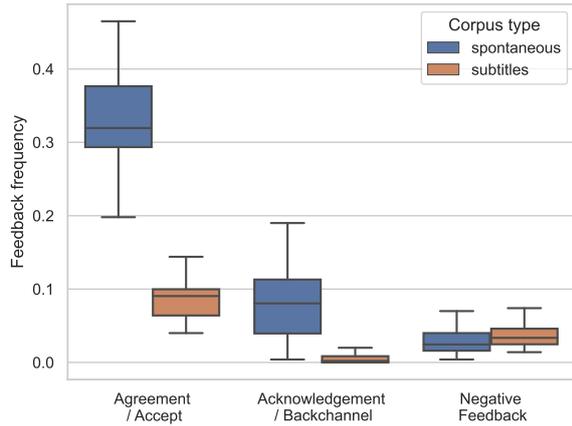


Figure 4: Frequency of communicative feedback depending on the corpus type and category of feedback, based on the predictions of the fine-tuned Gemma 2 model trained on human annotations.

sufficiently broad set of feedback types (see above), it is likely to overestimate the proportion of communicative feedbacks. To mitigate this bias, we fine-tune an instruction-tuned Gemma 2 model (Gemma Team et al., 2024) to predict the probability of an utterance including one of the three defined feedback categories. The fine-tuning relied on LoRA (Hu et al., 2021) and included as instructions the annotation guidelines also provided to the human experts. The full set of 24 dialogue samples was used for the fine-tuning, each utterance being provided in its local dialogue context. For non-English utterances, we also include in the prompt an English translation of the utterance and its context, obtained using Google Translate.

The fine-tuned Gemma2 LLM was then applied to all corpora to predict whether their utterances contained one of the three categories of feedback defined above. The results are shown in Figure 4. The proportions of communicative feedback are somewhat lower in the actual corpora than in the annotated samples (which is expected given how the dialogue samples were derived), but the overall trends remain similar to Figure 3.

## 7 Conversational Feedback in Synthetic Dialogues

We conclude by investigating the occurrence of communicative feedback in synthetic dialogues generated with autoregressive language models. More precisely, we wish to analyse whether the communicative feedback generated by those models are closer to the patterns found in real interac-

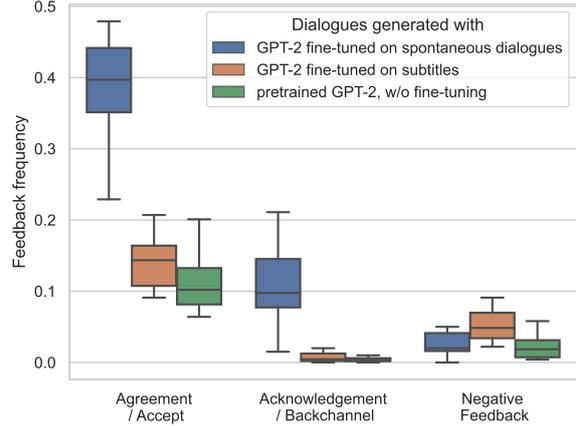


Figure 5: Frequency of communicative feedback in synthetic dialogues generated using GPT-2 models, either applied without fine-tuning or after fine-tuning on corpora of spontaneous interactions or subtitles.

tions or to scripted dialogues such as subtitles.

To this end, we use available GPT-2 models (Radford et al., 2019) for the eight covered languages<sup>5</sup>. The use of GPT-2 models is motivated by practical considerations and the need to obtain pre-trained models for each of the eight languages. For each corpus, we derive a fine-tuned version of its corresponding GPT-2 model by further training the model on the corpus dialogues. To account for the corpus size differences, the number of epochs is adjusted to ensure that the total number of gradient updates is similar across all corpora.

The GPT-2 models are then employed to produce synthetic dialogues (100 dialogues of about 50 turns per model). For the fine-tuned models, all turns are automatically generated, while for the base models, the following dialogue start is used as context: *Hi! – Hi, how are you? – Fine, and you?* to bias the model towards the generation of dialogues. Finally, the LLM annotator from the previous section is applied on those synthetic dialogues to estimate their frequency of communicative feedback.

The results are shown in Figure 5. We observe that the synthetic dialogues generated with the standard GPT-2 models without any further fine-tuning are much closer to the ones derived from subtitles than to those derived from spontaneous interactions when it comes to communicative feedback. This is

<sup>5</sup>The following pre-trained models are employed: gpt2-base (English), gpt-fr-cased-small (French), german-gpt2 (German), gpt2-small-italian (Italian), PULI-GPT-2 (Hungarian), norwegian-gpt2 (Norwegian), gpt2-chinese-cluecorpussmall (Mandarin Chinese), and japanese-gpt2-medium (Japanese).

notably the case for positive and neutral feedback. The occurrence of negative feedback is, however, not as common as in subtitles. Although the above results were obtained here using only GPT-2 pre-trained models, we expect to find similar patterns for other (and more recent) LLMs.

## 8 Conclusion and Future Work

As evidenced in this paper, movie and TV subtitles exhibit notable linguistic differences to actual spontaneous dialogues in the amount and type of conversational feedback they include. Based on a collection of corpora of both spontaneous dialogues and subtitles across eight languages, we provide both lexical statistics and dialogue act estimates derived with a fine-tuned dialogue act tagger. We show that the proportion of conversational feedback is considerably lower in subtitles than in spontaneous dialogues across the corpora included. Furthermore, the type of conversational feedback also differs, as negative feedback is proportionally more frequent in subtitles. This is corroborated by manual annotations of 24 dialogue samples from the selected corpora, and the use of a fine-tuned LLM trained on those annotations. Finally, we also show that dialogues generated from language models are closer to scripted dialogue than real interactions in their use of communicative feedback. Beyond their linguistic interest, these results can provide useful insights for the development of conversational models, as those are often trained on scripted dialogues and might therefore struggle both to understand communicative feedback from the user and to produce such feedback themselves.

## Acknowledgments

This work was carried out within the *GraphDial* project (<https://graphdial.nr.no/>), supported by the Research Council of Norway. The work was initiated during the bilateral PHC-Aurora project *French Norwegian Research Effort on Applied Dialogue Modelling*. Laurent Prévot would also like to acknowledge continuous support from the Institute for Language Communication and the Brain (ILCB ANR-16-CONV0002) and the Excellence Initiative of Aix-Marseille University (A\*MIDEX AAP-ID-17-46-170301-11.1). Hendrik Buschmeier was supported by the German Research Foundation (DFG) in the Collaborative Research Center TRR 318/1 2021 ‘Constructing Explainability’ (438445824). We would like to thank Hiro Yamazaki for helping

with the Japanese lexical items.

## References

- Jens Allwood. 1988. Om det svenska systemet för språklig återkoppling. In Per Linell, Viveka Adelswärd, Torbjörn Nilsson, and Per A. Pettersson, editors, *Svenskans Beskrivning 16*, volume 1, pages 89–106. Linköping University, Tema Kommunikation, Linköping, Sweden.
- Jens Allwood, Stefan Kopp, Karl Grammer, Elisabeth Ahlsén, Elisabeth Oberzaucher, and Markus Koppensteiner. 2007. [The analysis of embodied communicative feedback in multimodal corpora: A prerequisite for behaviour simulation](#). *Language Resources and Evaluation*, 41:255–272.
- Jens Allwood, Joakim Nivre, and Elisabeth Ahlsén. 1992. [On the semantics and pragmatics of linguistic feedback](#). *Journal of Semantics*, 9:1–26.
- Anne H. Anderson, Miles Bader, et al. 1991. The HCRC map task corpus. *Language and Speech*, 34:351–366.
- Nicholas Asher, Julie Hunter, Mathieu Morey, Farah Benamara, and Stergos Afantenos. 2016. [Discourse structure and dialogue acts in multiparty dialogue: The STAC corpus](#). In *Proceedings of the 10th International Conference on Language Resources and Evaluation*, pages 2721–2727.
- Agnes Axelsson, Hendrik Buschmeier, and Gabriel Skantze. 2022. [Modelling feedback in interaction with conversational agents – a review](#). *Frontiers in Computer Science*, 4:744574.
- Todd Berliner. 1999. Hollywood movie dialogue and the "real realism" of John Cassavetes. *Film Quarterly*, 52(3):2–16.
- Philippe Blache, Roxane Bertrand, Gaëlle Ferré, Berthille Pallaud, Laurent Prévot, and Stéphane Rauzy. 2017. The corpus of interactional data: A large multimodal annotated resource. *Handbook of Linguistic Annotation*, pages 1323–1356.
- Harry Bunt. 1994. Context and dialogue control. *Think Quarterly*, 3(1):19–31.
- Hendrik Buschmeier and Stefan Kopp. 2018. [Communicative listener feedback in human-agent interaction: Artificial speakers need to be attentive and adaptive](#). In *Proceedings of the 17th International Conference on Autonomous Agents and Multiagent Systems*, pages 1213–1221, Stockholm, Sweden.
- Jean Carletta. 2007. Unleashing the killer corpus: experiences in creating the multi-everything AMI Meeting Corpus. *Language Resources and Evaluation*, 41:181–190.
- Neda Chepinchikj and Celia Thompson. 2016. [Analysing cinematic discourse using conversation analysis](#). *Discourse, Context & Media*, 14:40–53.

- Christopher Cieri, David Miller, and Kevin Walker. 2004. The Fisher corpus: A resource for the next generations of speech-to-text. In *LREC*, volume 4, pages 69–71.
- Herbert H. Clark. 1996. *Using Language*. Cambridge University Press, Cambridge, UK.
- Herbert H. Clark and Edward F. Schaefer. 1989. [Contributing to discourse](#). *Cognitive Science*, 13:259–294.
- Cristian Danescu-Niculescu-Mizil and Lillian Lee. 2011. Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs. In *Proceedings of the 2nd Workshop on Cognitive Modeling and Computational Linguistics*, pages 76–87, Portland, OR, USA.
- Mark Davies. 2021. The TV and movies corpora: Design, construction, and use. *International Journal of Corpus Linguistics*, 26:10–37.
- Yasuharu Den and John Fry. 2000. CallHome Japanese corpus (in Japanese). *Journal of the Phonetic Society of Japan*, 4(2):24–30.
- 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, MN, USA.
- Christina Dideriksen, Morten H Christiansen, Kristian Tylén, Mark Dingemanse, and Riccardo Fusaroli. 2022. Quantifying the interplay of conversational devices in building mutual understanding. *Journal of Experimental Psychology: General*.
- Mark Dingemanse and Andreas Liesenfeld. 2022. [From text to talk: Harnessing conversational corpora for humane and diversity-aware language technology](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5614–5633, Dublin, Ireland. Association for Computational Linguistics.
- Mark Dingemanse, Seán G. Roberts, Julija Baranova, Joe Blythe, Paul Drew, Simeon Floyd, Rosa S. Gisladottir, Kobin H. Kendrick, Stephen C. Levinson, Elizabeth Manrique, Giovanni Rossi, and Nick J. Enfield. 2015. [Universal principles in the repair of communication problems](#). *PLoS ONE*, 10:e0136100.
- Jens Edlund, Mattias Heldner, and Antoine Pelcé. 2009. Prosodic features of very short utterances in dialogue. In Vainio Martti, Aulanko Reijo, and Olli Aaltonen, editors, *Nordic Prosody. Proceedings of the Xth Conference, Helsinki 2008*, pages 57–68. Peter Lang, Frankfurt am Main, Germany.
- Raquel Fernández, Jonathan Ginzburg, and Shalom Lappin. 2007. [Classifying Non-Sentential Utterances in Dialogue: A Machine Learning Approach](#). *Computational Linguistics*, 33(3):397–427.
- Riccardo Fusaroli, Kristian Tylén, Katrine Garly, Jakob Steensig, Morten H Christiansen, and Mark Dingemanse. 2017. Measures and mechanisms of common ground: Backchannels, conversational repair, and interactive alignment in free and task-oriented social interactions. In *the 39th Annual Conference of the Cognitive Science Society (CogSci 2017)*, pages 2055–2060. Cognitive Science Society.
- 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 Mikula, 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 Klimentko, 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](#).
- John J Godfrey, Edward C Holliman, and Jane McDaniel. 1992. Switchboard: Telephone speech corpus for research and development. In *Acoustics, Speech, and Signal Processing, IEEE International Conference on*, volume 1, pages 517–520. IEEE Computer Society.
- Charles Goodwin. 1986. Between and within: Alternative sequential treatments of continuers and assessments. *Human studies*, 9(2-3):205–217.
- Jan Gorisch, Corine Astésano, Ellen Gurman Bard, Brigitte Bigi, and Laurent Prévot. 2014. [Aix map](#)

- task corpus: The French multimodal corpus of task-oriented dialogue. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 2648–2652, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Henrik Gottlieb. 2012. Subtitles: readable dialogue? *Eye tracking in audiovisual translation*, pages 37–82.
- Agustín Gravano, Julia Hirschberg, and Štefan Beňuš. 2012. Affirmative cue words in task-oriented dialogue. *Computational Linguistics*, 38(1):1–39.
- Mattias Heldner, Jens Edlund, and Julia Hirschberg. 2010. Pitch similarity in the vicinity of backchannels. In *Proceedings of INTERSPEECH 2010*, pages 3054–3057, Makuhari, Japan.
- Mattias Heldner, Anna Hjalmarsson, and Jens Edlund. 2013. Backchannel relevance spaces. In *Proceedings of Nordic Prosody XI*, pages 137–146, Tartu, Estonia.
- Julian Hough and David Schlangen. 2016. Investigating fluidity for human-robot interaction with real-time, real-world grounding strategies. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 288–298, Los Angeles, CA, USA.
- Christine Howes and Arash Eshghi. 2021. Feedback relevance spaces: Interactional constraints on processing contexts in Dynamic Syntax. *Journal of Logic, Language and Information*, 30(2):331–362.
- 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.
- HZSK. 2010. HAMATAC. The Hamburg MapTask Corpus.
- Tim Isbister, Fredrik Carlsson, and Magnus Sahlgren. 2021. Should we stop training more monolingual models, and simply use machine translation instead? In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 385–390, Reykjavik, Iceland (Online). Linköping University Electronic Press, Sweden.
- Gail Jefferson. 1972. Side Sequences. *Studies in social interaction*.
- Gail Jefferson. 1984. Notes on a systematic deployment of the acknowledgement tokens “Yeah”; and “Mm Hm”. *Paper in Linguistics*, 17(2):197–216.
- Janne Bondi Johannessen, Kristin Hagen, Joel James Priestley, and Lars Nygaard. 2007. An advanced speech corpus for Norwegian. In *Proceedings of the 16th Nordic Conference of Computational Linguistics (NODALIDA 2007)*, pages 29–36, Tartu, Estonia. University of Tartu, Estonia.
- Dan Jurafsky, Elizabeth Shriberg, Barbara Fox, and Traci Curl. 1998. Lexical, prosodic, and syntactic cues for dialog acts. In *Discourse Relations and Discourse Markers*.
- Göran Kjellmer. 2009. Where do we backchannel?: On the use of mm, mhm, uh huh and such like. *International Journal of Corpus Linguistics*, 14(1):81–112.
- Andreas Liesenfeld and Mark Dingemanse. 2022. Bottom-up discovery of structure and variation in response tokens (‘backchannels’) across diverse languages. In *Proceedings of Interspeech 2022*, pages 1126–1130, Incheon, Korea.
- 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).
- Pierre Lison, Jörg Tiedemann, and Milen Kouylekov. 2018. OpenSubtitles2018: Statistical Rescoring of Sentence Alignments in Large, Noisy Parallel Corpora. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Yang Liu, Kun Han, Zhao Tan, and Yun Lei. 2017. Using context information for dialog act classification in DNN framework. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2170–2178, Copenhagen, Denmark. Association for Computational Linguistics.
- Zofia Malisz, Marcin Włodarczak, Hendrik Buschmeier, Joanna Skubisz, Stefan Kopp, and Petra Wagner. 2016. The alico corpus: Analysing the active listener. *Language Resources and Evaluation*, 50:411–442.
- Stefano Mezza, Alessandra Cervone, Evgeny Stepanov, Giuliano Tortoreto, and Giuseppe Riccardi. 2018. ISO-standard domain-independent dialogue act tagging for conversational agents. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3539–3551, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Philippe Muller and Laurent Prévot. 2003. An empirical study of acknowledgment structures. In *7th workshop on the semantics and pragmatics of dialogue*.
- Bill Noble and Vladislav Maraev. 2021. Large-scale text pre-training helps with dialogue act recognition, but not without fine-tuning. In *Proceedings of the 14th International Conference on Computational Semantics (IWCS)*, pages 166–172, Groningen, The Netherlands (online). Association for Computational Linguistics.
- Tim Paek and Eric Horvitz. 2000. Grounding criterion: Toward a formal theory of grounding. Technical Report MSR-TR-2000-40, Microsoft Research, Redmond, WA, USA.

- Laurent Prévot, Jan Gorisch, and Roxane Bertrand. 2016. [A CUP of CoFee: A large collection of feedback utterances provided with communicative function annotations](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3180–3185, Portorož, Slovenia. European Language Resources Association (ELRA).
- Laurent Prévot, Jan Gorisch, and Sankar Mukherjee. 2015. [Annotation and classification of French feedback communicative functions](#). In *Proceedings of the 29th Pacific Asia Conference on Language, Information and Computation*, pages 298–306, Shanghai, China.
- Laurent Prévot, Pierre Magistry, and Pierre Lison. 2019. [Should we use movie subtitles to study linguistic patterns of conversational speech? a study based on French, English and Taiwan Mandarin](#). In *Third International Symposium on Linguistic Patterns of Spontaneous Speech*.
- Matthew Purver. 2004. *The Theory and Use of Clarification Requests in Dialogue*. Ph.D. thesis, King's College, University of London, London, UK.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). Technical report, OpenAI.
- Christoph Rühlemann. 2020. [What dialog is absent from constructed dialog?](#) *English Text Construction*, 13:132–157.
- Renata Savy and Francesco Cutugno. 2009. [CLIPS diatopic, diamesic and diaphasic variations in spoken Italian](#). In *Proceedings of the 5th Corpus Linguistics Conference: CL2009*.
- Emanuel A Schegloff. 1982. Discourse as an interactional achievement: Some uses of 'uh huh' and other things that come between sentences. *Analyzing discourse: Text and talk*, 71:71–93.
- Emanuel A Schegloff. 2000. Overlapping talk and the organization of turn-taking for conversation. *Language in society*, 29(1):1–63.
- Elizabeth Shriberg. 1996. Disfluencies in switchboard. In *Proceedings of the International Conference on Spoken Language Processing*, volume 96, pages 11–14.
- Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000a. [Dialogue act modeling for automatic tagging and recognition of conversational speech](#). *Computational Linguistics*, 26:339–373.
- Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000b. [Dialogue Act Modeling for Automatic Tagging and Recognition of Conversational Speech](#). *Computational Linguistics*, 26(3):339–373.
- David R. Traum. 1994. *A Computational Theory of Grounding in Natural Language Conversation*. Ph.D. thesis, University of Rochester, Rochester, NY, USA.
- Khiet P. Truong, Ronald Poppe, Iwan de Kok, and Dirk Heylen. 2011. [A multimodal analysis of vocal and visual backchannels in spontaneous dialogs](#). In *Proceedings of INTERSPEECH 2011*, pages 2973–2976, Florence, Italy.
- Thomas Visser, David R. Traum, David DeVault, and Rieks op den Akker. 2014. [A model for incremental grounding in spoken dialogue systems](#). *Journal on Multimodal User Interfaces*, 8:61–73.
- Tamás Váradi. 2003. [A Budapesti Szociolingvisztikai Interjú](#). In Ferenc Kiefer, editor, *A magyar nyelv kézikönyve*, pages 339–360. Akadémiai Kiadó, Budapest.
- Nigel Ward. 2006. [Non-lexical conversational sounds in American English](#). *Pragmatics & Cognition*, 14:129–182.
- Barbara Wheatley. 1996. [CALLHOME Mandarin Chinese Transcripts](#). *FTP FILE*. Philadelphia, USA: Linguistic Data Consortium.
- Ramin Yaghoubzadeh, Karola Pitsch, and Stefan Kopp. 2015. [Adaptive grounding and dialogue management for autonomous conversational assistants for elderly users](#). In *Proceedings of the 15th International Conference on Intelligent Virtual Agents*, pages 28–38, Delft, The Netherlands.
- Victor H. Yngve. 1970. On getting a word in edgewise. In Mary Ann Campbell et al., editors, *Papers from the Sixth Regional Meeting of the Chicago Linguistic Society*, pages 567–577. Chicago Linguistic Society, Chicago, IL, USA.

## A Conversational Feedback Lexical Statistics

### Cue Word Lists

In Table 3, we present the list of cue words used for computing the lexical statistics in Section 4. **Content warning:** the lists contain potentially offensive language.

Language	FB	Lexical cues
de	+	ja, jaa, jaha, jap, jep, jo, joa, aha, hey, ach, achso, okay, ok, richtig, sicher, verstehe, cool, wow, klar, gut, definitiv, absolut, genau, natürlich, ja ja, jaja, ja okay, okay ja, ja genau, ja klar, ja gut, gut okay, ah ja, ja richtig, aber sicher, aber klar, na klar, ich weiß, weiß ich, das stimmt, du hast recht, sie haben recht, ja genau richtig, vermutlich, ja vermutlich, aber wirklich
	-	nein, nee, nö, niemals, stimmt nicht, das glaube ich nicht, glaube nicht, das glaub ich nicht, glaub nicht, vermutlich nicht
	?	wirklich, bitte, entschuldige, häh, was, wo, warum welchen, welcher, welche, welches, echt, bist du sicher, sind sie sicher
	=	mhm, m, mm, hm, ähm, mh, oh, äh
en	+	yes, yeah, yep, okay, oh, right, alright, good, ok, sure, ah, nice, cool, exactly, absolutely, true, great, oh wow, right right, oh okay, oh yeah, yeah right, um-hum yeah, that's great, yes yes, yeah yeah, uh-huh yeah, that's right, right yeah, oh yes, i see, i know, that right, that's true, that's good, all right, of course, got it, is he, oh that's nice, oh that's good, well that's nice, oh i see, oh that's great, yeah that's true, well that's good, well that's great, right that's right, oh yeah yeah, that sounds good, yeah that's right, yeah yeah yeah, yeah oh yeah, oh yeah oh, well that's true, i guess so, yeah i agree, yeah it is, i think so, oh i know, yeah i know, it really is, it is, i agree, definitely, i do too, you bet, you're right, it does, i think so too, that's it, i think you're right, i know it, i agree with you, it was, i agree with that, they are, deal, indeed, obviously, clearly, precisely, certainly, no doubt, so do I, I guess so, they really are, it did, they were, they did, me too, to me too, for me too
	-	no, wait, gosh, nope, my goodness, oh no, but um, but uh, stop it, oh my goodness, oh my gosh, wait a minute, oh my god, not really, not much, no way, shit, fuck, oh no
	?	what, really, oh really, why not, you sure, is that right
	=	um-hum, uh-huh, huh-uh, uh, hum, hm, hey, well, wow, um, huh, mh, mmhmm, m, um-hum um-hum, oh uh-huh, uh-huh uh-huh, um-hum um-hum um-hum, oh, ooh, hmm, mm, mmm
fr	+	oui, ouais, ok, ah, voilà, bien, d'accord, super, parfait, exactement, ah ouais, ouais ouais, et ouais, d'accord, ah oui, oui oui, c'est ça, eh ouais, ah ouais, je sais, très bien, je comprends, bien sûr, ouais ouais ouais, ah ouais ouais, c'est vrai, ah ouais d'accord, ah d'accord, ah ouais OK, ah ouais ok, ah oui oui, ah ben oui, tu m'étonnes, c'est bien, sans doute, tout à fait, absolument, vachement, je suis d'accord, moi aussi, c'est vrai, c'est juste, c'est exactement ça
	-	non, putain, pff, si, merde, oh putain, non non, mon dieu, oh mon dieu, je sais pas, non non non, pas trop, pas vraiment, pas possible
	?	hein, quoi, vraiment, comment ça
	=	ah, mh, euh, oh, han, ben, bon, hm, hum, peut-être, m, mh mh, mh ouais, ah bon, mh mh mh, eh, hé, hey
no	+	ja, jo, å ja, ok, oi, greit, presis, wow, riktig, sant, nettopp, absolutt, jepp, definitivt, åpenbart, deal, selvfølgelig, sikkert, akkurat, god, bra, helt sikkert, jeg vet, jeg skjønner, helt riktig, det stemmer, klart det, uten tvil, det er riktig, det er greit, det er sant, det er det, jeg er enig, du har rett, det gjør det, jeg tror det, jeg vet det, det var det, det gjør jeg, jeg antar det, det gjorde det, det gjør jeg også, det tror jeg også, jeg tror du har rett, jeg er enig med deg, jeg er enig i det, de er det, de var, det gjorde de, meg også, til meg også, for meg også
	-	nei, faen, javel, herregud, ikke helt, ikke mulig, ikke i det hele tatt
	?	virkelig, hva, hæ
	=	m, mhm, mh, hmm, mm, mmm, mmhmm, hm, uh-huh, ikke sant
hu	+	igen, tényleg, úgy van, helyes, jogos, igaz, valóban, pontosan, tudom, rendben, ok, oké, oksi, okés, okszi, igen az, de az, bizony, természetesen, határozottan, feltétlenül, mindenképp, egyetértek, szerintem is, ó igen, hogyne, tényleg az, én is, nekem is, engem is, tőlem is, bennem is, igazad van, naná, mi az hogy, meghiszem azt, biztosra veheted, biztos lehets benne, jó, ja, szerintem igen, szerintem is, én is így gondolom, én is úgy gondolom, ennyi, ez az, így van, úgy van, szerintem igazad van, szerintem igazatok van, tudom, jól tudom, egyetértek, az volt, ez volt, de, azok, igen, azok, megegyeztünk, egyértelműen, azt hiszem, kétségtelenül, biztosan, persze, értem, tudod, stimmel, valóban, hát igen, hát dehogyne
	-	nem, nem igazán, nem létezik, a francba, a fenét, ne, a csodát, hogy a csodába, hát nem
	?	ó tényleg, micsoda, tényleg, miért ne, biztos
	=	aha, hú, ú, ó, óh, hű, ja, mhm, mm, mmm, hmm, hmmm, wow, azta, ejha, nahát, ühüm

it	+	ehi, okay, okay, ok, sì, si, vabbè, ecco, perfetto, wow, esatto, certamente, esattamente, assolutamente, sicuramente, decisamente, ovviamente, precisamente, di sicuro, sono d'accordo, concordo, eccellente, grandioso, ottimo, certo, infatti, fantastico, magnifico, naturalmente, giusto, bene, già, lo ben so, ah ah, ah ha, vero, é vero, lo so, lo è, davvero, vero, oh sì, lo è veramente, anch'io, anche io, hai ragione, d'accordo, va bene, benissimo, bello, buono, penso di sì, credo di sì, mi sa di sì, mi pare di sì, anche secondo me, lo penso anch'io, è così, penso che tu abbia ragione, penso tu abbia ragione, credo che tu abbia ragione, credo tu abbia ragione, mi sa che hai ragione, sono d'accordo con te, sono d'accordo con voi, lo era, lo è stato, lo è stata, sono d'accordo con ciò, lo sono, senza dubbio, a posto, ci sto, lo sono stati, lo erano, anche a me
	-	oddio, merda, no, non proprio, non molto, non è possibile, cazzo, oh no, macché
	?	come, davvero, cosa
	=	eh, Mm-hmm, hmm, mmm, mh, eh, mhmh, eh, m, hm, ah, oh, beh, uh-huh, mmh, eeh
ja	+	そう, はい, ええ, そうか, はあ, どうぞ, 本当, は, あっ, ああ, あ, ね
	=	うん, ふーん, えっ, へえ, うーん, ふん, え, う
	-	ううん, いいえ, いや, いえ, ない, 全くない, ちょっと
	?	何
zh	+	okay, yeah, yes, ok, 对, 哦, 好, 是, 有, 真的, 还行, 然,, 太好了, 耶, 行, 一定, 没错, 那好, 对了, 真好, 好啊, 好吧, 可以, 太棒了, 太棒了, 好极了, 说得对, 没问题, 我同意, 懂了, 一样, 我也是, 不错, 是啊, 就是这样, 当然可以
	-	不, 沒有, 不起, 不是
	?	啊, 是吗, 什麼, 什么, 为什么
	=	hey, oh, 嘿, 嗯, 呃, 哼, 哈, 嘘, 喔, 呵呵, 噢, 哇, 哦, 哟, 咦

Table 3: Lists of cue phrases employed in the lexical overview of Section 4. We distinguish between four core categories of feedback, namely positive feedback/acknowledgment (+), neutral/continuer (=), negative (-), and clarification request (?).

### Lexical Statistics plots

Figures 6 – 12 present statistics for utterance and feedback types as well as common feedback-related lexical items for different languages. Figure 13 shows politeness keywords and emojis in our English and French corpora.

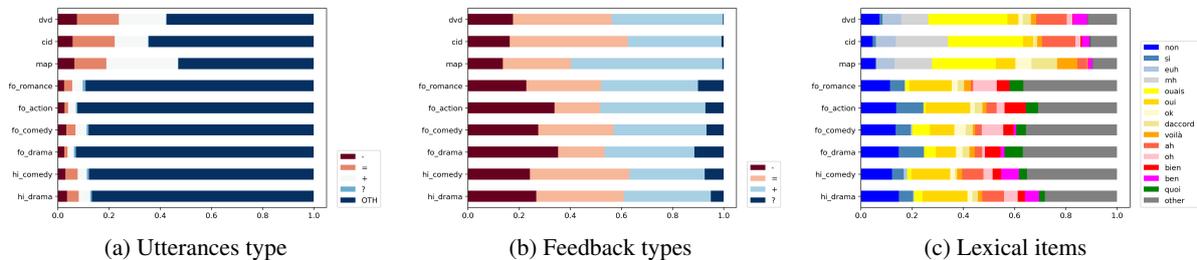


Figure 6: French across genres (rule-based, based on cue word lists).

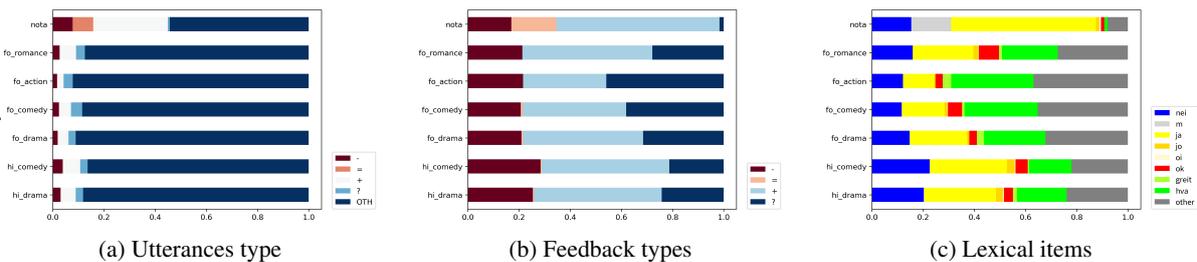


Figure 7: Norwegian across genres (rule-based, based on cue word lists).



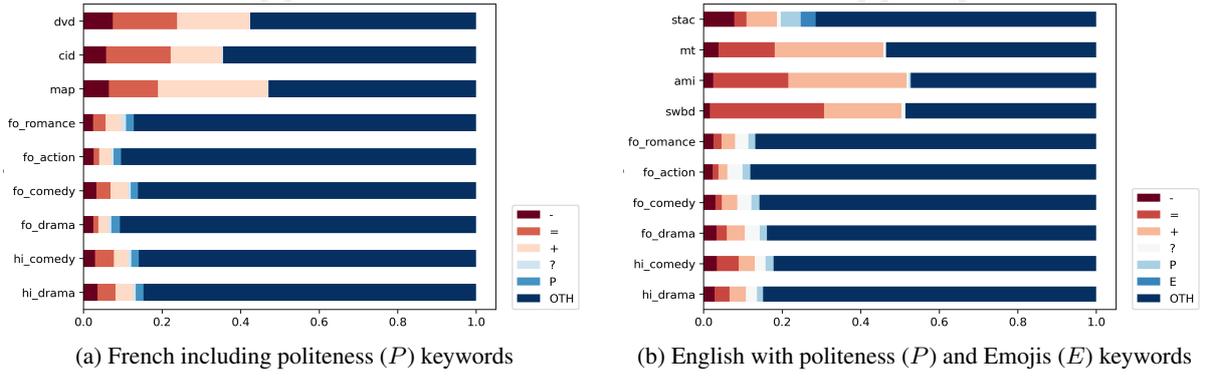


Figure 13: Short utterance distribution including politeness and emojis.

## B Detailed Dialogue Act Tagging Results

### Dialogue Act Grouping

Table 4 shows the distribution of instances per mapped dialogue act group in the DAMSL-Switchboard (SWBD) corpus.

DA group	# inst.	SWBD labels
<b>Forward looking</b>	109,382	sd, fx/sv, bf, na, ny <sup>e</sup> , arp, nd, no, cc, co, oo, ad, qr/qy, qw, qw <sup>d</sup> , qh, qo
<b>Backchannel</b>	41,017	b, bk, bh, bf, br
<b>Assessment</b>	15,727	aa, fe/ba
<b>Yes/no answer</b>	4,324	ny, nn
<b>Other</b>	40,124	<i>all other categories</i>
<b>Total</b>	<b>210,574</b>	

Table 4: Instances created from the DAMSL-SWBD corpus with labels mapped to coarse-grained dialogue act groups.

### Results per Corpus

Tables 5 and 6 present the results of our dialogue act tagger per (sub)corpus used. Here, we only make a binary distinction by grouping the feedback-relevant classes *Backchannel* and *Assessment* into a single *Feedback* category. The number of utterances refers to the final version of the data after pre-processing with meta-linguistic information removed.

Lang	Corpus	# utt	# feedback	% feedback
<b>de</b>	action_foreign	12,760	1,703	13.35
	action	12,134	1,637	13.49
	comedy_foreign	12,627	1,849	14.64
	comedy	16,152	2,369	14.67
	crime_foreign	12,589	1,245	9.89
	crime	11,817	1,581	13.38
	drama_foreign	14,669	1,350	9.2
	drama	11,460	1,452	12.67
	romance_foreign	13,499	1,500	11.11
	romance	11,809	1,596	13.52
<b>en</b>	action	11,094	1,437	12.95
	action_foreign	12,908	1,448	11.22
	comedy	13,948	1,665	11.94
	comedy_foreign	13,533	1,677	12.39
	crime	14,990	1,700	11.34

	crime_foreign	13,911	1,267	9.11
	drama	14,944	1,729	11.57
	drama_foreign	10,243	1,041	10.16
	romance	16,132	2,166	13.43
	romance_foreign	15,521	1,698	10.94
<b>fr</b>	action_foreign	11,236	1,119	9.96
	action	12,406	1,453	11.71
	comedy_foreign	17,239	1,788	10.37
	comedy	13,932	1,913	13.73
	crime_foreign	12,159	1,017	8.36
	crime	10,821	1,003	9.27
	drama_foreign	10,002	804	8.04
	drama	11,094	1,313	11.84
	romance_foreign	12,043	1,360	11.29
	romance	13,959	1,604	11.49
<b>hu</b>	action_foreign	12,781	1,377	10.77
	comedy_foreign	15,031	1,998	13.29
	comedy	14,692	2,462	16.76
	crime_foreign	13,620	1,655	12.15
	drama_foreign	13,138	1,400	10.66
	drama	7,872	1,103	14.01
	romance_foreign	13,771	1,611	11.7
<b>it</b>	action_foreign	12,010	1,585	13.2
	action	7,703	826	10.72
	comedy_foreign	15,055	2,058	13.67
	comedy	15,363	1,777	11.57
	crime_foreign	12,454	1,320	10.6
	crime	13,885	1,479	10.65
	drama_foreign	17,444	2,289	13.12
	drama	12,838	1,467	11.43
	romance_foreign	14,702	1,696	11.54
	romance	14,573	1,549	10.63
<b>ja</b>	action_foreign	11,245	967	8.6
	action	3,007	443	14.73
	comedy_foreign	16,173	1,777	10.99
	comedy	15,675	2,555	16.3
	crime_foreign	16,296	1,311	8.04
	drama_foreign	14,201	997	7.02
	drama	11,410	1,204	10.55
	romance_foreign	14,042	1,210	8.62
	romance	1,780	145	8.15
<b>no</b>	action_foreign	10,480	892	8.51
	action	1,855	290	15.63
	comedy_foreign	14,406	1,834	12.73
	comedy	11,957	1,199	10.03
	crime_foreign	12,788	1,137	8.89
	crime	9,863	853	8.65
	drama_foreign	12,031	1,202	9.99
	drama	6,688	589	8.81
	romance_foreign	12,830	1,313	10.23
	romance	4,197	399	9.51
<b>zh</b>	action_foreign	11,570	967	8.36
	action	2,722	159	5.84
	comedy_foreign	14,692	1,564	10.65
	comedy	13,587	1,034	7.61
	crime_foreign	10,778	795	7.38
	crime	11,182	697	6.23
	drama_foreign	14,527	1,330	9.16
	drama	9,567	743	7.77
	romance_foreign	13,362	1,079	8.08
	romance	11,440	700	6.12

Table 6: Number and frequency of communicative feedback phenomena predicted by the BERT-based dialogue act tagger on our subtitle corpora. Non-English datasets were automatically translated into English before inference.

<b>Lang</b>	<b>Corpus</b>	<b># utt</b>	<b># feedback</b>	<b>% feedback</b>
<b>de</b>	Hamburg MapTask	4,012	1,126	28.07
<b>en</b>	AMI	83,085	20,044	24.12
	Fisher	2,117,748	421,069	19.88
	HCRC MapTask	26,949	8,366	31.04
	STAC	5,841	514	8.8
<b>fr</b>	CID	12,326	1,754	14.23
	Aix-DVD	7,578	1,323	17.46
	French MapTask	6,046	1,226	20.28
<b>hu</b>	BUSZI-2	30,979	6,125	19.77
<b>it</b>	CLIPS	24,289	4,461	18.37
<b>ja</b>	Japanese CallHome	38,701	13,432	34.71
<b>no</b>	NoTa-Oslo	85,506	16,861	19.72
<b>zh</b>	Chinese CallHome	17,853	2,251	12.61

Table 5: Number and frequency of communicative feedback phenomena predicted by the BERT-based dialogue act tagger on spontaneous dialogue corpora. Non-English datasets were automatically translated into English with the Google Translate API before inference.

# Exploring the Use of Natural Language Descriptions of Intents for Large Language Models in Zero-shot Intent Classification

Taesuk Hong<sup>1,2</sup>, Youbin Ahn<sup>1</sup>, Dongkyu Lee<sup>1</sup>, Joongbo Shin<sup>1,†</sup>,  
Seungpil Won<sup>1</sup>, Janghoon Han<sup>1</sup>, Stanley Jungkyu Choi<sup>1</sup>, Jungyun Seo<sup>1,2</sup>

<sup>1</sup>LG AI Research, <sup>2</sup>Sogang University

lino.taesuk@gmail.com,

{youbin.ahn, movingkyu, jb.shin, seungpil.won,  
janghoon.han, stanleyjk.choi, seojy}@lgresearch.ai

## Abstract

In task-oriented dialogue systems, intent classification is crucial for accurately understanding user queries and providing appropriate services. This study explores the use of intent descriptions with large language models for unseen domain intent classification. By examining the effects of description quality, quantity, and input length management, we identify practical guidelines for optimizing performance. Our experiments using FLAN-T5 3B demonstrate that 1) high-quality descriptions for both training and testing significantly improve accuracy, 2) diversity in training descriptions doesn't greatly affect performance, and 3) off-the-shelf rankers selecting around ten intent options reduce input length without compromising performance. We emphasize that high-quality testing descriptions have a greater impact on accuracy than training descriptions. These findings provide practical guidelines for using intent descriptions with large language models to achieve effective and efficient intent classification in low-resource settings.

## 1 Introduction

In task-oriented dialogue systems, mapping user utterances to a predefined set of intents is crucial and is known as 'intent classification.' This process is essential because it helps determine the service that the user requires, making it the foundational step in fulfilling the user's goal via a chatbot (Bang et al., 2023; Sung et al., 2023; Zhang et al., 2021a, 2022). Due to the vast range of domains where chatbots can be utilized and the limited availability of intent classification data, research on transferring intent classifiers to unseen domains under low-resource conditions is very active (Zhang

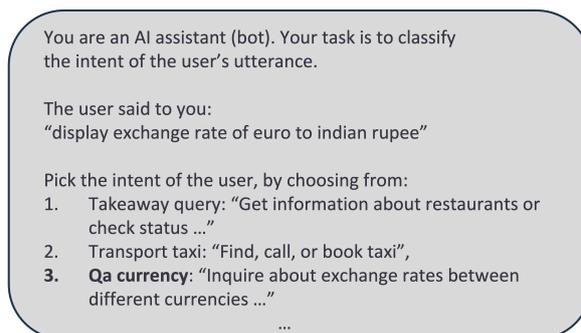


Figure 1: An example of an intent classification input for a large language model that includes intent descriptions. The figure was adopted from Parikh et al. (2023).

et al., 2021b; Mueller et al., 2022; Kuo and Chen, 2023).

Parikh et al. (2023) proposed an in-context learning classification method using large language models to classify intents in unseen domains. They provided detailed intent descriptions as inputs to compensate for the lack of user query examples for each intent. Figure 1 illustrates how descriptions are included in the in-context learning input. Provided descriptions can capture the subtle semantic nuances and exceptions that are challenging to address with intent names alone. However, the paper does not clarify the quality or quantity of descriptions that should be used for training or inference, leaving practitioners without concrete guidelines. This paper aims to provide specific guidelines on the effective and efficient use of intent descriptions during training and testing for intent classification with large language models in unseen domains.

This study specifically explores how to utilize intent descriptions in large language models through the following aspects: 1) *Effect of description quality*: Using Chat-

<sup>†</sup>Corresponding author

GPT (OpenAI, 2023), the study collects intent descriptions for the CLINC150 (Larson et al., 2019), HWU64 (Liu et al., 2021), and BANKING77 (Casanueva et al., 2020) datasets. Three sets of descriptions are collected: **dependent descriptions**, which consider semantic differences between intents, **independent descriptions**, generated without considering semantic differences, and **cleansed descriptions** manually filtered to address the subtle semantics. The impact of description quality on training and testing is investigated. 2) *Impact of description quantity*: The study examines the effect of increasing the number of descriptions used for training on intent classification accuracy. 3) *Input length Management*: Usage of off-the-shelf rankers, selecting the most probable intent options based on the similarity between the user query and descriptions, is examined to address the input length issue caused by descriptions. The optimal number of intent options to select that balances the trade-off between input length and performance is investigated. The study uses FLAN-T5 3B (Chung et al., 2022) which is an instruction-tuned model of T5 3B model (Raffel et al., 2019).

Our findings and contributions can be summarized as follows:

- Fine-tuning is required for effective understanding of descriptions in large language models and high-quality descriptions improve classification accuracy for both training and testing.
- Enhancing the quality of test descriptions has a more significant impact on accuracy than improving those used for training.
- Using a ranker to reduce to around ten classification achieves similar performance to using all options.

## 2 Method for Analysis

### 2.1 Quality-varied Description Generation and Filtering

To investigate the impact of description quality on intent classification using large language models, three different qualities of descriptions were collected using ChatGPT (gpt-3.5-turbo) via the OpenAI API. Prompts used for the

<https://openai.com/index/openai-api>

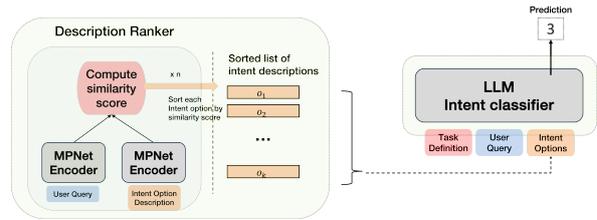


Figure 2: An off-the-shelf ranker scores the similarity between the user query and each description, selecting the top ‘ $k$ ’ intent options for intent classification input.

API calls can be found in the Appendix A.

### Independent Description Generation

The nuanced differences between distinct intents pose challenges for intent classification. For **independent descriptions**, prompts were crafted to include only a single intent and three user query examples specific to that intent, excluding other intents. Consequently, the collected description may lack comparative context, resulting in relatively lower quality. Prompts for each intent was called seven times to collect a total of seven **independent descriptions** per intent.

**Dependent Description Generation** In contrast, **dependent descriptions** include all possible intents within the prompt to ensure that the generated description uniquely distinguishes itself from others. Thus, these descriptions are considered relatively higher quality. For each intent, seven unique **dependent descriptions** was collected using API call.

**Human-Cleansed Description** Since the automatically collected descriptions may not fully capture differences between intents, manual review was added. One description per intent was carefully filtered to ensure clear distinction from other intents. This final filtering aimed to produce highest quality descriptions among our control-group for description quality. Henceforth, we will refer to this type of description as a **cleansed description**.

### 2.2 Description-Based Intent Option Ranker

Including intent descriptions increases the input length proportionate to the number of intent options. Given a model with a maximum length of 1024 tokens, descriptions of just ten

words per intent for 100 intents would exceed this limit. To address this, a description-based ranker was used to optimize input length. The off-the-shelf mpnet-base-v2 (Song et al., 2020; Reimers and Gurevych, 2019) model was employed. Figure 2 shows how this ranker integrates into the intent classification architecture. It calculates the similarity between user queries and intent descriptions, sorts intent options by similarity, and passes the top- $k$  intents to the intent classifier. This paper experimentally determines the optimal  $k$  to maintain high performance while reducing input length.

### 2.3 Fine-Tuning Large Language Models for Intent Classification

Consider the user’s utterance of  $i$ -th instance as  $u_i$  and intent options as  $\{o_{i1}, o_{i2}, \dots, o_{in}\}$ . Descriptions for each intent are denoted as  $\{d_{i1}, d_{i2}, \dots, d_{in}\}$ . All intent options are organized as

1.  $o_1: d_1,$
2.  $o_2: d_2,$
- ...
- $n. o_n: d_n.$

Replacing this option text with a predefined instruction template forms the input *INST*. The training objective for FLAN-T5 and Llama-2-Chat is defined as:

$$L(\theta) = - \sum_i^N \log p(y_i | INST; \theta), \quad (1)$$

where  $y_i$  is the correct intent index mapped to the  $i$ -th instance,  $N$  is the total number of instances, and  $\theta$  represents model parameters. An example of an input as an instruction format, *INST*, can be found in Appendix E.

## 3 Experiments

### 3.1 Datasets

We used the publicly recognized intent classification datasets CLINC150, HWU64, and BANKING77. For training, we divided ten domains of the CLINC150 dataset in half and trained on 75 intents from five domains. The remaining 75 intents from the other five domains were reserved for testing. This domain split simulates an unseen domain scenario for the intent classification test. Detailed statistics for the datasets are provided in Appendix F.

<https://sbert.net/>

Table 1: Rows lower in the table represent higher description quality used during training. Similarly, columns further to the right indicate higher description quality used during testing. The accuracy for CLINC dataset is reported.

		Types of Descriptions Used in Testing		
		without descriptions	independent descriptions	dependent descriptions
Types of Descriptions Used in Training	without descriptions	<b>84.28%</b> ±3.95%	84.15%±3.67%	90.55%±2.16%
	independent descriptions	81.93%±4.05%	85.64%±3.53%	90.97%±1.45%
	dependent descriptions	82.1%±3.56%	<b>86.99%</b> ±3.09%	<b>91.75%</b> ±1.91%

### 3.2 Impact of Description Quality on Intent Classification Training and Testing

Table 1 examines how description quality affects training and testing in intent classification models. When testing without descriptions, model trained without descriptions achieves the highest performance at 84.28%, while the performances of models trained with **independent** and **dependent descriptions** drop by 2.35% and 2.18% absolute points, respectively. However, when models trained with descriptions are tested with descriptions (specifically, **independent descriptions**), scores improve by 1.49% and 2.84% over the model trained without descriptions, respectively. This indicates that descriptions not only help models understand the detailed semantics of intents to improve classification accuracy but that fine-tuning models to understand descriptions enhances their ability to leverage them in testing.

The score improvements of the model trained with **dependent descriptions** over the model trained with **independent descriptions** demonstrate that fine-tuning with higher-quality descriptions optimizes their effective use in classification. This result supports the premise of this research that improving description quality is crucial and should not be left to random selection. In testing, higher-quality descriptions can boost performance, and their influence is more significant than in training. The model trained with **dependent descriptions** starts at 82.1% when tested without descriptions, improves by 4.89% when tested with **independent descriptions**, and achieves an additional 4.76% increase when tested with **dependent descriptions**. The improvement in testing quality has a larger impact

Table 2: The middle row shows models trained using a single type of **dependent description** per intent. In contrast, the top row represents models trained using five different **dependent descriptions** per intent, alternating during training. The bottom row shows models trained with a single manually filtered **cleansed description**.

		Types of Descriptions Used in Testing	
		1 dependent description	1 cleansed description
Types of Descriptions Used in Training	5 dependent descriptions	82.05%±3.15%	85.28%±1.69%
	1 dependent description	82.52%±2.44%	85.71%±1.61%
	1 cleansed description	83.4%±3.71%	85.79%±1.69%

than that in training. Notably, a model trained with **independent descriptions** but tested with **dependent descriptions** scored 90.97%, while one trained with **dependent descriptions** but tested with **independent descriptions** scored only 86.99%. This clearly shows that testing is particularly sensitive to description quality.

### 3.3 Impact of Description Quantity on Intent Classification Training

This experiment evaluates the effect of the quantity and quality of descriptions on model training. The result is shown in Table 2. The results reveal little to no difference between models trained with multiple descriptions and those trained with just a single description. In fact, performance tends to decline with the inclusion of varied descriptions. However, training with higher-quality descriptions – **cleansed descriptions** – resulted in the highest performance. This highlights the importance of training with a higher-quality description, even if only one, rather than relying on multiple descriptions of varying quality.

### 3.4 Optimizing ‘ $k$ ’ for Efficient Intent Classification with Ranker

This experiment investigates the optimal value of ‘ $k$ ’ for a ranker, determining the number of intent options to include in the intent classification input. Figure 3 demonstrates the performance trends as  $k$  increases. In the CLINC dataset, starting at approximately 44.21% accuracy with  $k$  set to 1, performance improves consistently as  $k$  increases, peaking at around 90% when  $k$  reaches around 13.

These results indicate that using a descrip-

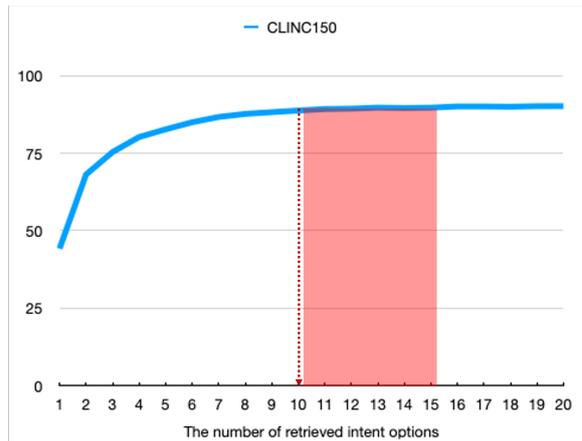


Figure 3: Graph depicts the intent classification accuracy on CLINC dataset converges as  $k$  becomes near 10.

tion ranker with the top  $k$  around 10 intent options provides near-optimal performance. The CLINC dataset, with 75 intent options and descriptions of 10 to 20 tokens each, requires around 1,200 to 1,300 tokens in total. By retrieving only the top 10 descriptions, the required input length drops to 300 to 400 tokens, reducing the input size by roughly 75%. This demonstrates that the approach proposed in this study significantly optimizes instruction-tuned models, enhancing their efficiency by minimizing the input length required for classification. For the HWU and BANKING dataset, the similar trend is shown and it can be found in Appendix B.

## 4 Conclusions

This paper thoroughly explored the impact of intent description quality and quantity on zero-shot intent classification using large language models while addressing the challenges of increased input length. The results show that fine-tuned models with descriptions are more effective for intent classification with descriptions. Additionally, higher-quality descriptions for both training and testing enhance performance, particularly during testing. Using an off-the-shelf ranker to reduce input length by selecting the top ten intent options minimizes input length without significant trade-offs in performance. Overall, this study provides practical guidelines for leveraging intent descriptions with large language models to address intent classification in low-resource settings.

## References

- Namo Bang, Jeehyun Lee, and Myoung-Wan Koo. 2023. [Task-optimized adapters for an end-to-end task-oriented dialogue system](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7355–7369, Toronto, Canada. Association for Computational Linguistics.
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. [Efficient intent detection with dual sentence encoders](#). In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45, 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](#). *Preprint*, arXiv:2210.11416.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. [Don’t stop pretraining: Adapt language models to domains and tasks](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.
- Hui-Chi Kuo and Yun-Nung Chen. 2023. [Zero-shot prompting for implicit intent prediction and recommendation with commonsense reasoning](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 249–258, Toronto, Canada. Association for Computational Linguistics.
- Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. [An evaluation dataset for intent classification and out-of-scope prediction](#). 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 1311–1316, Hong Kong, China. Association for Computational Linguistics.
- Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2021. [Benchmarking Natural Language Understanding Services for Building Conversational Agents](#), pages 165–183. Springer Singapore, Singapore.
- Tingting Ma, Qianhui Wu, Zhiwei Yu, Tiejun Zhao, and Chin-Yew Lin. 2022. [On the effectiveness of sentence encoding for intent detection meta-learning](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3806–3818, Seattle, United States. Association for Computational Linguistics.
- Aaron Mueller, Jason Krone, Salvatore Romeo, Saab Mansour, Elman Mansimov, Yi Zhang, and Dan Roth. 2022. [Label semantic aware pre-training for few-shot text classification](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8318–8334, Dublin, Ireland. Association for Computational Linguistics.
- OpenAI. 2023. Chatgpt (gpt-3.5-turbo). Accessed via <https://chat.openai.com>.
- Soham Parikh, Mitul Tiwari, Prashil Tumbade, and Quaizar Vohra. 2023. [Exploring zero and few-shot techniques for intent classification](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 744–751, Toronto, Canada. 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](#). *CoRR*, abs/1910.10683.
- 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.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. [MpNet: masked and permuted pre-training for language understanding](#). In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS ’20*, Red Hook, NY, USA. Curran Associates Inc.
- Mujeen Sung, James Gung, Elman Mansimov, Nikolaos Pappas, Raphael Shu, Salvatore Romeo, Yi Zhang, and Vittorio Castelli. 2023. [Pre-training intent-aware encoders for zero- and few-shot intent classification](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10433–10442, Singapore. Association for Computational Linguistics.
- 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](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*. Preprint, arXiv:2307.09288.

Haode Zhang, Haowen Liang, Yuwei Zhang, Li-Ming Zhan, Xiao-Ming Wu, Xiaolei Lu, and Albert Lam. 2022. *Fine-tuning pre-trained language models for few-shot intent detection: Supervised pre-training and isotropization*. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 532–542, Seattle, United States. Association for Computational Linguistics.

Haode Zhang, Yuwei Zhang, Li-Ming Zhan, Jiaxin Chen, Guangyuan Shi, Xiao-Ming Wu, and Albert Y.S. Lam. 2021a. *Effectiveness of pre-training for few-shot intent classification*. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1114–1120, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Jianguo Zhang, Trung Bui, Seunghyun Yoon, Xiang Chen, Zhiwei Liu, Congying Xia, Quan Hung Tran, Walter Chang, and Philip Yu. 2021b. *Few-shot intent detection via contrastive pre-training and fine-tuning*. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1906–1912, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

## Appendix

### A Prompt for Description Generation using ChatGPT

**Prompt for Independent Description Generation** The independent description generation prompt is illustrated as follows:

#### Independent-Description Generation Prompt

```
Intent Name: {intent name}
Few-Shot Queries: {q1}, {q2}, {q3}
Instruction:
The above is the list of intents and their examples.
Now, I want you to create unique descriptions for the intent. Make the description of the intent, 'intent name'. Here, make the description that encompasses the provided few-shot queries. Also, don't use the given use cases examples of intent for the description. Make the descriptions no longer than 10 words. I want you to return the result as following format of json:
```

```
List({
  "{intent}": "description"
})
DO NOT return any words other except for the requested format of the result.
```

**Prompt for Dependent Description Generation** Dependent description generation prompt has the following format:

#### Dependent-Description Generation Prompt Example

```
Intent Name: {intent name}
Few-Shot Queries: {q1}, {q2}, {q3}
...
Intent Name: {intent name}
Few-Shot Queries: {q1}, {q2}, {q3}
```

```
Instruction:
The above is the list of intents and their examples.
Now, I want you to create unique descriptions for each intent. This time, please make the description of the intent, '{intent}'. Here, the most important thing is that each description of intents is distinct and separate to each other. Don't make one description of intent to be inclusive to another. For example, if you have an intent, 'find restaurant', 'restaurant', don't make the description of each of them to be 'Find a available restaurant' and 'every acts related to restaurant' so that the former one is inclusive to the latter one. Also, don't use the given use cases examples of intent for the description. Make the descriptions longer than 10 words. Generate as long as possible. I want you to return the result as following format of json:
```

```
List({
  "{intent}": "description"
})
DO NOT return any words other except for the requested format of the result.
```

## B Optimizing ‘k’ for Efficient Intent Classification with Intent Option Ranker in HWU and BANKING datasets

For HWU, initial accuracy is 26.51 points with  $k$  set to 1, rising to almost maximum of 80 points around  $k$  equals 15. Lastly, in the BANKING dataset, the accuracy begins at 50.83 points and reaches around 80 points near the optimal  $k$  value of 10.

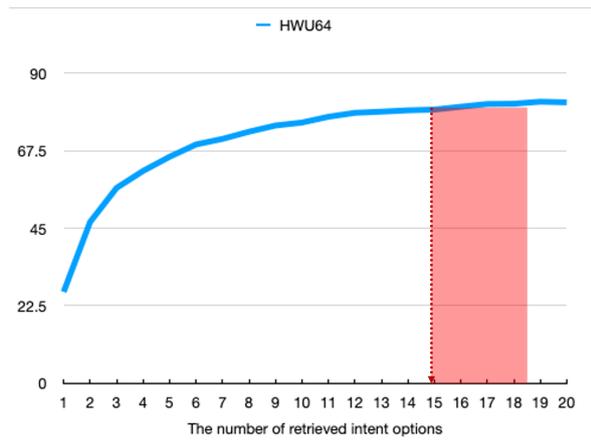


Figure 4: The graph shows that intent classification accuracy on the HWU dataset converges as  $k$  approaches near 15.

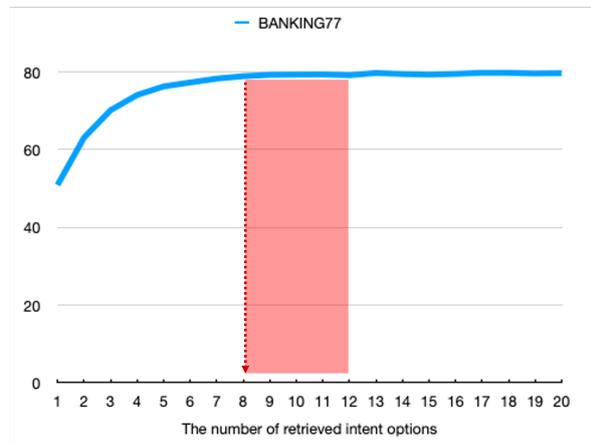


Figure 5: The Graph depicts the intent classification accuracy on BANKING dataset converges as  $k$  approaches near 10.

## C Baseline Experiments

In Table 3, we compare our model, labeled as FLAN-T5-ranker (ours), with state-of-the-art models presented by Sung et al. (2023), which performed intent classification on the CLINC,

Table 3: Out-domain intent classification accuracy compared to state-of-the-art models and baselines. Zero and few-shot accuracy results are reported as percentages. The datasets had 50, 25, and 27 intent options for CLINC, HWU, and BANKING datasets respectively. We followed the same configuration but trained on CLINC data only outside the 50 test set intents.

	CLINC $N=50$		HWU $N=25$		BANKING $N=27$	
	K=0	K=1	K=0	K=1	K=0	K=1
L-BERT <sub>TAPT</sub> (Gururanganet et al., 2020)	79.5	86.5	63.1	69.4	70.1	78.5
L-SBERT <sub>Paraphrase</sub> (Maet et al., 2022)	84.5	90.9	67.5	75.5	77.4	82.8
L-PIE (Sung et al., 2023)	86.5	91.8	70.6	77.4	77.6	82.9
FLAN-T5 (ours)	<b>97.53</b>	<b>97.62</b>	<b>87.92</b>	<b>87.22</b>	<b>84.72</b>	<b>85.52</b>
FLAN-T5-ranker (ours)	96.46	96.26	86.23	85.92	84.88	85.34
Llama-2-Alpaca (ours)	96.38	96.91	85.92	86.07	83.61	84.10
Llama-2-Alpaca-ranker (ours)	96.15	96.24	85.71	86.01	83.95	84.44

HWU, and BANKING datasets. Please refer to the original paper for details on the baselines: L-BERT<sub>TAPT</sub>, L-SBERT<sub>Paraphrase</sub>, and L-PIE.

Using the FLAN-T5 3B model fine-tuned with dependent descriptions and tested with the top-10 ranked cleansed descriptions per option, our zero-shot approach outperformed L-PIE by 9.96, 15.63, and 7.28 points for the CLINC, HWU, and BANKING datasets, respectively. When trained on one sample per intent (one-shot learning), our model showed improvements of 4.44, 8.61, and 1.54 points over L-PIE for those datasets. The significant gap between our model and the state-of-the-art may be attributed to size differences, but these results demonstrate the objectivity of our findings and the model’s superior performance over existing models.

Our model without the ranker, labeled FLAN-T5 (ours), shows slightly better performance than the version using a ranker, but the difference is minimal.

We also trained another well-known instruction-tuned model, Meta’s Llama-2-Chat 7B (Touvron et al., 2023). This model was initially instruction-tuned with the Stanford Alpaca dataset (Taori et al., 2023) and further fine-tuned using intent classification data. Our model, referred to as Llama-2-Alpaca-ranker (ours), achieved accuracy comparable to our state-of-the-art FLAN-T5 model. Notably, our proposed method of using a ranker did not negatively impact performance and even provided slight improvements on the BANKING dataset. This confirms that using a ranker can not only reduce the burden of handling long inputs but also maintain effective performance in zero-shot intent classification.

## D Training Detail

We use the *HuggingFace* implementation for fine-tuning FLAN-T5 models. In training FLAN-T5 model, AdamW optimizer with the learning rate  $2e - 5$  is used in training. The learning rate is gradually decayed during training with a cosine scheduler. The model is trained for 2 epochs and the batch size is 64. Every FLAN-T5 model performance reported in this work is the model of the final epoch. We run experiments with 4 NVIDIA A100 GPUs.

## E Instruction Input Example

The Figure 6 shows an example of an input to the model for fine-tuning intent classification task. We manually crafted ten instruction templates following the FLAN v2 format (Chung et al., 2022) for the intent classification task. The input consists of a section that instructs the model to classify the given intent, a section with the user query, and another with the intent options.

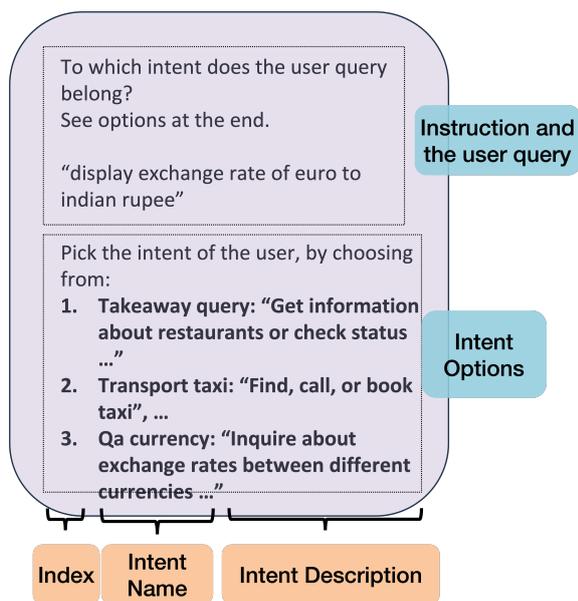


Figure 6: An input example of an instruction format.

## F Dataset Statistics

The Table 4 provides statistics for the training and testing datasets of CLINC, HWU, and BANKING. For the CLINC and HWU datasets, the domains were split in half for different seeds,

while for the BANKING dataset, all intents were split in half. The numbers below the dashed line represent the number of instances for each seed. The 'Seen domain' column corresponds to the training data, and the 'Unseen domain' column corresponds to the testing instances.

Table 4: The statistics for the training and testing datasets of CLINC, HWU, and BANKING.

seed	Seen domain	Unseen domain		
	CLINC	CLINC	HWU	BANKING
42	credit cards, banking, auto and commute, meta, utility	home, travel, work, kitchen and dining, small talk	music, recommendation, news, email, general, iot, transport, qa, date-time	banking
	7,500	2,250	669	1,560
52	auto and commute, banking, work, utility, kitchen and dining	home, meta, travel, credit cards, small talk	music, cooking, iot, play, transport, qa, date-time, social, weather	banking
	7,500	2,250	524	1,560
62	meta, kitchen and dining, credit cards, utility, work	home, travel, auto and commute, banking, small talk	alarm, music, audio, recommendation, general, play, lists, qa, cooking	banking
	7,500	2,250	621	1,560

# Voice and choice: Investigating the role of prosodic variation in request compliance and perceived politeness using conversational TTS

Éva Székely<sup>1</sup>, Jeff Higginbotham<sup>2</sup>, Francesco Possemato<sup>2,3</sup>

<sup>1</sup>Division of Speech, Music and Hearing, KTH Royal Institute of Technology, Sweden

<sup>2</sup>Department of Communicative Disorders and Sciences, University at Buffalo, NY, USA

<sup>3</sup>Centre for Language and Cognition, Rijksuniversiteit Groningen, The Netherlands

szekely@kth.se, cdsjeff@buffalo.edu, f.possemato@rug.nl

## Abstract

As conversational Text-to-Speech (TTS) technologies become increasingly realistic and expressive, understanding the impact of prosodic variation on speech perception and social dynamics is crucial for enhancing conversational systems. This study explores the influence of prosodic features on listener responses to indirect requests using a specifically designed conversational TTS engine capable of controlling prosody, and generating speech across three different speaker profiles: female, male, and gender-ambiguous. We conducted two experiments to analyse how naturalistic variations in speech rate and vocal effort impact the likelihood of request compliance and perceived politeness. In the first experiment, we examined how prosodic modifications affect the perception of politeness in permission- and action requests. In the second experiment participants compared pairs of spoken requests, each rendered with different prosodic features, and chose which they were more likely to grant. Results indicate that both faster speech rate and higher vocal effort increased the willingness to comply, though the extent of this influence varied by speaker gender. Higher vocal effort in action requests increases the chance of being granted more than in permission requests. Politeness has a demonstrated positive impact on the likelihood of requests being granted, this effect is stronger for the male voice compared to female and gender-ambiguous voices.

## 1 Introduction

The importance of pragmatics in the development of conversational technologies is becoming increasingly critical (Levinson, 2024). As Text-to-Speech (TTS) systems achieve greater realism in speech generation, a significant gap persists in understanding the pragmatic effects these technologies have within interactions. By modeling prosodic features based on empirical research, conversational TTS can be made more engaging and effective in

a variety of interactive contexts. The necessity for human-oriented pragmatics in these systems is particularly evident in scenarios requiring compliance to requests. Moreover, understanding the subtleties of how politeness is conveyed in language and speech is crucial for comprehending the factors influencing request compliance, especially in conversational systems that utilise TTS.

People use a variety of spoken strategies to manage the potential threats that communication can pose to their own and others' self-esteem and autonomy (Brown and Levinson, 1987). Indirect requests are employed to mitigate face-threatening acts, demonstrate respect for the listener's autonomy, and maintain the fundamentally cooperative and prosocial nature of human communicative behaviour and respect between interlocutors (Rossi et al., 2023). They are a fundamental aspect of polite discourse, reflecting the speaker's sensitivity to social dynamics and the listener's ability to interpret and respond to nuanced communicative cues (Drew and Couper-Kuhlen, 2014).

Initiating actions, such as requests, can be seen as a basic form of social coercion (Enfield et al., 2019). Requests have a bearing on the sequential organisation of the ensuing talk, while also restricting the agency of the requestee, and even threatening their autonomy (Soubki and Rambow, 2024). The linguistic structure of requests influences politeness and compliance. For example, the choice between using an imperative form, which might seem direct and blunt, and opting for a more conditional or interrogative form can alter the level of imposition perceived by the interlocutor. Chalfoun et al. (2024) emphasise the strategic use of politeness markers like 'please' in everyday requests, demonstrating how these markers are employed to manage face-threats in ill-fitted interactional contexts, particularly when requests could be seen as intrusive or when they encounter resistance from the requestee. Research on modal constructions

in requests reveals further nuances in how these requests are framed and understood in different contexts (Steensig and Heinemann, 2014). Modal verbs like ‘could’ or ‘might’ introduce a level of uncertainty or optionality into the request, thereby softening it and enhancing its politeness.

Enfield (2014) highlights the importance of the “infrastructure” that underpins requests, which includes the social and interactional contexts influencing how requests are made and received. Understanding the interplay between prosodic features and sociolinguistic norms is essential for designing effective conversational agents.

In our study, we develop and employ a prosody-controllable gender-ambiguous TTS system as a research tool to conduct controlled experiments assessing the role of prosodic variation in request compliance. This approach allows us to isolate the impact of vocal traits from gender biases without relying on the ability of voice actors to consistently reproduce prosodic variations.

The key contributions of this research are the following: We pioneer the use of a gender-ambiguous neural TTS built on spontaneous speech in perceptual studies, which allows for an unprecedented exploration of how gender perception influences listener responses to prosodic variations. Moreover, this study provides empirical evidence on how natural variations in speech rate and vocal energy influence the listener’s perception of politeness, and their responsiveness to indirect requests. Our findings illustrate that the impact of prosodic variations can differ based on the speaker’s gender profile, contributing to a more tailored approach in the design of TTS systems to accommodate diverse user interactions.

## 2 Background

### 2.1 Prosody in social signaling

Prosody contributes significantly to signaling speaker attitudes and interpersonal stances (Ward, 2019). Various aspects of stance can be predicted from prosodic features with significant accuracy beyond mere chance (Ward et al., 2017). Politeness strategies and their impact on compliance are not only influenced by the linguistic content but are also significantly modulated by prosodic features such as intonation, pitch, and speech rate. Research indicates that variations in these prosodic features can critically affect listeners’ perceptions and their subsequent responses to requests (Kendrick and

Drew, 2014). Trott et al. (2023) explore how prosodic features help in disambiguating English indirect requests, highlighting the complex interplay between acoustic signals and intended meanings in speech. They find that prosodic cues such as duration, pitch, and pitch slope significantly correlate with a speaker’s intent, influencing how listeners interpret pragmatically ambiguous utterances.

Vergis and Pell (2020) explore the effects of linguistic structure, imposition, and prosody on the perception of politeness in requests. Their findings show that prosody significantly affects politeness ratings, with prosodic features e.g. appropriate intonation and pitch enhancing perceived politeness. The study highlights that not only the content but the manner of speech delivery plays a critical role in social interactions. Similarly, Caballero et al. (2018) examine the acoustic cues of politeness, demonstrating that prosodic variations such as changes in pitch, intonation, and speech rate are essential for conveying politeness. Specifically, they found that higher pitch, increased pitch range, and a melodic intonation contour are perceived as more polite, whereas rude request displayed slower speech rate, lower pitch and tended to fall in pitch. Their analysis of verbal requests shows that while a specific prosody of politeness may not exist, these features significantly influence how politeness is perceived, with certain prosodic patterns leading to higher politeness ratings. Gryllia et al. (2018) investigated the role of pragmatics and politeness in prosodic variability in Greek wh-questions. Their study showed that context and social factors, such as the power and solidarity between interlocutors, influence prosodic patterns. The findings suggest that prosodic modifications are not merely stylistic but are pragmatically motivated to achieve desired social outcomes, such as politeness or authority.

### 2.2 Spontaneous TTS as a research tool

Voice talents can effectively use prosodic cues to convey subtle pragmatic nuances across various speech acts (Hellbernd and Sammler, 2016). At the same time, the reliance on actors to generate experimental stimuli introduces variability, as personal interpretations of how specific utterances should be delivered may differ. While analysing speech patterns in corpora of ecologically valid, spontaneous speech data avoids this bias, this approach often lacks the necessary control over linguistic content and prosodic realisations needed for conducting rigorously controlled experiments.

An emerging alternative methodology is to use state-of-the-art TTS built on spontaneous speech data to create experimental stimuli. This method combines advantages from both traditional approaches, relying on the authenticity of natural speech for modelling, and providing the controllability required for experimental rigor. Several previous works employed prosody-controllable neural TTS as a research tool in controlled listening experiments, with the aim to discover new knowledge about various aspects of speech perception. Székely et al. (2017) investigated the interaction of vocal effort and hesitation disfluencies in synthesised speech, focusing on how these factors influence the perception of uncertainty. Székely et al. (2019) discovered using spontaneous TTS that filled pauses improved the perception of speaker authenticity and engagement. Elmers et al. (2023) looked into the perceptual impact of tongue clicks using neural TTS, revealing that their inclusion can alter perceived speaker confidence. O’Mahony et al. (2024) extends this methodology with a corpus-based approach to investigate the prosody and pragmatic functions of the discourse marker "well".

As the capabilities and controllability of spontaneous TTS systems continue to evolve, this methodology is gaining increased attention for its potential to uncover new insights into how subtle signals in speech are interpreted by listeners. These detailed insights increase our understanding about speech perception in general, and they are particularly applicable in dialogue systems and Augmentative Communication Technologies (ACT), since these applications directly employ TTS.

In the current study, we use spontaneous conversational TTS as a research tool, and we advance this methodology by training a prosody-controllable multi-speaker TTS system that can generate male, female, and perceptually gender-ambiguous (Sutton, 2020) TTS. Using this we investigate the impact of prosodic features on listener’s willingness to comply with an indirect request.

### 3 Overview of the method

In this study, we develop and utilise a multi-speaker TTS model built upon two corpora of spontaneous speech, to investigate the impact of prosodic variations on listener responses to indirect requests.

**TTS model development:** A spontaneous conversational TTS system is engineered to include gender-ambiguous voice capabilities, using a mod-

ified Tacotron 2 architecture (Shen et al., 2018; Székely et al., 2023b) that allows for the dynamic control of prosodic features at the utterance level. This setup enables exploring how different prosodic renderings affect listener perceptions.

**Stimuli design and synthesis:** Stimuli for the experiments were designed to directly address our research questions regarding the social dynamics of request-making in conversation. Using an interactive interface, we synthesised these stimuli, ensuring each varied systematically in prosody according to predefined settings. This approach allowed precise control over the acoustic and prosodic variables of interest.

**Verification of stimuli:** Before deployment in experiments, all stimuli are tested for naturalness, gender ambiguity, and the presence of significant acoustic-prosodic differences across conditions, using objective measures. These verifications support the reliability and validity of the stimuli used in the subsequent online listening tests.

**Experimental setup:** The experiments are conducted as online listening tests on a crowd-sourcing platform, where participants are presented with synthesised speech samples. This method facilitates the collection of data on how listeners perceive and react to variations in speech delivery within an imagined conversational context.

## 4 Text-to-Speech Synthesis

### 4.1 Corpora

Two corpora of spontaneous conversational speech were used to build the TTS model. The first is a multimodal multi-party dataset called AptSpeech, described in Kontogiorgos et al. (2018). This dataset comprises 15 multi-party interactions involving a single moderator, a male speaker of General American English, and two distinct participants per session, engaged in a collaborative task. The speech data from the moderator was used to create a TTS corpus, along with additional recordings of reading newspaper articles and the Arctic sentences (Kominek and Black, 2004). The complete corpus has a duration of approximately 8 hours: 2h 26min of reading and 5h 40min of spontaneous speech. The second corpus was created from 14h 43min of conversational podcast recordings of a female speaker of General American English, who consented to make the recordings available for TTS research purposes. The speaker supplemented the material with 1h 52 min of reading non-fiction.

Both corpora were segmented into breath groups (stretches of speech delineated by two breath events) using the method proposed by Székely et al. (2020). Automatic Speech Recognition (ASR) was used to transcribe the utterances. The transcriptions were annotated for spontaneous speech events, such as filled pauses, breathing and repetitions, as well as turn-internal pauses and turn endings, in order to be able to produce these behaviours at synthesis. To balance the corpora, the number of breath groups from each style per speaker was set to the minimum of the two speakers, 480 breath groups of read speech and 2788 breath groups of spontaneous speech.

## 4.2 Prosody-controllable gender-ambiguous conversational TTS

We developed a prosody-controllable multi-speaker TTS model specifically for the purpose of the experiment. Our method follows (Székely et al., 2023a) closely, with the main difference being the use of spontaneous conversational speech corpora. We use a modified Tacotron 2 (Shen et al., 2018) TTS architecture which allows for features appended to the encoder output which can be controlled on a gradient at synthesis time (Székely et al., 2023a). The publicly available pre-trained gender-ambiguous model trained on 20 hours of speech data<sup>1</sup> is used as base, from which the existing speaker embedding is dropped and training is reinitialised as a warm start with the new spontaneous corpora. The prosodic features speech rate and energy were added to the encoder, using values normalised across the corpora. In this architecture, the appended features can be controlled at synthesis time on utterance- or phrase level. Speech rate is calculated as syllables per second, including silences, which is different from articulation rate, which excludes silences. As a result, slower speech rate values at synthesis time result in insertion of longer pauses as well as a slowed down articulation rate. Energy is calculated as the Root Mean Square of power of the signal, using a window of 20ms and step size of 5ms (Suni et al., 2017). Because energy is an acoustic feature that correlates with other prosodic features in spontaneous conversational speech, increasing the energy feature of the TTS results in a natural increase of pitch and articulatory effort, as well as more pronounced emphasis patterns. This results in a prosodic rendering that perceptually translates to

<sup>1</sup><https://github.com/evaszekely/ambiguous>

a quality which can be described as “speaking up”. This combination of prosodic features is sometimes referred to as upgraded, salient, or marked prosody (Selting, 1996) in the fields of Conversation Analysis (Sidnell and Stivers, 2013) and Interactional Linguistics (Couper-Kuhlen and Selting, 2018). In order to be descriptive without implying a direct emotional connotation, we will be referring to the prosodic realisation that this method of feature control in spontaneous neural TTS produces through the modification of the energy input feature as a level of *vocal effort*.

The model was trained for 45k iterations on 4GPUs (batch size 28). The speech signal is decoded from the output using the neural vocoder HiFi-GAN (Kong et al., 2020). The model published by the authors is fine-tuned on the combined corpora for 180k iterations on 4GPUs (batch size 28).

## 5 Experiments

### 5.1 Hypotheses

The present study aims to systematically investigate the interactions between prosodic features, perceived politeness, request compliance and gender by positing several hypotheses:

**Prosodic variation hypothesis (H1):** Changes in prosodic features speech rate and vocal effort influence the perceived politeness and compliance rates of requests. Faster speech rate and higher vocal effort are predicted to enhance perceived politeness and likelihood of compliance.

**Gender perception hypothesis (H2):** The gender perception of a TTS voice (male, female, or gender-ambiguous) mediates the effect of prosodic variations on politeness and compliance.

**Request type hypothesis (H3):** Listener responses are influenced by the interaction between prosodic features and the type of request: permission versus action requests.

### 5.2 Stimuli

We designed 8 indirect requests that are formulated to be considerate of the listener’s capacity to grant them, possibly at a minor inconvenience, yet also allowing room for a polite refusal. All sentences are similar in length and they all contain politeness markers that express a positive face and attempt to mitigate the controlling the threat to autonomy expressed by the utterance (e.g.: *would you mind, can I please*). Whilst slight differences

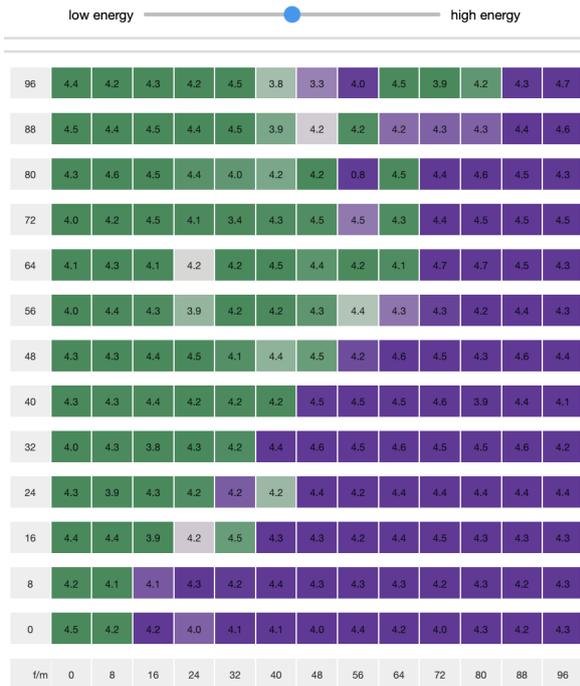


Figure 1: An example of the So-to-Speak interface, adapted to display gender-ambiguity ratings on a color gradient (green = female, purple = male) and automatic evaluation of naturalness MOS on a scale of 1-5. On the axes x and y, the percentage of respectively the male and female speaker embedding input is displayed. Samples play upon clicking on a cell. Moving the bar on top updates the grid to display a new set of samples corresponding to the setting. Step size and feature range are adjustable. In this example, speech rate is set to fast, speaker embedding and energy input vary.

in the verbal formulations might impact the willingness of the listener to comply, the decision to vary the types of politeness formulas is driven by our goal to avoid experimental monotony and to realistically reflect the variation in everyday speech. We included two types of requests, 4 sentences in each: *Permission Requests*: requests that seek authorisation or consent to perform an action. *Action Requests*: requests that involve asking for a particular service to be provided. The sentences included in the experiment are listed in Table 1.

Stimuli were created using an adaptation of the So-to-Speak interface (Székely et al., 2023b), which is an open source exploratory platform designed to help researchers interact with multi-dimensionally controllable TTS systems<sup>2</sup>. The interface enables the synthesis and playback of hundreds of samples simultaneously, displayed on an interactive grid, varying both low level prosodic

<sup>2</sup>[https://github.com/evaszekely/So\\_To\\_Speak](https://github.com/evaszekely/So_To_Speak)

Table 1: Sentences synthesised for the experiment

Permission Requests	
	<i>Is it ok if I switch off these lights?</i>
	<i>Would you mind if I opened the windows?</i>
	<i>Do you mind if I adjust the thermostat?</i>
	<i>Can I please take down these posters?</i>
Action Requests	
	<i>Would you please turn down the music?</i>
	<i>Would you mind changing the channel?</i>
	<i>Could you please turn off the air conditioning?</i>
	<i>Do you mind closing the curtains?</i>

features and high level style controls. Automatic estimates of naturalness Mean Opinion Scores (MOS) (Huang et al., 2022) are presented for each sample. For this work, we created an adaptation of the interface where the output of an automatic gender classifier (Rizhinashvili et al., 2022) is displayed on a color gradient. Figure 1 shows an example of the interface.

For the speech rate feature, two settings were chosen: *fast* and *normal*, using normalised speech rate values from the corpora, where *fast* is defined as 2 std higher than the mean. To create stimuli displaying different levels of vocal effort, three settings, *low*, *medium* and *high* were synthesised, using normalised values of the energy feature, as described in Section 4. Note that because of the corpus-driven approach to prosody control, these input features impact other characteristics in the speech samples as well, such as the increased presence of reduced articulation in fast speech rate, and higher f0 and more pronounced emphasis patterns in high energy settings. For the three gender types, three different speaker embeddings were used as input to the multi-speaker model: *male*, *female*, and *gender-ambiguous*, amounting to a total of 144 speech samples.

To verify the stimuli, we used three objective measures: Gender-ambiguity was automatically evaluated using a gender-classifier first described by Rizhinashvili et al. (2022) and adapted by Székely et al. (2023a) trained on the LibriTTS dataset (Zen et al., 2019). Additionally, all samples were evaluated through an automatic MOS rating introduced by Huang et al. (2022) which has been shown to correlate highly with perceptual ratings of naturalness. This test ensured that all stimuli had a minimum of 4.5 MOS rating. A third test was

carried out to ensure that speech rate and energy features are indeed significantly different across conditions. This test was deemed necessary because of the way the prosody-control in the TTS architecture is designed, the features are not explicitly modified, rather they are an input to the TTS model. For this reason, we validate with acoustic measurements on the samples, that the output of the TTS reflected the change in input values. The value ranges measured on speech rate and energy when these are varied in the inputs for the different settings are significantly different. This was further confirmed by a series of one-sided paired t-tests over stimuli between each combination of settings (all  $p < 0.01$ ). The measurements of speech rate, energy and  $f_0$  are found in Appendix A. Note that these values are specific to the individual speakers' own register, as represented in the training data. As such, they should not be used as independent references. The audio samples are available online<sup>3</sup>.

### 5.3 Experiment 1: politeness ratings

In order to not prime participants into specific behavior patterns regarding politeness and compliance, we conducted separate experiments concerning these two aspects. The first experiment was specifically designed to investigate how variations in prosody influence perceptions of politeness. This experiment sought to isolate politeness as variable, assessing its effects through a structured rating system. Stimuli were presented one per trial, listeners were asked to rate *“How polite does this request sound to you?”* on a scale of 1-5 (where 1 = very impolite, 2 = impolite, 3 = neutral, 4 = polite, 5 = very polite). To avoid any bias that might arise from participants recognising the experiment's focus on potential gender differences, each gender type (male, female, and gender-ambiguous) was rated by a separate group.

### 5.4 Experiment 2: request compliance

The goal of this experiment was to determine how different levels of speech rate and vocal effort (increased pitch and energy) contribute to conveying a tone of voice which makes it more likely to result in an indirect request being granted. The situational context presented to the participants was the following: *“Imagine that each request causes you a minor inconvenience—for example, if asked to open a window, consider that you are already feeling a*

<sup>3</sup><https://www.speech.kth.se/tts-demos/sigdial2024-request>

*bit cold and would prefer it closed. However, depending on how the request is conveyed, you might be more inclined to accommodate the speaker if it seems particularly important to them.”* Participants were presented with pairs of stimuli where each pair consisted of the same request rendered with different prosodic features. The task required participants to listen to each pair and decide based on the tone, which version they were more likely to grant: (a) = **A** much more likely, (b) = **A** more likely, (c) = both equally likely, (d) = **B** more likely, (e) = **B** much more likely). The pairwise design was chosen for 2 main reasons: firstly, to mitigate the effect of the differences of lexical content and formulations and topics among the individual sentences.

To gain further insight into what aspects of the speech samples people considered important while listening, at the end of the experiment, participants were asked the question: *“Could you tell us what helped you make your decisions?”*. The same between-subjects design was used as in Experiment 1: a different group of participants was recruited for each gender type.

## 6 Results

### 6.1 Experiment 1: politeness ratings

This listening test was completed by 90 people, 30 in each gender condition. The experiment took on average 10 minutes to complete and participants were paid £12 per hour. Everyone was asked to confirm that they were using headphones or earphones while listening to the stimuli. Participants' age ranged between 23 and 69, 45 identified as female and 45 as male.

A linear regression analysis was performed to evaluate the influence of speech rate, vocal effort, type of request, and gender on the average politeness rating of the various speech stimuli. For this analysis, the request type was coded 0 for permission requests and 1 for action requests. Gender also received an ordinal coding as this reflects the way

Table 2: Regression analysis of factors affecting politeness ratings

Variable	Coefficient	Std. Error	t-value	P-value
Constant	3.2705	0.050	65.729	<0.001
VoiceGender	0.1212	0.035	3.444	0.001
RequestType	-0.3950	0.057	-6.875	<0.001
SpeechRate	0.3184	0.057	5.542	<0.001
VocalEffort	0.4390	0.035	12.479	<0.001

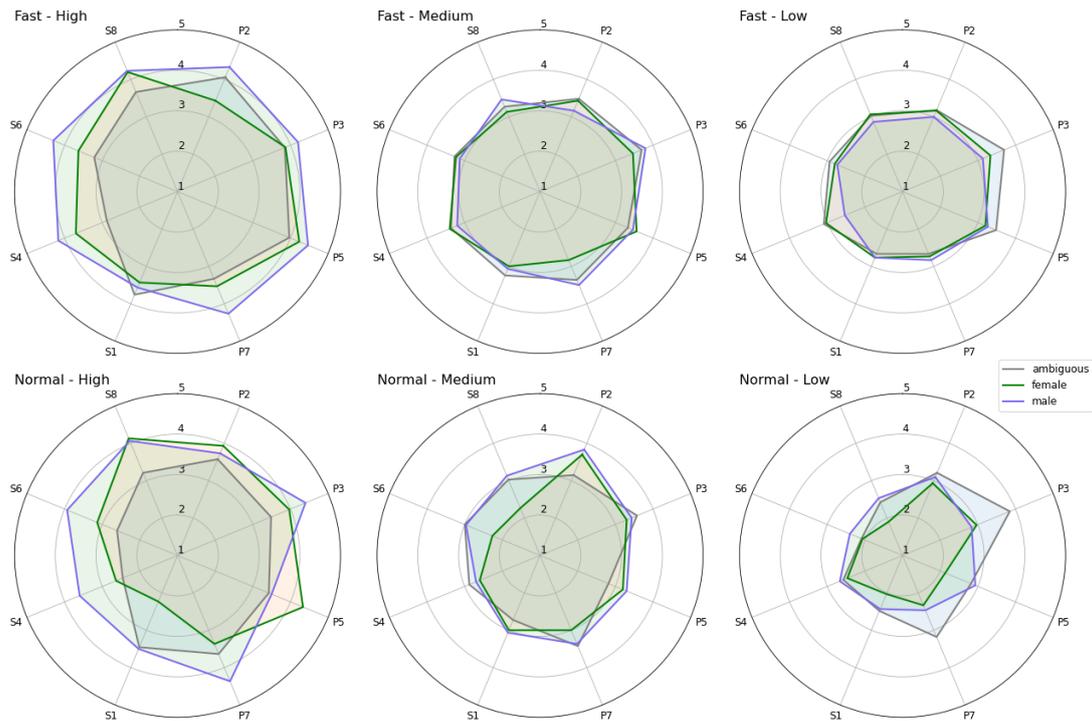


Figure 2: Results of the politeness rating per stimulus. Stimulus names on the left starting with **P**, indicate *permission requests*, and stimulus names on the right, starting with **S** indicate *action requests*. Ratings range from the center of the circle (1 = very impolite) to the edge of the circle (5 = very polite).

the TTS voices were generated: -1 for female, 0 for ambiguous and 1 for male. The results, included in Table 2 show that all four variables are significant in explaining the variance in average politeness rating. The model, based on R-squared explains 63.9% of the variance in average politeness ratings between the stimuli.

Results of this rating experiment are illustrated in Figure 2. Speech rate, vocal effort and gender have an increasing impact on perceived politeness. Action requests are rated significantly less polite than permission requests. The result that lower pitch and slower speech rate is considered less polite confirms the findings of Caballero et al. (2018).

## 6.2 Experiment 2: request compliance

90 native speakers of English, recruited through the Prolific<sup>4</sup> platform completed the study. Participants' age ranged between 23 and 75, 42 identified as female and 48 as male. Recruitment of participants followed the same setup as in Experiment 1. The experiment took on average 18 minutes to complete. Results of the test are in Figure 3 as the proportion of the pairs in which the stimulus for that condition was preferred over another,

excluding no-preference cases. Confidence intervals are calculated based on the standard error of the proportion of preferences and then applying the normal distribution's critical value to get 95% confidence intervals. We evaluated the effect of differences in speech rate and vocal effort between the two samples for each voice on the individual preference results, also controlling for participant gender. Linear regression models were applied separately for each voice. Results in Table 3 show that increases in both speech rate and vocal effort had a significant and consistently positive impact on the preference rating. The gender of the participants did not significantly influence preference for any of the voices.

Answers to the follow-up question about what helped listeners make their decision revealed that participants' decisions were influenced by their perceptions of politeness in the spoken requests. Several participants indicated that a friendlier or gentler tone made them more inclined to grant the requests, whereas harsh or demanding tones tended to deter compliance. This feedback highlights that, alongside the prosodic features conveying request importance to the speaker, the perceived politeness - or lack thereof - conveyed through prosody is a secondary factor in participants' decision-making.

<sup>4</sup>[app.prolific.com](http://app.prolific.com)

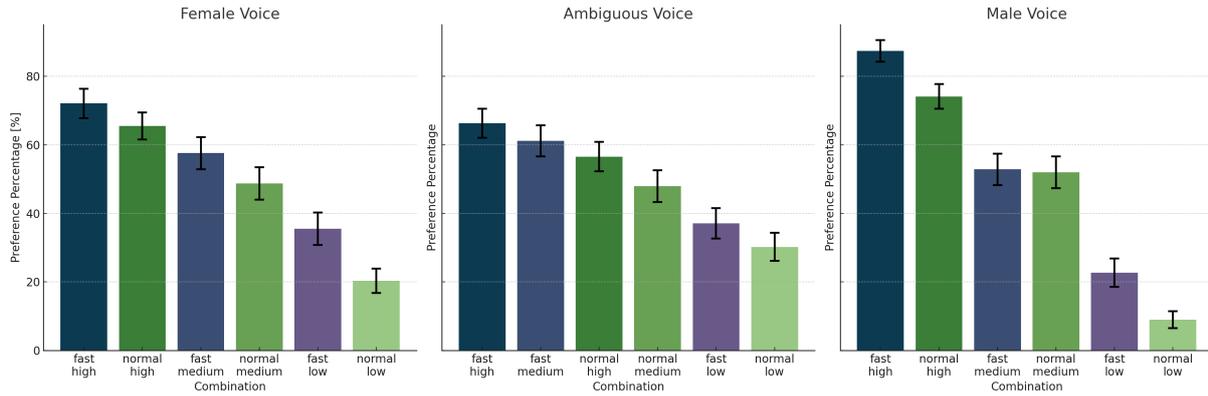


Figure 3: Percentage of evaluations where a stimulus with given speech rate: *fast* (blues), *normal* (greens)) and vocal effort: *high*, *medium*, *low* (darker to lighter shades) was preferred over other combinations, excluding ties.

Table 3: Regression Analysis Results for Compliance

Variable	Coeff.	P-Value	95% CI
<b>Female</b>			
Intercept	-0.0325	0.522	[-0.132, 0.067]
$\Delta$ Speech Rate	0.2596	<0.001	[0.143, 0.376]
$\Delta$ Vocal Effort	0.3991	<0.001	[0.354, 0.444]
Participant	0.0569	0.322	[-0.056, 0.170]
<b>Ambiguous</b>			
Intercept	-0.1206	0.046	[-0.239, -0.002]
$\Delta$ Speech Rate	0.3815	<0.001	[0.251, 0.512]
$\Delta$ Vocal Effort	0.3216	<0.001	[0.270, 0.373]
Participant	-0.0451	0.473	[-0.168, 0.078]
<b>Male</b>			
Intercept	-0.0076	0.867	[-0.097, 0.082]
$\Delta$ Speech Rate	0.2480	<0.001	[0.148, 0.348]
$\Delta$ Vocal Effort	0.6213	<0.001	[0.583, 0.660]
Participant	-0.0028	0.954	[-0.098, 0.092]

### 6.3 Influence of perceived politeness on request compliance

Combining findings from both experiments we can examine the impact of perceived politeness differences between paired stimuli on request compliance. The average perceived politeness for each sample from Experiment 1 is introduced as an explanatory variable in the analysis of the results of Experiment 2.

The significant results of the ordinal logistic regression explaining compliance, considering main

Table 4: Significant Effects in the Combined Analysis

Variable	Coeff.	P-value	95% CI
$\Delta$ Speech Rate	0.2764	<0.001	[0.135, 0.418]
$\Delta$ Vocal Effort	0.3224	<0.001	[0.205, 0.440]
$\Delta$ Politeness	0.8803	<0.001	[0.671, 1.090]
$\Delta$ Proj.* Request	0.1738	0.008	[0.045, 0.303]
$\Delta$ Pol.* Request	-0.3904	<0.001	[-0.605, -0.176]

and interaction effects are presented in Table 4. Controlling for the difference in perceived politeness between the stimuli, increasing speech rate or vocal effort still have a significant positive impact on the likelihood that a request is granted. On their own, the request type and gender of the voice do not show a significant effect. For action requests, increases in vocal effort are more effective, while the effect of politeness is more limited.

## 7 Discussion

One of the limitations in our study is that our experiments utilised only one voice per gender. To enhance the generalisability of our findings, future work will explore the use of voice conversion technologies to create a wider variety of stimuli across different gender profiles. Additionally, while this study primarily focused on prosodic features such as pitch, energy and speech rate, there are numerous other features in request articulation that warrant exploration. These include voice quality, placement of emphasis, the strategic use of pauses, and utterance-final intonation which can influence the perception of requests. Moreover, as [Levinson \(2024\)](#) points out, utterances are unlikely to be action-determinate by virtue of their form alone. Experiment 2 addresses this to an extent by presenting listeners with an imaginary scenario, but it is important to acknowledge the inherent limitations of controlled listening experiments in simulating the complex dynamics of wider social contexts. Consequently, the findings from this experiment should be further evaluated in more realistic, interactive scenarios where deeper contextual embeddings can be implemented.

Reflecting the rate of technological advancements, we expect to see an increasing demand for personalised, conversational TTS to represent and display the identity of individuals in real and virtual environments. One group, individuals with disabilities (e.g., cerebral palsy, autism, adult-onset disorders) who rely on computer-based Augmentative Communication Technologies, already use TTS to engage in real-time spoken conversations. The lack of pragmatically appropriate and effective TTS to accomplish various conversational tasks, including indirect requests, is a common critique of commercial ACTs. Our findings specifically show that adjustments to prosodic features such as speech rate and vocal effort significantly impact the perceived politeness of requests, and also affect compliance rates. This is particularly important for the design of TTS systems in ACTs, where effectively conveying requests in a polite manner is essential for users with communication challenges. Our hope is that incorporating these insights, developers can better equip conversational systems to meet the varied communication demands of individuals, ensuring more respectful and successful interactions across both real and virtual settings.

## 8 Conclusions

This study has demonstrated that the perception of politeness significantly enhances the likelihood of requests being granted. The effectiveness of changing politeness through prosody is stronger for the male voice compared to female and gender-ambiguous voices. Additionally, higher vocal effort in action requests significantly increases the chances of compliance, more so than in permission requests. This highlights the significant role that prosodic manipulation of TTS can play in enhancing the effectiveness of communicative acts within spoken dialogue systems to accommodate diverse user interactions more effectively.

## Acknowledgments

This work was supported by This research was supported by the Swedish Research Council project Perception of speaker stance (VR-2020-02396), the Riksbankens Jubileumsfond project CAPTivating (P20-0298), the Engelke Family Foundation and National Institute on Disability, Independent Living, and Rehabilitation Research (NIDILRR grant number 90DPCP0007).

## References

- Penelope Brown and Stephen C Levinson. 1987. *Politeness: Some universals in language usage*, volume 4. Cambridge university press.
- Jonathan A. Caballero, Nikos Vergis, Xiaoming Jiang, and Marc D. Pell. 2018. *The sound of im/politeness*. *Speech Communication*, 101:14–27.
- Andrew Chalfoun, Giovanni Rossi, and Tanya Stivers. 2024. The magic word? face-work and the functions of please in everyday requests. *Social Psychology Quarterly*, page 01902725241245141.
- Elizabeth Couper-Kuhlen and Margret Selting. 2018. *Interactional linguistics: Studying language in social interaction*. Cambridge University Press.
- Paul Drew and Elizabeth Couper-Kuhlen. 2014. Requesting – from speech act to recruitment. In Paul Drew and Elizabeth Couper-Kuhlen, editors, *Requesting in social interaction*, pages 1–34. John Benjamins Publishing Company.
- Mikey Elmers, Johannah O’Mahony, and Éva Székely. 2023. Synthesis after a couple PINTs: Investigating the role of pause-internal phonetic particles in speech synthesis and perception. In *Proc. Interspeech*, pages 4843–4847.
- N. J. Enfield. 2014. Human agency and the infrastructure for requests. In Paul Drew and Elizabeth Couper-Kuhlen, editors, *Requesting in social interaction*, pages 35–54. John Benjamins Publishing Company.
- Nicholas J Enfield, Tanya Stivers, Penelope Brown, Christina Englert, Katariina Harjunpää, Makoto Hayashi, Trine Heinemann, Gertie Hoymann, Tiina Keisanen, Mirka Rauniomaa, et al. 2019. Polar answers. *Journal of Linguistics*, 55(2):277–304.
- Stella Gryllia, Mary Baltazani, and Amalia Arvaniti. 2018. The role of pragmatics and politeness in explaining prosodic variability. In *Proc. Speech Prosody*, pages 158–162. Speech Prosody Special Interest Group.
- Nele Hellbernd and Daniela Sammler. 2016. Prosody conveys speaker’s intentions: Acoustic cues for speech act perception. *Journal of Memory and Language*, 88:70–86.
- Wen-Chin Huang, Erica Cooper, Yu Tsao, Hsin-Min Wang, Tomoki Toda, and Junichi Yamagishi. 2022. The VoiceMOS Challenge 2022. In *Proc. Interspeech*, pages 4536–4540.
- Kobin H. Kendrick and Paul Drew. 2014. The putative preference for offers over requests. In Paul Drew and Elizabeth Couper-Kuhlen, editors, *Requesting in social interaction*, pages 87–114. John Benjamins Publishing Company.

- John Kominek and Alan W Black. 2004. The CMU arctic speech databases. In *Proc. SSW*, pages 223–224.
- Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. Hifi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in Neural Information Processing Systems*, 33:17022–17033.
- Dimosthenis Kontogiorgos, Vanya Avramova, Simon Alexanderson, Patrik Jonell, Catharine Oertel, Jonas Beskow, Gabriel Skantze, and Joakim Gustafson. 2018. A multimodal corpus for mutual gaze and joint attention in multiparty situated interaction. In *Proc. LREC*, pages 119–127.
- Stephen C Levinson. 2024. The dark matter of pragmatics: Known unknowns. *Elements in Pragmatics*.
- Victoria S McKenna and Cara E Stepp. 2018. The relationship between acoustical and perceptual measures of vocal effort. *The Journal of the Acoustical Society of America*, 144(3):1643–1658.
- Johannah O’Mahony, Catherine Lai, and Éva Székely. 2024. "Well", what can you do with messy data? Exploring the prosody and pragmatic function of the discourse marker "well" with found data and speech synthesis. In *Proc. Interspeech*.
- Davit Rizhinashvili, Abdallah Hussein Sham, and Gholamreza Anbarjafari. 2022. Gender neutralisation for unbiased speech synthesising. *Electronics*, 11(10):1594.
- G. Rossi, M. Dingemans, S. Floyd, et al. 2023. [Shared cross-cultural principles underlie human prosocial behavior at the smallest scale](#). *Scientific Reports*, 13:6057.
- Margret Selting. 1996. Prosody as an activity-type distinctive cue in conversation: The case of so-called ‘astonished’. *Prosody in conversation: Interactional studies*, 12:231.
- Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. 2018. Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions. In *Proc. ICASSP*, pages 4779–4783.
- Jack Sidnell and Tanya Stivers. 2013. *The handbook of conversation analysis*. John Wiley & Sons.
- Adil Soubki and Owen Rambow. 2024. Intention and face in dialog. In *Proc. LREC-COLING*, pages 9143–9153.
- Jakob Steensig and Trine Heinemann. 2014. The social and moral work of modal constructions in granting remote requests. In Paul Drew and Elizabeth Couper-Kuhlen, editors, *Requesting in social interaction*, pages 145–170. John Benjamins Publishing Company.
- Antti Suni, Juraj Šimko, Daniel Aalto, and Martti Vainio. 2017. Hierarchical representation and estimation of prosody using continuous wavelet transform. *Computer Speech & Language*, 45:123–136.
- Selina Jeanne Sutton. 2020. Gender ambiguous, not genderless: Designing gender in voice user interfaces (VUIs) with sensitivity. In *Proc. CUI*, pages 1–8.
- Éva Székely, Joakim Gustafson, and Ilaria Torre. 2023a. Prosody-controllable gender-ambiguous speech synthesis: a tool for investigating implicit bias in speech perception. In *Proc. Interspeech*, pages 1234–1238.
- Éva Székely, Gustav Eje Henter, Jonas Beskow, and Joakim Gustafson. 2019. Spontaneous conversational speech synthesis from found data. In *Proc. Interspeech*, pages 4435–4439.
- Éva Székely, Gustav Eje Henter, Jonas Beskow, and Joakim Gustafson. 2020. Breathing and speech planning in spontaneous speech synthesis. In *Proc. ICASSP*, pages 7649–7653.
- Éva Székely, Joseph Mendelson, and Joakim Gustafson. 2017. Synthesising uncertainty: The interplay of vocal effort and hesitation disfluencies. In *Proc. Interspeech*, pages 804–808.
- Éva Székely, Siyang Wang, and Joakim Gustafson. 2023b. So-to-Speak: an exploratory platform for investigating the interplay between style and prosody in tts. In *Proc. Interspeech*, pages 2016–2017.
- Sean Trott, Stefanie Reed, Dan Kaliblotzky, Victor Ferreira, and Benjamin Bergen. 2023. The role of prosody in disambiguating English indirect requests. *Language and Speech*, 66(1):118–142.
- Nikos Vergis and Marc D Pell. 2020. Factors in the perception of speaker politeness: The effect of linguistic structure, imposition and prosody. *Journal of Politeness Research*, 16(1):45–84.
- Nigel G Ward. 2019. *Prosodic patterns in English conversation*. Cambridge University Press.
- Nigel G Ward, Jason C Carlson, Olac Fuentes, Diego Castan, Elizabeth Shriberg, and Andreas Tsiartas. 2017. Inferring stance from prosody. In *Proc. Interspeech*, pages 1447–1451.
- Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu. 2019. LibriTTS: A corpus derived from LibriSpeech for Text-to-Speech. In *Proc. Interspeech*.

## A Acoustic measurements on the experimental stimuli

In the tables the output ranges are recorded of the measured speech rate, energy and f0 levels for the different input settings for speech rate and vocal effort respectively. Energy and f0 are selected as acoustic measurements as these have been demonstrated as significant predictors of listeners' perception of vocal effort (McKenna and Stepp, 2018). ANOVA tests were performed for each voice to validate the statistical difference of the output measurements for the different levels of input setting.

Table 5: Speech Rate Range (syl/s) by Voice and Input Speech Rate Level with ANOVA p-values

voice	female	ambiguous	male
fast	4.99 - 6.77	5.30 - 7.07	5.38 - 7.29
normal	3.83 - 5.15	4.15 - 5.25	4.25 - 6.00
<i>p-value</i>	<i>0.0</i>	<i>0.0</i>	<i>0.0</i>

Table 6: f0 Range (Hz) by Voice and Input Energy Level with ANOVA p-values

voice	female	ambiguous	male
high	212 - 269	177 - 261	112 - 159
medium	162 - 207	139 - 178	95 - 112
low	123 - 153	104 - 137	82 - 90
<i>p-value</i>	<i>0.0</i>	<i>0.0</i>	<i>0.0</i>

Table 7: Energy Range (RMS power) by Voice and Input Energy Level with ANOVA p-values

voice	female	ambiguous	male
high	0.071 - 0.124	0.070 - 0.119	0.058 - 0.091
medium	0.070 - 0.127	0.067 - 0.103	0.050 - 0.080
low	0.059 - 0.094	0.057 - 0.084	0.045 - 0.070
<i>p-value</i>	<i>0.00217</i>	<i>0.00034</i>	<i>0.00001</i>

# A Dialogue Game for Eliciting Balanced Collaboration

**Isidora Jeknić**

Saarland University  
Saarbrücken, Germany  
jeknic@lst.uni-saarland.de

**David Schlangen**

University of Potsdam  
Potsdam, Germany  
david.schlangen@uni-potsdam.de

**Alexander Koller**

Saarland University  
Saarbrücken, Germany  
koller@coli.uni-saarland.de

## Abstract

Collaboration is an integral part of human dialogue. Typical task-oriented dialogue games assign asymmetric roles to the participants, which limits their ability to elicit naturalistic role-taking in collaboration and its negotiation. We present a novel and simple online setup that favors balanced collaboration: a two-player 2D object placement game in which the players must negotiate the goal state themselves. We show empirically that human players exhibit a variety of role distributions, and that balanced collaboration improves task performance. We also present an LLM-based baseline agent which demonstrates that automatic playing of our game is an interesting challenge for artificial systems.

## 1 Introduction

Language use is a highly collaborative process that involves constant negotiation and cooperation between interlocutors, with the ultimate goal of facilitating mutual understanding (Clark and Wilkes-Gibbs, 1986; Clark, 1996; Grice, 1989). An improved understanding of these negotiation processes would benefit the development of future systems for effective human-AI cooperation and go beyond today's rigid division of roles between dialogue systems and users (Dafoe et al., 2020, 2021).

Collaborative dialogue is often studied through dialogue games involving reconstruction, where one player has information about a target configuration and guides the other player towards it (Clark and Wilkes-Gibbs, 1986; Zarriß et al., 2016; Kim et al., 2019; Lachmy et al., 2022). This approach places the players in fixed roles (instruction giver/follower), which is in contrast to the fluid and implicit negotiation of these roles in naturally occurring collaborative dialogue, in turn limiting the ability of such games to elicit negotiation about collaborative roles.

In this paper, we directly address this issue by introducing a collaborative game designed to elicit

dialogues with more flexible role-taking – a 2D object placement game in which the target configuration is not predetermined, but must be negotiated by the players. The players use online chat to jointly decide how to arrange movable objects, without seeing each other's boards. The initially symmetric roles ensure a level playing field between players in terms of environment knowledge and the goal state. By describing the target state as only "an identical placement", we transfer the task of goal state selection onto the players, which, in turn, enables us to study the task-solving approach that they choose to take.

We observe that players indeed exhibit a variety of collaboration strategies in this dialogue game, further illustrated by a metric we define, the dominance score, representing the degree to which one player controls the gameplay. Only a minority of player dyads choose an asymmetric strategy in which one player always dominates; this strategy is also associated with systematically lower scores than more balanced strategies. Finally, we describe a baseline computational agent for this game. It achieves a significantly lower average score than a human player using a limited collaboration strategy, indicating that natural and effective collaboration in balanced games like ours is an interesting avenue for future research.<sup>1</sup>

## 2 Background

**Collaboration in dialogue.** In situated dialogue, common ground and shared context are paramount to avoiding misunderstandings (Clark, 1996; Brown-Schmidt and Heller, 2018). Interlocutors commonly engage in a collaborative effort, i.e., negotiation, to establish these commonalities, and, ultimately, a joint purpose (Clark, 1996), frequently through a coordinated referencing approach (Clark and Wilkes-Gibbs, 1986). In order

<sup>1</sup>Our code and data are available at: <https://github.com/coli-saar/placement-game>.

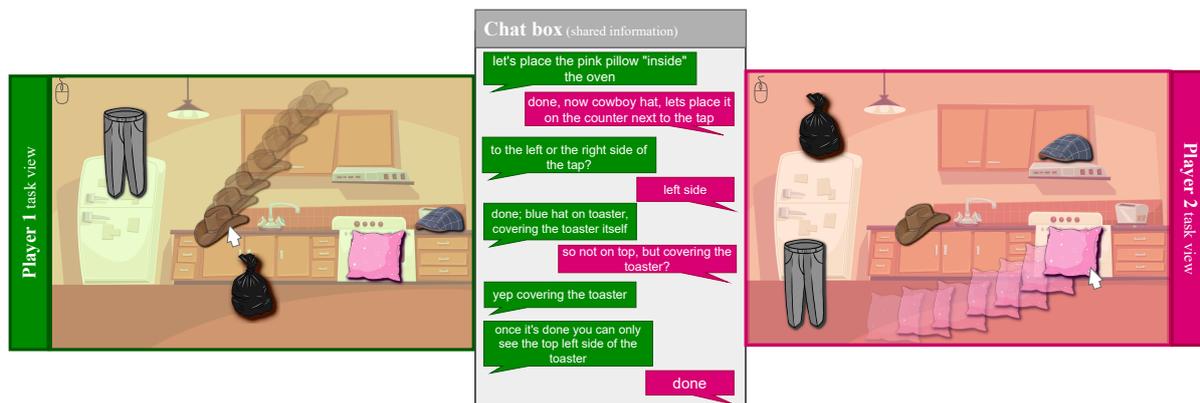


Figure 1: A reconstructed task view of both players illustrating the shared information (middle, chat box) and information only available to each respective player (left and right). Additionally, illustrates an instance of the back and forth strategy (as described in Section 4.1).

to arrive at a unified goal, the interlocutors must work together and coordinate their actions over time (Pickering and Garrod, 2013). This ongoing coordination process can lead to the acquisition of new knowledge, including how to coordinate better (Schlangen, 2023).

**Collaborative games.** There have been many reconstruction game environments developed for the purpose of studying collaboration and negotiation (e.g., Zarriß et al. (2016); Kim et al. (2019); Pacella and Marocco (2022); Narayan-Chen et al. (2019); for overview see Suglia et al. (2024)). However, all previously cited environments assume a predetermined target state to which one player must guide the other. This inherently places the participants on different levels dependent on the role they are assigned (instructor vs. follower), determined by the information they are given. We go beyond this by removing these constraints and allowing for more balanced task-solving approaches, which are necessary for a holistic study of collaboration (Schlangen et al., 2018b).

**Human-computer collaboration.** Here we refer to all collaborative situations in which “agents may be able to achieve joint gains or avoid joint losses” (Dafoe et al., 2020, 8). In the field of human-computer dialogue systems, the most frequent such agents are instruction-giving (Koller et al., 2010; Köhn et al., 2020; Sadler et al., 2024; Janarthanam and Lemon, 2010; Narayan-Chen et al., 2019; Zarriß et al., 2016), or instruction-following (Hill et al., 2020; Chan et al., 2019). While they do involve a level of first-hand human-computer interaction and dialogue necessary for completing a given task, both cases are characterized by a built-

in asymmetry, analogous to the aforementioned reconstruction games. In order to ensure successful and robust human-computer cooperation, and facilitate trust, it is integral for inherently collaborative systems (e.g., assistants) to be able to handle balanced collaboration, as well (Dafoe et al., 2020).

### 3 Collaborative object-placement game

We developed a collaborative, 2D object placement game that can be played by two players over the Internet. In each round, the two players see an identical, static background, upon which movable objects have been placed in random positions that are different for the two players (see Figure 1). The goal of the game is for the players to place each object in the same position by dragging it with the mouse. Players cannot see each other’s scene; they can only communicate through a chat window.

Each pair of players played two rounds of the game together, with a kitchen background in the first round and a living room in the second (see Appendix A.1 for more images). This allowed us to study how their collaboration strategies evolved as they became more familiar with each other.

We make the game available online by integrating it into Slurk (Schlangen et al., 2018a; Götze et al., 2022), which is a dialogue collection platform built to deal with server-side client events and API calls, ensuring participants could play the game online; additionally, it provides a straightforward and customizable logging system, as well as an off-the-shelf front-end interface with a built-in chat box feature.

All the images are “cartoonish” illustrations of real rooms and objects, in order to facilitate natural-

language communication while creating a “game” feeling. There were a total of five movable objects: a pillow, pair of pants, trash bag, flat cap, and cowboy hat. We found five items to strike a good balance between rich interactions and efficient gameplay. Our game implementation prevented placing objects on top of each other in order to enforce nontrivial reference to locations through background landmarks.

The players were scored jointly, based on the mean Manhattan distance between identical objects. The closer the two common objects were placed on the grid, i.e., the smaller the distance between them was, the higher the score the pair received. The score was normalized on a scale from 0 to 100, contributing to the typical game “feel”. Participants with very high scores (>99) got awarded a bonus.

## 4 Game playing strategies

We gathered a dataset of 71 games by crowdsourcing participants via Prolific. We used this data to analyze human dialogue behavior in a collaborative environment.

### 4.1 Collaboration strategies

The participants in our dataset exhibited a number of distinct collaboration strategies, manually detected based on the players’ contribution to the task-solving process. Examples of each strategy can be found in Appendix A.3. Crucially, what we call the “Leader” strategy – in which one player always dominates the collaboration – is a minority.

**Leader.** (33.8%<sup>2</sup>) One party predominantly leads, the other predominantly follows. It includes different situations: the explicit case (the players outwardly decide who should give instructions), cases where one player imposes the leader role and the other accepts it, or those where one player has to prompt the other for placements. The leader may or may not remain consistent across the two rounds—a swap in leadership was observed in 25% of all games, whereas 67% of games had a consistent leader. The remaining 8% were “miscommunication” cases, where both users attempt to maintain the leader role. Linguistically, we observe shorter utterances with imperative voice on the leader’s side, as typically seen in instruction-giving dialogue. Regarding the follower role, it is primarily

<sup>2</sup>The brackets contain the percentage of total games employing each strategy.

characterized by messages containing only acceptance phrases and clarification interrogatives.

**Back and forth.** (35.2%) Both parties participate actively in solving the task, and the problem solving load is split between the two players. It contains the explicit case (the parties decided to each present a new placement for alternating objects), and the more natural case (one party suggests a new placement, the other accepts and follows up with a suggestion for another object). In contrast to the Leader strategy, there is not a distinguishable leader among the player pairs who opted for the Back and forth strategy. Moreover, their messages typically contain more hedging, e.g., ending demonstrative sentences with question marks, or hedging placement suggestions with tokens such as “maybe”. Among their dialogues, we also observe more static object personification, as reflected in the placement of movable objects—for example, placing one of the hats on top of the fridge “as if it’s wearing it”. These traits all contribute to an overall much more relaxed dialogue.

**Grip Tightening.** (11.3%) The players move from a Back and forth to a Leader strategy. We observe this approach either in cases where the first round does not go as smoothly as expected (resulting in one user taking the leader role onto themselves), or when the players have established a successful task-solving approach in the first round which can be carried out sufficiently well and more efficiently by only one player in the subsequent round.

**Grip Loosening.** (19.7%) The players move from a Leader to a Back and forth strategy. The first round typically contains a user that did not fully understand the task or was reluctant to communicate, resulting in the other player having to take the initiative and lead the game. The initially reluctant user would catch on by the end of the first round, and be more willing and ready to engage in a back and forth in the second round.

### 4.2 Dominance scores

Subsequently, we calculate a *dominance score* for each player in each round of a game, capturing the extent to which one player dominates the way in which gameplay decisions are made. We assign a high dominance score to a player with high verbosity (mean message length) and high volume (percentage of messages sent, out of 100).

More specifically, let A be the player with

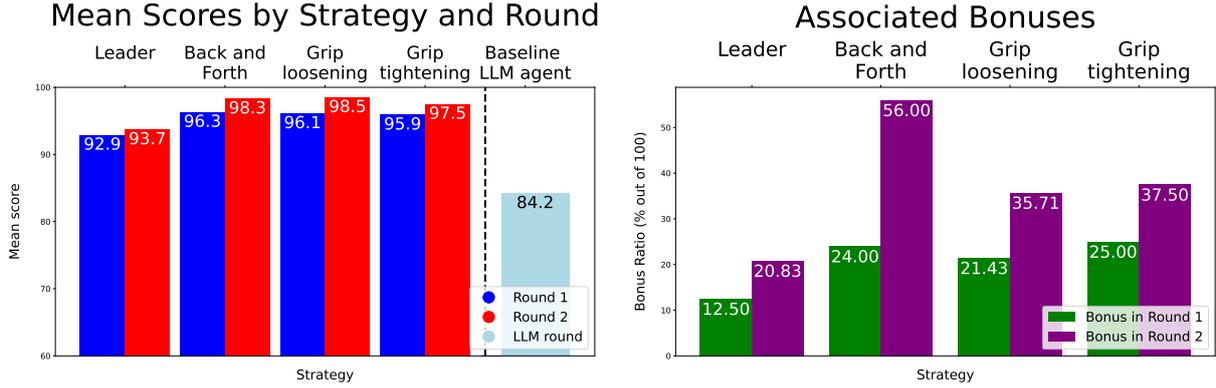


Figure 2: Overview of strategies; **left** graph shows the mean scores in each round for each strategy (out of 100), while the **right** graph shows the distribution of bonuses (score > 99) per strategy in each round (expressed in %).

Strategy	Round 1	Round 2
<b>leader</b>	1.47	2.37
<b>back and forth</b>	1.17	0.98
<b>grip tightening</b>	0.88	1.70
<b>grip loosening</b>	1.42	0.99
<b>LLM batch</b>	2.02	-

Table 1: Mean difference in the two players’ dominance scores for each round (columns) in each strategy (the first 4 rows). The last row corresponds to the mean difference of the baseline reactive LLM agent when playing with a human player described in Section 5.

the higher volume and B the other player. We let  $RD = (\text{volume}_A - \text{volume}_B) / (\text{volume}_A + \text{volume}_B)$  be the relative volume advantage of player A. Then we define

$$\begin{aligned} \mathcal{D}_A &= \text{verbosity}_A \cdot L(RD) \\ \mathcal{D}_B &= \text{verbosity}_B \cdot (1 - L(RD)), \end{aligned}$$

where  $L(x) = 1/(1 + e^{-x})$  is the logistic function, so as to dampen large differences and emphasise smaller ones, enhancing the robustness of the score.

We observe distinct patterns in each strategy’s mean dominance score difference and its development across the two rounds (see Table 1, rows 1–4), corresponding to their qualitative descriptions: in the Leader case, one player has a much higher dominance score than the other in both rounds, whereas in the Back and forth case, it is low across both rounds. In the Grip tightening case, the dominance score difference is significantly higher in the second round than the first, indicating a change from a more balanced to an asymmetric approach, while the opposite is true in the Grip loosening case.

### 4.3 Impact of strategy on task success

Figure 2 breaks down game performance by strategy. The figure on the left shows mean scores in each round for the four collaboration strategies; the figure on the right plots the proportion of games that received a bonus (score of 99 or more). It is clear that the Leader strategy underperforms with respect to the others, with Back and forth providing the greatest boost of bonus games from the first to the second round. This illustrates that our placement game is played most effectively by pairs who take a balanced approach to collaboration.

A key difference between our game and earlier reconstruction games is that our game forces the players to negotiate a goal state rather than being able to navigate to a predefined one. Moreover, the partial observability of the environment greatly impedes a leading player’s ability to monitor the other player’s actions and gauge the success of their leadership. Together, these features of our game seem to effectively encourage balanced play.

## 5 Baseline LLM agent

Our game is intended as a testbed for computational agents that collaborate effectively with humans. To gauge how challenging it is for such agents, we evaluated a simple baseline agent based on LLMs.

The agent enforces a Leader strategy, with the human player as the leader, by asking the human player for instructions in the first message and remaining passive and reactive otherwise. It uses an LLM to perform simple semantic parsing of the human’s instruction into triples of the form (object to move, landmark in the scene, spatial relation) and then uses simple handwritten rules to map such triples into  $(x, y)$  positions in the scene. For in-

stance, if the centerpoint of the fridge is at position  $(x, y)$ , the description “above the fridge” will be resolved to  $(x, y - 10)$ . We use few-shot instruction giving with GPT 3.5 Turbo Instruct (OpenAI, 2023); see Appendix A.2 for details.

In an online evaluation with ten human participants (referred to as “the LLM batch”), the agent obtained a mean score of 84.2 (left plot of Fig. 2). This shows that the task is within reach of LLM-based agents; at the same time, the agent considerably lags behind even the human-human Leader strategy, suggesting that effective collaboration remains a challenge. Moreover, we calculate the mean dominance score difference and report a score of 2.02 (see Table 1, last row). This is in line with the difference observed within the Leader strategy, further solidifying the comparability of our setup.

## 6 Conclusion and future work

We have presented a 2D object placement game which is suitable for eliciting dialogues with varied collaboration strategies. This is in contrast to earlier dialogue games, in which one player typically takes the lead. The key innovation of our game is that players must negotiate their joint goal state. A baseline computational agent achieves a task performance that is within reach of, but still considerably below human performance, indicating that variants of our game would be an interesting and challenging platform for investigating human-computer collaboration.

In the future, it would be interesting to explore even more balanced versions of the game, e.g. by adding rules that increase the cost of failed collaboration. Another avenue of future research is to investigate the interplay of collaboration strategy and mutual adaptation of the player’s lexica.

Additionally, it would be insightful to empirically verify the dominance distribution of the other aforementioned dialogue games’ outputs, as well as to further investigate the contributing factors to the occurrence of mixed-leader dialogues beyond the symmetric roles and lack of predefined goal state. Moreover, the dominance score is a useful operationalization of the collaborative imbalance between players, but it is an approximation that does not actually take the content of the players’ chat messages into account. It might be interesting to refine this measure in the future, e.g. by having the messages evaluated by an LLM. Nevertheless,

the post-hoc manual analysis of the games in the four strategies indicates that the dominance score captures differences in collaboration strategy well.

Lastly, the LLM agent presented in Section 5 is a relatively simple baseline. It is conceivable that a more intricate LLM model would close the gap to human performance, at least to the Leader strategy. We leave the exploration of such models, and of more intricate versions of our game that would remain challenging for them, for future research.

## Ethics Statement

We do not see any particular ethics challenges with the research reported here.

## Acknowledgements

We thank Sebastiano Gigliobianco for help with the backend of Slurk. We also thank the reviewers for their helpful comments. This research was partially funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID KO 2916/3-1.

## References

- Sarah Brown-Schmidt and Daphna Heller. 2018. [Perspective-taking during conversation](#). In Shirley-Ann Rueschemeyer and M. Gareth Gaskell, editors, *The Oxford Handbook of Psycholinguistics*, pages 549–572. Oxford University Press.
- Harris Chan, Yuhuai Wu, Jamie Kiros, Sanja Fidler, and Jimmy Ba. 2019. [ACTRCE: Augmenting experience via teacher’s advice for multi-goal reinforcement learning](#).
- Herbert H Clark. 1996. *Using language*. Cambridge University Press.
- Herbert H. Clark and Deanna Wilkes-Gibbs. 1986. [Referring as a collaborative process](#). *Cognition*, 22:1–39.
- Allan Dafoe, Yoram Bachrach, Gillian Hadfield, Eric Horvitz, Kate Larson, and Thore Graepel. 2021. [Cooperative AI: machines must learn to find common ground](#).
- Allan Dafoe, Edward Hughes, Yoram Bachrach, Tatum Collins, Kevin R. McKee, Joel Z. Leibo, Kate Larson, and Thore Graepel. 2020. [Open problems in Cooperative AI](#).
- Jana Götze, Maike Paetzel-Prüsmann, Wencke Liermann, Tim Diekmann, and David Schlangen. 2022. [The slurk interaction server framework: Better data](#)

- for better dialog models. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4069–4078, Marseille, France. European Language Resources Association.
- Paul Grice. 1989. *Studies in the way of words*. Harvard University Press.
- Felix Hill, Sona Mokra, Nathaniel Wong, and Tim Harley. 2020. [Human instruction-following with deep reinforcement learning via transfer-learning from text](#).
- Srinivasan Janarthanam and Oliver Lemon. 2010. [Learning to adapt to unknown users: Referring expression generation in spoken dialogue systems](#). In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 69–78, Uppsala, Sweden. Association for Computational Linguistics.
- Jin-Hwa Kim, Nikita Kitaev, Xinlei Chen, Marcus Rohrbach, Byoung-Tak Zhang, Yuandong Tian, Dhruv Batra, and Devi Parikh. 2019. [CoDraw: Collaborative drawing as a testbed for grounded goal-driven communication](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6495–6513, Florence, Italy. Association for Computational Linguistics.
- Arne Köhn, Julia Wichlacz, Álvaro Torralba, Daniel Höller, Jörg Hoffmann, and Alexander Koller. 2020. [Generating instructions at different levels of abstraction](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2802–2813, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Alexander Koller, Kristina Striegnitz, Donna Byron, Justine Cassell, Robert Dale, Johanna Moore, and Jon Oberlander. 2010. *The First Challenge on Generating Instructions in Virtual Environments*, pages 328–352. Springer Berlin Heidelberg, Berlin, Heidelberg.
- Royi Lachmy, Valentina Pyatkin, Avshalom Manevich, and Reut Tsarfaty. 2022. [Draw Me a Flower: Processing and Grounding Abstraction in Natural Language](#). *Transactions of the Association for Computational Linguistics*, 10:1341–1356.
- Anjali Narayan-Chen, Prashant Jayannavar, and Julia Hockenmaier. 2019. [Collaborative dialogue in minecraft](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, page 5405–5415. Association for Computational Linguistics.
- OpenAI. 2023. OpenAI GPT-3.5 API [turbo-instruct]. Available at: <https://platform.openai.com/docs/models/gpt-3-5-turbo>.
- Daniela Pacella and Davide Marocco. 2022. [Understanding negotiation: A text-mining and nlp approach to virtual interactions in a simulation game](#). *Applied Sciences*, 12(10).
- Martin J. Pickering and Simon Garrod. 2013. [An integrated theory of language production and comprehension](#). *Behavioral and Brain Sciences*, 36:329–347.
- Philipp Sadler, Sherzod Hakimov, and David Schlangen. 2024. [Learning communication policies for different follower behaviors in a collaborative reference game](#). *Presented at the Cooperative Multi-Agent Systems Decision-making and Learning Workshop*. This workshop is part of the Thirty-Eight AAAI Conference on Artificial Intelligence (AAAI-24).
- David Schlangen. 2023. [What A situated language-using agent must be able to do: A top-down analysis](#). *CoRR*, abs/2302.08590.
- David Schlangen, Tim Diekmann, Nikolai Ilinykh, and Sina Zarrieß. 2018a. [slurk—a lightweight interaction server for dialogue experiments and data collection](#). In *Proceedings of the 22nd Workshop on the Semantics and Pragmatics of Dialogue (AixDial/semdial 2018)*.
- David Schlangen, Nikolai Ilinykh, and Sina Zarrieß. 2018b. [MeetUp! A Task For Modelling Visual Dialogue](#). In *Short Paper Proceedings of the 22nd Workshop on the Semantics and Pragmatics of Dialogue (AixDial / semdial 2018)*.
- Alessandro Suglia, Ioannis Konstas, and Oliver Lemon. 2024. [Visually grounded language learning: a review of language games, datasets, tasks, and models](#). *Journal of Artificial Intelligence Research*, 79:173–239.
- Sina Zarrieß, Julian Hough, Casey Kennington, Ramesh Manuvinakurike, David DeVault, Raquel Fernández, and David Schlangen. 2016. [PentoRef: A corpus of spoken references in task-oriented dialogues](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 125–131, Portorož, Slovenia. European Language Resources Association (ELRA).

Step	Description
Step 1 †	verify if the message contains a set of instructions
Step 2 †	- parse the message - for each group (target, landmark, direction): 1. extract the term 2. map the term to one of the predefined allowed terms
Step 3 *	change the position of the objects according to Step 2 based on predefined constraints

Table 2: A table showcasing the logic the baseline agent followed in order to complete the task. LLM-based steps are labeled with †, whereas rule-based ones have a \*.

## A Appendix

### A.1 Game environment design.

Figure 3 depicts the two background images used for the two rounds.

### A.2 Baseline agent.

**System logic.** Each message that the human sends is analyzed by the agent following the steps from Table 2. First, the agent determines if the user’s message contained instructions, by using the input message together with the base (first prompt) and passing it to an LLM. If this step results in a TRUE, the system moves on to step 2, consisting of two extraction steps: in the first one, the agent extracts the movable (target) object and static (reference) object, and in the second one, it extracts the placement direction of the target in reference to the static object. Table 3 contains an overview of allowed terms for each extraction category. These entities are used in Step 3, which is a rule-based altering of the agent’s world state, following the rule set from Table 4 and hard-coded positions of the reference objects (this information corresponds to information available to the human, i.e., seeing one’s own board). The next section of the appendix contains the base prompts.

target	landmark	direction
pillow	fridge	on
cowboy	toaster	next to
cap	lamp	above
pants	oven	below
garbage	stove	
	counter	
	sink	

Table 3: All allowed terms per group; the extracted objects from the message are mapped to one term from each list.

	new $x[t]$	new $y[t]$
<b>on</b>	$x[r]$	$y[r]$
<b>next to</b>	$x[r] + 10$	$y[r]$
<b>above</b>	$x[r]$	$y[r] - 10$
<b>below</b>	$x[r]$	$y[r] + 10$

Table 4: The movement constraints for position manipulation. The first column contains the directions; the second and third columns refer to the target object (t)’s new x and y coordinates with respect to the reference landmark (r).

**Prompts.** Here we provide the prompts we used for the LLM part of the agent.

1. The base of the prompt used to extract the placement location, in reference to a static object.

''you are playing a game with another player in which you have to follow their instructions about where to put certain objects. i will give you a message and i want you to tell me if it contains a set of instructions. don't provide explanation, just give me the output (True or False).

examples:

[user 1]: place the lamp on the fridge  
[you]: True

[user 1]: can you put the knife in the drawer?  
[you]: True

[user 1]: do you have a toaster?  
[you]: False

[user 1]: what objects do you have?  
[you]: False

[user 1]: let's place the pan on top of the

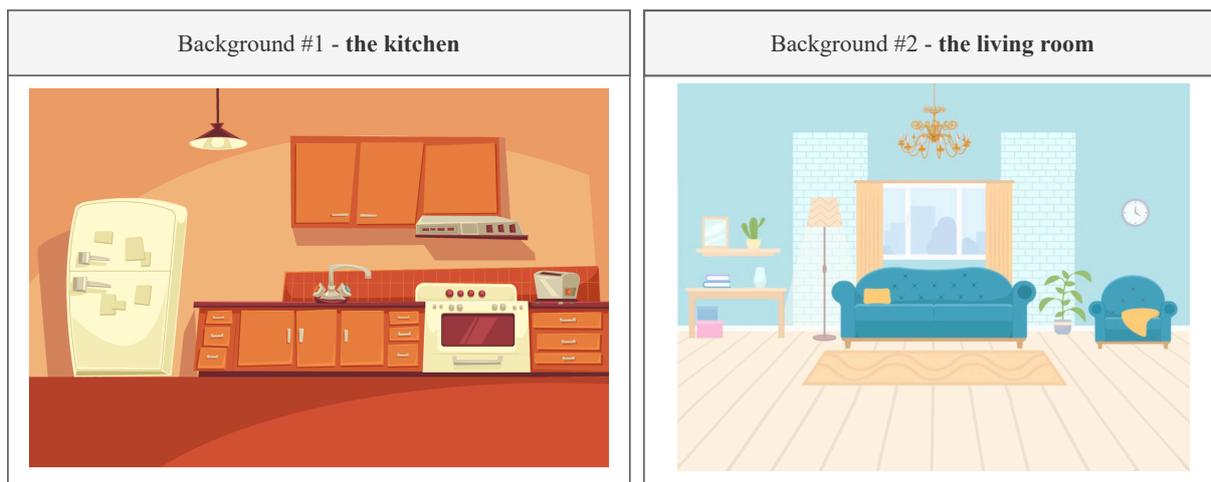


Figure 3: The background images for the two rounds.

lamp

[you]: True

[user 1]: put hat on sink

[you]: True

[user 1]: lamp on toilet

[you]: True''

**2.** The base of the prompt used to extract the static (reference) and movable (target) object.

'''i will give you a set of instructions and i want you to extract two things: one, the object that should be moved. then, i want you to compare it to the following four words and return the one it is most close to. the objects are: garbage, cowboy, cap, pants, pillow. next, i want you to extract the location where the object should be placed. then, match the output place with one of the possible places: fridge, counter, toaster, lamp, stove, oven, sink. don't provide explanation, just give me the output.

for example:

user 1: put the pillow to the right of the fridge

you: pillow, fridge

user 1: put the jeans on the stove

you: pants, stove

user 1: let's place the cushion on the ceiling light

you: pillow, lamp

user 1: place the garbagebag in the upper right corner of the counter  
you: garbage, counter

user 1: cowboy hat to the left of the water faucet  
you: cowboy, sink

user 1: the other hat on the right behind the pants  
you: cap, toaster

user 1: garbage bag on top of lamp stand  
you: garbage, lamp

user 1: let's place the blue hat on the toaster  
you: cap, toaster

user 1: put peaky blinders hat in the oven  
you: cap, oven'''

**3.** The base of the prompt used to extract the placement location, in reference to a static object.

'''i will give you a set of instructions and i want you to extract the key spatial word or phrase. then, i want you to compare it to the following four words and return the one it is most close to. the words are: above, below, next to, on. don't provide explanation, just give me the output. for example:

[user 1]: put the knife to the right of the

fridge

[you]: next to

[user 1]: put the pan above the oven

[you]: above

[user 1]: place the toilet paper in the upper right corner of the cupboard

[you]: on

[user 1]: cowboy hat to the left of the water faucet

[you]: next to

[user 1]: the cowboy hat on the right behind the pants

[you]: next to

[user 1]: pillow under the sink

[you]: below

[user 1]: garbage bag on top of lamp stand

[you]: above'''

### **A.3 Strategy examples**

Figures 4 to 7 illustrate examples of different strategies, namely:

- leader — Figure 4
- back and forth — Figure 5
- grip tightening — Figure 6
- grip loosening — Figure 7

hey  
2023-10-16T17:59:07.142205

hello  
2023-10-16T17:59:12.910876

lets place things on corners of the picture  
2023-10-16T17:59:46.262485

ok starting with what?  
2023-10-16T18:00:03.312811

pants top left, pillow top right, brown hat bottom left, grey hat bottom right, garbage bag bottom middle  
2023-10-16T18:00:31.182863

make sure to push the objects even if they get half hidden  
2023-10-16T18:01:10.661669

done  
2023-10-16T18:01:11.798818

**User 1 decides on all placements, User 2 accepts; Room 2854, Round 1**

hi  
2023-10-17T16:07:14.796980

Hi  
2023-10-17T16:07:21.811195

Where is the black bag on your screen  
2023-10-17T16:07:38.779261

its infront of the stove, on the floor  
2023-10-17T16:08:12.895898

Ok I have moved mine there  
2023-10-17T16:08:29.499296

Where is the pink cushion  
2023-10-17T16:08:34.443209

okay put the cushion on top of the fridge  
2023-10-17T16:09:05.350633

Done  
2023-10-17T16:09:18.378798

The trousers?  
2023-10-17T16:09:26.923354

put them on the sink  
2023-10-17T16:09:57.536071

Done.  
2023-10-17T16:10:11.278247

put the cowboy hat on the toaster  
2023-10-17T16:10:15.824850

I have one hat left  
2023-10-17T16:10:29.059877

put it on the stove  
2023-10-17T16:10:49.824955

Ok  
2023-10-17T16:10:54.379354

Done. We should be matching now  
2023-10-17T16:11:04.811169

**User 2 prompts User 1 for placements Room 2859, Round 1**

where are you place the items?  
2023-10-11T13:28:32.916503

hat above the toaster  
2023-10-11T13:28:56.887417

cool. trash bag on top of the fridge  
2023-10-11T13:29:42.043520

pillow above the stove  
2023-10-11T13:30:06.323539

hat just below the light  
2023-10-11T13:30:24.328903

and pant under the hat next to the fridge  
2023-10-11T13:31:02.888616

all good  
2023-10-11T13:31:31.371282

ok submit  
2023-10-11T13:31:54.225210

**User 2 asserts dominance, User 1 accepts; Room 2812, Round 1**

Figure 4: Leader strategy example; "User 1" refers to the one whose messages are pink, and "User 2" to the one whose messages are yellow.

hey
2023-10-16T16:29:12.764679
hi!
2023-10-16T16:29:18.721795
where would you like to place the objects ?
2023-10-16T16:29:39.488996
cowboy hat on top of the fridge
2023-10-16T16:30:27.180591
okay cool ! i moved it
2023-10-16T16:30:49.129561
same , i put it on the left edge of top of fridge
2023-10-16T16:31:25.918818
you can choose next item!
2023-10-16T16:31:47.549732
next , pillow ontop of the stove?
2023-10-16T16:32:08.958765
done!
2023-10-16T16:32:21.640095
pants on top of toaster
2023-10-16T16:32:36.974553
would you like to go again
2023-10-16T16:32:39.551996
sure, pants on top of toaster
2023-10-16T16:32:55.595353
perfect ! done , other hat on sink tap
2023-10-16T16:34:00.791119
done
2023-10-16T16:34:15.983679
trash bag?
2023-10-16T16:34:22.084877
where should we place it
2023-10-16T16:34:51.258232
on the light?
2023-10-16T16:35:07.891711
done
2023-10-16T16:35:37.043945
nice one!
2023-10-16T16:35:41.135892
gonna submit so
2023-10-16T16:35:48.439486
have a nice day
2023-10-16T16:35:51.594732
same to you !
2023-10-16T16:35:56.782567
The users take turns choosing the placements for objects; Room 2834, Round 1

lets place the garbage bag on top of the table
2023-10-08T15:48:10.370017
and lets place the pillow on the couch
2023-10-08T15:48:20.564849
right in the middle
2023-10-08T15:48:26.305140
which couch the small or bigger one?
2023-10-08T15:48:43.867241
the big one at the centre
2023-10-08T15:49:08.232056
okay done
2023-10-08T15:49:19.561887
blue hat on top of the lamp?
2023-10-08T15:49:30.419448
done, pants on top of the pot plant?
2023-10-08T15:50:36.482054
how about in front of?
2023-10-08T15:50:55.885612
Until the whole plant is covered?
2023-10-08T15:51:29.273680
and the cowboy hat on top of the clock/watch
2023-10-08T15:51:35.206579
yes on the pants
2023-10-08T15:51:45.324661
done
2023-10-08T15:52:07.546620
The users discuss object placement; Room 2765, Round 1

Figure 5: Back and forth strategy example; "User 1" refers to the one whose messages are pink, and "User 2" to the one whose messages are yellow.

hi
2023-10-16T19:46:30.071951
hi
2023-10-16T19:46:35.630598
let's start. let's place pillow under the lamp so it lightly touches the bulb
2023-10-16T19:47:22.433228
sure, pillow centered below lamp, slightly touching bulb
2023-10-16T19:48:13.729320
I have my garbage bag at the third big bottom doors, the single ones
2023-10-16T19:49:09.420380
single doors left of oven? Put cowboy hat on the oven gas remover
2023-10-16T19:50:10.011742
yes that's correct, putting my cowboy hat on that gray "gas remover" above oven
2023-10-16T19:50:41.714464
keep it centered
2023-10-16T19:50:58.841330
okay
2023-10-16T19:51:07.172360
lets position the trousers so it's using the line of the double doors at the bottom below the sink between legs, touching floor
2023-10-16T19:51:53.959139
done, top is reaching top of tiles
2023-10-16T19:52:20.190738
correct
2023-10-16T19:52:25.412314
last one the weird hat let's put it above the toaster so it's kinda shielding it
2023-10-16T19:53:08.404857
sure
2023-10-16T19:53:18.720050
The back and forth round; Room 2858, Round 1

here we go again
2023-10-16T19:54:07.907593
so trousers on the lamp as before touching floor leg line using lamp
2023-10-16T19:54:24.930624
sure sorry for the delay
2023-10-16T19:55:54.830008
hat shielding the lamp
2023-10-16T19:56:12.440636
the trousers' lamp
2023-10-16T19:56:19.662731
sure
2023-10-16T19:56:25.386617
which hat?
2023-10-16T19:56:34.823048
ah yes, the gray one
2023-10-16T19:56:42.525235
let's put pillow touching clock by it's center
2023-10-16T19:57:03.370789
like just above the small sofa
2023-10-16T19:57:27.940211
ok done
2023-10-16T19:57:37.317533
garbage in the flower standing on the floor
2023-10-16T19:58:05.313584
bag is touching floor
2023-10-16T19:58:23.098991
sure
2023-10-16T19:58:29.063850
and last one the cowboy let's put it on the yellow rectangular pillow
2023-10-16T19:58:54.664132
so half of the rectangle is visible
2023-10-16T19:59:14.835796
at an angle
2023-10-16T19:59:18.897982
just before it's blocked by the trousers
2023-10-16T19:59:34.888023
done?
2023-10-16T20:00:03.738153
yeah we can submit I think
2023-10-16T20:00:11.636827
submit
2023-10-16T20:00:22.566365
User 1 takes the lead; Room 2858, Round 2

Figure 6: Grip tightening strategy example; "User 1" refers to the one whose messages are pink, and "User 2" to the one whose messages are yellow.

hi	2023-10-16T16:29:36.096124
hi	2023-10-16T16:29:42.344270
what objects would you like to move	2023-10-16T16:30:04.950744
I think the black trash plastic	2023-10-16T16:30:33.794131
where would you like to move it too	2023-10-16T16:31:00.990567
Next to the fridge.	2023-10-16T16:31:28.068412
right or left	2023-10-16T16:31:42.000224
Left	2023-10-16T16:31:55.140308
lets put the pillow next to it	2023-10-16T16:32:31.647656
okay	2023-10-16T16:32:38.897917
what would you like to move next	2023-10-16T16:33:15.390177
Pants	2023-10-16T16:33:33.354289
where too	2023-10-16T16:33:41.495783
Stove	2023-10-16T16:33:49.611243
left or right	2023-10-16T16:34:01.556577
Right	2023-10-16T16:34:07.096673
awesome lets put the cowboy hat on the right side of the tap	2023-10-16T16:34:35.686751
lets put the golf hat on top of the fridge	2023-10-16T16:34:55.284423
okay	2023-10-16T16:34:56.532633
I will press submit now thank you	2023-10-16T16:35:36.348405
Great	2023-10-16T16:35:37.076501

User 1 leads the round;  
Room 2833, Round 1

Lets put the pillow on the right side of the couch	2023-10-16T16:36:30.094552
awesome done	2023-10-16T16:36:42.704305
The trash under the table.	2023-10-16T16:37:00.970195
done	2023-10-16T16:37:28.868594
lets put the cowboy hat on the lamp	2023-10-16T16:37:56.839922
Maybe the pants on the other couch	2023-10-16T16:37:57.122863
okay done	2023-10-16T16:38:08.605548
pants is on the other couch	2023-10-16T16:38:13.882111
What about the other hat?	2023-10-16T16:38:30.680392
lets put it next to the pillow	2023-10-16T16:38:41.543764
Great	2023-10-16T16:38:47.637049
awesome ill hit submit thank you	2023-10-16T16:39:01.427098

More conversational round;  
Room 2833, Round 2

Figure 7: Grip loosening strategy example; "User 1" refers to the one whose messages are green, and "User 2" to the one whose messages are pink.

# Improving Speech Recognition with Jargon Injection

Minh-Tien Nguyen<sup>1,2</sup>, Phuoc-Dat Nguyen<sup>1</sup>, Tuan-Hai Luu<sup>1</sup>, Xuan-Quang Nguyen<sup>1</sup>,  
Tung-Duong Nguyen<sup>1</sup>, Jeff Yang<sup>1</sup>

<sup>1</sup>Cinnamon AI, 10th floor, Geleximco building, 36 Hoang Cau, Dong Da, Hanoi, Vietnam.  
{ryan.nguyen, barb, sam, albert, neo, jeff.yang}@cinnamon.is

<sup>2</sup>Hung Yen University of Technology and Education, Hung Yen, Vietnam.  
tienm@utehy.edu.vn

## Abstract

This paper introduces a new method that improves the performance of Automatic speech recognition (ASR) engines, e.g., Whisper in practical cases. Different from prior methods that usually require both speech data and its transcription for decoding, our method only uses jargon as the context for decoding. To do that, the method first represents the jargon in a trie tree structure for efficient storing and traversing. The method next forces the decoding of Whisper to more focus on the jargon by adjusting the probability of generated tokens with the use of the trie tree. To further improve the performance, the method utilizes the prompting method that uses the jargon as the context. Final tokens are generated based on the combination of prompting and decoding. Experimental results on Japanese and English datasets show that the proposed method helps to improve the performance of Whisper, specially for domain-specific data. The method is simple but effective and can be deployed to any encoder-decoder ASR engines in actual cases. The code and data are also accessible.<sup>1</sup>

## 1 Introduction

Automatic speech recognition (ASR) is the task of automatically transcribing input audio to output text (Radford et al., 2023; O’Shaughnessy, 2024). The output of ASR systems can be used in several applications such as intelligent personal assistants (McGraw et al., 2016; He et al., 2019), voice searches (Chiu et al., 2018), or meeting analyses (Yu et al., 2020; Song et al., 2020; Jung et al., 2023; Li et al., 2023; Rennard et al., 2023). Recently, the performance of end-to-end ASR models has been improved by several approaches such as connectionist temporal classification (Graves et al., 2006; Graves and Jaitly, 2014), recurrent neural network transducer (Graves, 2012), attention-based encoder-decoder (Chorowski et al., 2015; Chan et al., 2016;

Dong et al., 2018) with strong ASR engines (Gulati et al., 2020; Han et al., 2020). Among those, Whisper has shown strong performance for speech recognition (Radford et al., 2023). It was trained with 680,000 hours of labeled audio data with multitasking and multilingual learning.

Strong ASR engines such as Whisper have achieved promising results in English, yet, we observe the decent accuracy of ASR engines applied to actual business, especially for low-resource languages, e.g., Japanese. To fill the gap, there are two possible solutions for domain adaptation. The first well-known solution is to continuously fine-tune ASR engines with domain-specific data (Huang et al., 2021; Javed et al., 2022). However, creating training corpora (including speech and text) data is a non-trivial task that is time-consuming and labor-expensive. In many cases, the creation requires domain experts, especially for narrow specific domains, e.g., high-pressure gas incidents. Also, fine-tuning is a complex process that requires skilled practitioners (Radford et al., 2023). The second solution is to consider domain-specific data as a context and inject the context into the decoding phase of ASR engines (Pundak et al., 2018; Zhao et al., 2019; Alon et al., 2019; Le et al., 2021b,a; Sun et al., 2021; Han et al., 2022). Among them, biasing methods are simple and potential to inject a context into the ASR process. However, these methods are usually used with hybrid ASR (Pironkov et al., 2020) or CTC end-to-end models (Graves and Jaitly, 2014) which are behind the performance of encoder-decoder ASR models such as Whisper.

This paper addresses the problem of improving the performance of ASR engines by using jargon. The problem comes from the fact that in practical cases, only jargon (domain-specific terms) is provided by clients. The jargon only includes specific words and phrases without the availability of speech data and domain-specific text. It challenges pre-trained ASR models, e.g., Whisper, and cur-

<sup>1</sup><https://shorturl.at/YiBUr>

rent methods of contextual speech recognition that usually require both speech and text data. To address the problem, we introduce a new method that injects domain-specific knowledge in the form of jargon into the decoding phase of ASR engines. To do that, the method uses Whisper as the backbone and jargon represented as a trie tree as the domain-specific context. By utilizing it as a form of the tree to manipulate the beam search decoding process and a prompt to give instructions to Whisper, the method improves the performance on various datasets and the appearance of jargon in the final output. The method does not require speech data for fine-tuning ASR models, so that facilitates the deployment in actual cases. In summary, the paper makes two main contributions as follows.

- It introduces a method that injects the jargon into the beam search decoding by boosting the score of the beam that includes tokens in the jargon. The method is further supported by the initial prompt method offered by Whisper. The method is simple, effective, and easy to adapt with any encoder-decoder ASR engines.
- It validates the efficiency of the method on Japanese and English datasets. Experimental results show that domain knowledge injection helps to improve the quality of ASR engines.

## 2 Related Work

**ASR** The recent success of deep neural networks has been contributed to improve the performance of ASR. Approaches range from traditional methods such as connectionist temporal classification (CTC) (Graves et al., 2006; Graves and Jaitly, 2014), recurrent neural network transducer (Graves, 2012), attention-based encoder-decoder (Chorowski et al., 2015; Chan et al., 2016; Dong et al., 2018), to sequence-to-sequence models (Chiu et al., 2018). These approaches leverage the development of strong ASR engines (Gulati et al., 2020; Han et al., 2020; Radford et al., 2023) trained by the Transformer architecture (Vaswani et al., 2017) such as SeamlessM4T (Barrault et al., 2023) or Whisper (Radford et al., 2023). We used Whisper as the main backbone of our method because of its efficiency for ASR in domain-specific Japanese data.<sup>2</sup>

**Context-aware ASR** has recently been used to improve the quality of ASR (Williams et al., 2018;

<sup>2</sup>Whisper gives better performance than SeamlessM4T for domain-specific Japanese datasets in the internal testing.

Pundak et al., 2018; Zhao et al., 2019; Alon et al., 2019; Le et al., 2021a,b; Han et al., 2022; Jung et al., 2022). The context can be the text of testing data (Han et al., 2022) or a list of biasing phrases (Zhao et al., 2019; Alon et al., 2019; Pundak et al., 2018). There are two main directions. The first is to bias the decoding of ASR models by using shallow fusion methods (Zhao et al., 2019; Le et al., 2021b,a). The fusion methods create a finite state transducer (FST) created from the list of biasing phrases and use the FST to adjust the decoding process without adding any neural networks. In contrast, the second usually encodes the context by using the encoder and then uses attention to change the probability of tokens in the decoding phase (Han et al., 2022). Between the two directions, TCPGen (Sun et al., 2021) introduced a tree-constrained pointer generator that incorporates a list of biasing words into both attention encoder-decoder and transducer end-to-end ASR models.

The method of contextual speech recognition is perhaps the most relevant to our work (Williams et al., 2018). The method adjusts the output likelihoods of a neural network at each step in the beam search by a sequence probability computed from  $n$ -grams. While sharing the idea of using shallow fusion, our proposed method distinguishes two main points. First, we consider a small dictionary rather than using the text of testing data to create  $n$ -grams language models (LMs) as Williams et al. (2018). It makes our task to be more challenging. Second, in the decoding phase, we modify the probability of a token appearing in the jargon while Williams et al. (2018) just simply used shallow fusion with the probability of LMs. We follow the shallow fusion approach because it is simple but effective and is appropriate for ASR in business cases.

## 3 Proposed Method

### 3.1 Problem Statement

The problem is to improve the performance of an encoder-decoder ASR engine, e.g., Whisper by taking into account jargon when doing decoding. We define the jargon as a dictionary  $D$  that includes domain-specific tokens used in actual businesses, e.g., technical terms used in the high-pressure gas incident domain. Precisely, given the jargon  $D$ , we design a method that adjusts the decoding process of Whisper to inject  $D$  for speech recognition.

Figure 1 (right) introduces the proposed decoding method. The input speech is processed by an

ASR engine. The decoder adjusts the probability of tokens by considering the jargon represented in a trie tree. The initial prompt method is also combined to further improve the quality of ASR.

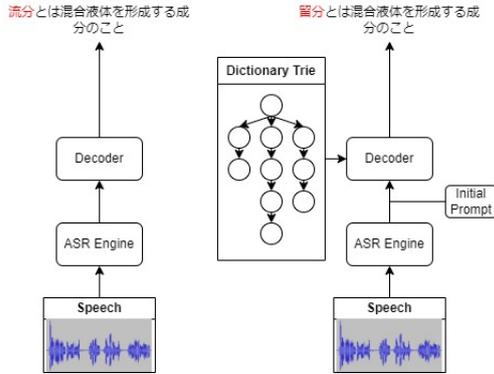


Figure 1: Whisper (left) and our proposed method (right). The inclusion of "留分" increases our model recognition. The Japanese sentence means "Part of distillation is the component that form the liquid mixture".

### 3.2 Whisper

There are several attention-based encode-decoder models have been released recently and achieved good results on various datasets, such as Seamless4MT (Barrault et al., 2023). While we use Whisper (Radford et al., 2023) as the backbone of our method, our method can be applicable to all models within this family. We select Whisper because of its strong performance of speech recognition in many languages, especially in Japanese. Whisper is a sequence-to-sequence model, based on the Transformer architecture. It consists of two transformers. As an encoder-decoder ASR engine, the first transformer encodes audio information and the second one encodes linguistic information. The model is jointly trained to predict a sequence of tokens by the decoder for many different speech-processing tasks including multilingual speech recognition, speech translation, spoken language identification, and voice activity detection. The decoder predicts each next token by conditioning on previously processed tokens with that token's probability of the model. During inference, Whisper generates a sequence of output labels given an input speech that maximizes the likelihood probability distribution.

$$y^* = \arg \max_y p(y|x) \quad (1)$$

with  $p(y|x)$  being the output probability distribution from Whisper.

### 3.3 Jargon Injection

The strong performance of pre-trained attention-based encoder-decoder ASR models, e.g., Whisper facilitates the deployment of ASR models, yet, domain adaptation makes a challenge for the deployment in actual cases. To fill this gap, a straightforward method is to fine-tune ASR models with domain-specific data. However, this method requires training data that includes both speech and its transcription. In practical cases, creating this training data is a time-consuming and non-trivial task that requires domain experts. Therefore, we come up with another direction that directly injects domain-specific jargon for domain adaptation. This section shows the proposal of injecting jargon into the decoding process of Whisper.

**Jargon representation** In business cases, jargon refers to a set of tokens that is uncommon with out-domain people, causing difficulty for an ASR model to recognize correctly or often be mistaken with similar tokens, e.g., 留分 ("part of distillation") and 流分 (does not exist in Japanese dictionary). In almost all cases, the jargon is usually created by humans in specific domains.

To represent the jargon  $D$  for decoding, we use the trie tree, which is one of most efficient methods for representing a collection of strings (Le et al., 2021a). Each node in a tree is associated with a character. The root node is associated with an empty string and children of a node will share a common prefix with its parent node. Figure 2 shows an example of a trie tree that represents "eat", "can", "could", and "card" tokens. Representing

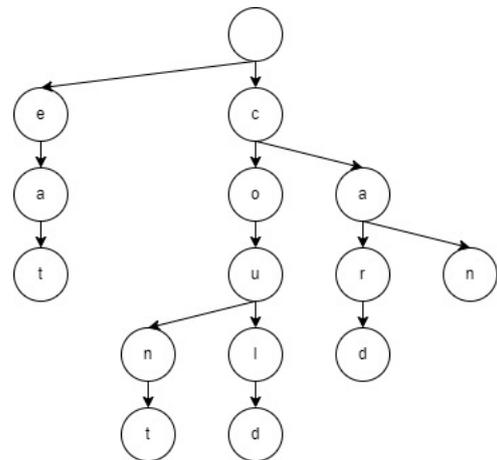


Figure 2: A trie sample made from "eat", "can", "could", "count", and "card" words.

representing  $D$  as a trie allows our method to efficiently

look up jargon during decoding. The trie stores a pointer for each sequence. Whenever the decoding method extends a sequence in the beam, it moves this pointer down to the tree. If it encounters either a leaf node or a node marked as jargon-boundary, the method permanently adds a score to the total score of that sequence using Eq. (4).

**Contextual beam search** The conventional method in Eq. (1) provides a good way for decoding the next tokens, yet, encoder-decoder ASR engines only support generating step-wise distribution  $p(y_n|y_{n-1}, \dots, y_0, x)$ . If we brute-force this formula, checking all possible sequences and computing probabilities from them, it would take a huge amount of time. That’s why, it’s more common to use beam search to approximate Eq. (1) while limiting the search space. Beam search is conducted as follows. Let  $k$  be the beam size,  $T$  is the number of distinct tokens generated from ASR engines, the beam search starts with one single sequence that consists of <eos>-equivalent tokens and its scores as  $p(\langle \text{eos} \rangle | x)$ . When decoding, the beam search extends the sequence by appending only top  $k$  tokens with the best scores among all possible  $T$  tokens that ASR models can generate, adding all  $k$  sequences back to our beam. Next, these  $k$  sequences are again appended with the next top  $k$  tokens, resulting  $k^2$  sequences in total. The beams are sorted according to their probability to retain only the best  $k$  sequences with the highest probability. The decoding process ends when the <eos>-equivalent token appears. This process is repeated until all sequences in beams are locked.

We extend the original beam search to define our contextual beam search. The idea is to adjust the probability of generated tokens by using the jargon  $D$ . In Eq. (1), the best sequence is only generated based on the output probability distribution from ASR. To take into account the additional context  $D$  for the beam search, Eq. (1) is modified as follows.

$$y^* = \arg \max_y \frac{p(y|x)}{p_D(y)^\alpha} \quad (2)$$

where  $p_D(y)$  is a distribution calculated based on the jargon  $D$  and  $\alpha$  is a tunable hyperparameter controlling the contribution of the context  $D$ . In practice, running beam search decoding by using Eq. (2) often causes numerical instability. Therefore, we transform Eq. (2) by using the logarithm version of the equation as follows.

$$y^* = \arg \max_y \log(p(y|x)) - \alpha \log(p_D(y)) \quad (3)$$

Eq. (3) shows our decoding method. It includes two main parts: one from the original ASR models and another from the jargon  $D$ . The main idea of this equation is to boost the score of sequence based on prior distribution  $p_D(y)$ . Since  $D$  is a list of specific tokens, Eq. (3) cannot calculate prior distribution by modeling it as a language model. Instead, we artificially boost the output probability whenever the process reaches the boundary of jargon. As a result, Eq. (3) is written as follows.

$$y^* = \arg \max_y \log(p(y|x)) - \alpha S_D(y, x) \quad (4)$$

with  $S_D(y, x)$  is defined as follows.

$$S_D(y, x) = \sum_{0 \leq i \leq j < n} \sum_{k=i}^j \log(p(y_k|x)), \text{ if } y_{i..j} \in D \quad (5)$$

with  $n$  being the length of sequence  $y$ . In short,  $S_D(y, x)$  is calculated by summing the output probability of every possible jargon in  $D$  containing in the sequence  $y$ . Therefore, all sequences that contain the jargon will have their boosted scores, increasing their chance of appearing in the final beam of the beam search decoder. At each step of decoding, calculating Eqs. (4) and (5) is computationally expensive because of storing and finding candidate tokens in  $D$ . Therefore, we use the trie data structure in Figure 2 for efficient string lookup.

**Initial prompt** The initial prompt is one of the prominent features of Whisper. The prompt serves as the previous context of the current speech frames. By appending a context  $C$  before decoding, Whisper is indirectly biased into the prompt. The context  $C$  is composed of two parts: an instruction and the jargon. The detailed context is shown in Section 4.2. To use  $C$ , Eq. (1) is written as follows.

$$y^* = \arg \max_y \log(p(y|x, C)) \quad (6)$$

**Combination** Since context text  $C$  is appended before decoding and does not change during the beam search, we combine the initial prompt and our method for the best result. The combination offers a global optimization that takes into account static information from the prompt (prefix) and dynamic probabilities from our decoding method (suffix). It forces Whisper to more focus on the jargon. Therefore, Eq. (4) is written as follows.

$$y^* = \arg \max_y \log(p(y|x, C)) - \alpha S_D(y, x, C) \quad (7)$$

with  $S_D(y, x, C)$  is defined similarly to Eq. (5).

## 4 Experimental Setup

### 4.1 Datasets

We validated the proposed method on three Japanese datasets: two in-house corpora and two benchmark dataset, and one English corpus.

**HGP** HGP is a smaller set of the original high-gas incident corpus published in 2022 by the High-Pressure Gas Safety Institute of Japan. From 18,171 incident cases, we extracted 1,500 incident reports in three industries: "general chemistry", "petrochemical", and "oil re- fining" (Inoue et al., 2023). HGP’s dictionary was made by automatically comparing the difference between the Whisper model output and the validation ground truth. We took only tokens that exist at least one time (or higher) in the Japanese dictionary. Due to the term,<sup>3</sup> we can not disclose the dataset.

**GPT synthesis** GPT synthesis is an in-house dataset that consists of 200 audios imitating HGP dataset’s style. Its content was both handcrafted and assisted with OpenAI’s ChatGPT. The dictionary was made manually by a domain expert, targeting technical terms in the incident domain.

**JNAS** (Itou et al., 1999) is a public Japanese dataset that consists of 16679 utterances, spoken by 306 speakers, with half of them are male (16,176 by reading Mainichi Newspaper and 503 from ATR 503 PB-Sentences). The JNAS’s dictionary was made with the same method as HGP. We compared the difference between the output of Whisper with the gold label of the validation dataset.

Table 1: Statistics of test sets of four databases. JA stands for Japanese and EN is English.

Dataset	Samples	Jargon	Domain	Lang
HGP	1467	150	Incident	JA
GPT-Syn	200	47	Incident	JA
JNAS	2253	157	Newspaper	JA
LibriSpeech	2620	146	Audiobook	EN

**LibriSpeech** (Panayotov et al., 2015) is a public English dataset that consists of approximately 1000 hours of audiobook recordings, mostly come from Project Gutenberg collection. In this paper we exclusively utilized the "clean" category of the dataset, which includes recordings from 20 male and 20 female speakers for both the "dev" and "test"

<sup>3</sup><https://shorturl.at/fnKNO>

subsets. The dictionary was automatically constructed using the same method applied to the HGP dataset, except that the comparison was performed at word-level instead of token-level.

### 4.2 Baselines

We compared the proposed method to strong baselines. The original **Whisper** (Radford et al., 2023) is the first baseline. It directly transcribes speech data to text without using our inject method. We used two versions: Whisper small (Whisper S) and Whisper medium (Whisper M) as strong baselines. We did not report the performance of the Whisper large model because of the tiny gap on testing datasets. **Initial prompt** uses the prompt as "はい、日本語で、*token1*、*token2*、..., *tokenN*の単語をすべて含むテキストを生成します。". ("Yes, it will generate a text containing all the words *token1*, *token2*, ..., *tokenN* in Japanese") for decoding. **BeamSearch + *n*-grams LM** (Williams et al., 2018) uses beam search decoding combined with a *n*-grams language model for decoding. **TCPGen** (Sun et al., 2021) incorporates a list of biasing words into both attention of encoder-decoder and transducer of ASR models. The probability of generated tokens is computed by using the probability distribution over a subset of output subword units conditioned by a prefix tre.

### 4.3 Evaluation Metrics

#### Character and word error rate (CER and WER)

We used CER and WER (Rix et al., 2001; Hu and Loizou, 2006), well-known metrics for evaluating ASR models, as the main metrics for evaluation. The CER was used for Japanese datasets and WER was used for the English dataset.

**Dictionary recognition rate** The CER or WER metrics can measure the improvement in terms of corrected predicted tokens over the gold label. However, the number of corrected characters is small compared to the total number of characters in a testing set. As a result, a small improvement in CER or WER may not represent the efficiency of the proposed method. We, therefore, define a new metric called **Dictionary Recognition Rate (DRR)** as follows.

$$DRR = \frac{\# \text{words in dictionary correctly recognized}}{\# \text{words in dictionary should be recognized}}$$

### 4.4 Implementation

The  $\alpha$  parameter was selected in the range of [0.1, 0.6] shown in Figure 3. The beam size was fixed

at five. Whisper small and medium versions were used as the main backbone for transcription. We did not report the performance of the Whisper large version due to very tiny improvements.

The method requires a single Tesla P100 GPU for running the backbone ASR model.

## 5 Results and Discussion

**Performance comparison** Table 2 shows the comparison of the proposed method to strong baselines. We can observe that the combination of our decoding with the initial prompting method obtains the best results in almost all cases on two evaluation metrics. The improvement may come from two possible reasons. First, our decoding method forces Whisper to more focus on domain-specific tokens in the dictionary. By adjusting the probability of tokens generated by the decoder of Whisper, the decoding process can replace generated tokens with those in the dictionary, e.g., 流分 (does not exist in Japanese dictionary) and 留分 (part of distillation). Second, the initial prompt method also provides a good indicator for the decoding. It uses a prompt that includes a short instruction and a set of tokens in the jargon for guiding the decoding process. This prompt forces the decoder of Whisper to more focus on tokens in the prompt when doing speech recognition. By using both methods, the combination can achieve the best results.

An important point is that the performance of initial prompting and our decoding methods is better than Whisper, the strong baseline of ASR on the four testing datasets. As mentioned, the initial prompt method directly embeds the jargon in the prompt that serves as the prefix when decoding audio segments. However, this method has a main limitation. The prompt is short (less than 250 tokens) and faces a challenge with a long dictionary in practical cases. In contrast, our decoding method directly adjusts the decoding process by using the trie tree representation. Tokens appearing in the jargon receive higher probabilities than those that are not in the jargon. As a result, our decoding method is more general and can work with any ASR engines and dictionaries. In addition, the initial prompt method does not show the efficiency on Librispeech. This is because this corpus is public and already included in the pre-training of Whisper, therefore, adding additional context in form of a dictionary is not necessary. The inclusion of Librispeech in the pre-training of Whisper

also results the high score of DRR on this dataset. Despite the limitation, in actual cases, the initial prompt method still can be used by using matching or selection methods for dealing with a large dictionary. For example, the selection method can select relevant tokens given an input speech for recognition to reduce the number of tokens in the prompt. In this case, the combination of our decoding and initial prompt methods can retain promising results. In the case that the initial prompt method is not available, our method can still output competitive scores. For example, it is the best on Librispeech and the second best of DRR on HGP. It confirms the efficiency of our proposed decoding method.

For the baselines, Whisper (small and medium versions) obtains promising results. As mentioned, by using a huge amount of data for pre-training, Whisper can work well on various domains in the multilingual setting. TCPGen does not show the efficiency on GPT synthesis and JNAS datasets because these datasets have only testing sets while TCPGen requires training data for adaptation. For Japanese, only HGP provides both training and testing sets. Therefore, we fine-tuned TCPGen on HPC and directly used the model for GPT Synthesis and JNAS testing sets. The Beam Search+ $n$ -grams LM method achieves competitive results even it is a quite simple method. It shows the role of beam search and language modeling for the contextual basing task of ASR. However, the performance of this method is still behind our methods.

Among the four datasets, the improvement on the two in-house datasets is larger than JNAS and Librispeech, two public corpora. There are two possible reasons. First, the two in-house datasets contain more domain-specific knowledge and terms that may not appear in the training data of Whisper. This, therefore, challenges the transcription of Whisper in domain-specific data. Second, JNAS and Librispeech are benchmark datasets that are easy to collect and use as the training data of Whisper. It explains the reason why Whisper operates well on benchmark ASR corpora. The proposed method is simple and effective. It does not increase much the running time due to the efficient searching method on a small dictionary using a trie tree.

**Human vs. automatic dictionary creation** Dictionaries created by humans are good at biasing, however, it's creation is costly and requires domain knowledge. We, therefore, investigated the behavior of injection methods using dictionaries created

Table 2: Performance comparison on three datasets. **Bold** text is the best and underline text is the second best. S (small) and M (medium) stand for the two versions of Whisper. This notation is also used for Table 3.

Method	HGP		GPT Synthesis		JNAS		LibriSpeech	
	CER	DRR	CER	DRR	CER	DRR	WER	DRR
Beam Search+ $n$ -grams LM	14.98	78.76	15.61	16.95	9.92	28.19	2.92	98.85
TCPGen (Whisper S)	22.90	68.59	24.61	13.60	19.98	24.77	3.23	<u>99.03</u>
TCPGen (Whisper M)	19.02	78.89	15.55	16.80	15.17	29.15	2.92	98.30
Whisper (S)	23.50	66.42	23.08	13.22	17.27	23.52	3.67	98.19
Whisper (M)	14.97	78.82	15.60	16.95	9.92	28.35	<u>2.82</u>	<b>99.25</b>
Initial prompt (Whisper S)	23.18	66.60	20.67	51.86	19.34	32.40	3.48	98.45
Initial prompt (Whisper M)	<u>12.05</u>	81.09	<u>13.09</u>	59.32	<u>9.22</u>	39.56	3.34	98.77
Ours (Whisper S)	22.82	71.86	22.40	29.49	16.96	31.15	3.65	98.21
Combined (Whisper S)	22.40	71.91	20.40	<u>79.66</u>	20.24	<u>41.31</u>	3.49	98.45
Ours (Whisper M)	14.66	<u>82.76</u>	14.86	25.76	9.60	34.74	<b>2.81</b>	<b>99.25</b>
Combined (Whisper M)	<b>11.57</b>	<b>85.08</b>	<b>12.24</b>	<b>82.37</b>	<b>9.06</b>	<b>46.26</b>	3.34	98.77

by humans or automation. As mentioned in Section 3.3, the dictionaries used in the experiments were created by humans. To assess the setting of using automatic dictionary creation, we used a method as follows. First, the speech data of testing sets of the datasets were transcribed by using Whisper. The transcription was then aligned with the gold texts to obtain wrong words based on word segmentation (we used Mecab<sup>4</sup> for Japanese). Finally, common words were removed by using Japanese or English dictionaries. Note that we could not create the dictionary by humans for Librispeech due to the large number of testing samples. Therefore, we did not report the observation on Librispeech.

Table 3 reports the scores of the injection methods using dictionaries created by humans or automation. The general trend shows that dictionaries created by humans help to improve the performance of ASR. For example, methods using dictionaries created by humans output better performance than those using dictionaries created automatically in four cases. It is easy to understand that domain experts can create closer and more correct domain-specific words that the original Whisper models can not transcribe correctly. In this sense, incorporating these words into the decoding of Whisper can improve CER. In contrast, the automatic creation method may create wrong words due to the accumulated error of ASR. It then affects the final scores of the injection methods. An interesting observation is that the gap between the two setting is small. It suggest that we can reduce human effort in creating domain-specific dictionaries

by applying appropriate automatic methods.

Table 3: The performance of the injection methods using automatic dictionary creation or human-involved dictionary creation. **IP** stands for Initial Prompt.

Method	HGP	GPT Syn	JNAS	Avg
Human creation				
IP (S)	23.18	20.67	19.34	21.06
IP (M)	12.05	13.09	9.22	<b>11.45</b>
Ours (S)	22.82	22.40	16.96	<b>20.72</b>
Combined (S)	22.40	20.40	20.24	21.01
Ours (M)	14.66	14.86	9.60	<b>13.04</b>
Combined (M)	11.57	12.24	9.06	<b>10.95</b>
Automatic creation				
IP (S)	21.20	21.13	19.86	<b>20.73</b>
IP (M)	16.01	14.28	15.99	15.42
Ours (S)	23.35	22.20	17.34	20.96
Combined (S)	21.33	20.34	20.22	<b>20.63</b>
Ours (M)	18.57	14.48	14.37	15.80
Combined (M)	12.60	13.08	14.79	13.49

**Parameter fine-tuning** Eqs. (4) and (7) use the parameter  $\alpha$  to control the contribution of the jargon, so we investigated the behavior of the model with different  $\alpha$  values. To do that, we tuned the parameter  $\alpha$  in the range of [0.1, 0.2, 0.3, 0.4, 0.5, 0.6]. Figure 3 shows the performance with different  $\alpha$  values using the Whisper medium version. We can observe that the CER slowly decreases until 0.4 and then quickly increases. The reason is that with a small value of  $\alpha$  the contribution of our decoding is tiny. When increasing  $\alpha$ , the decoding process is hallucinated by the jargon that leads to high CERs and WER. The change of per-

<sup>4</sup><https://github.com/elisa-aleman/Mecab-python>

formance on Librispeech is tiny compared to other Japanese datasets. As mentioned, the Librispeech corpus may be included in the pre-training of Whisper, that mitigates the contribution of our decoding method. Based on the observation, we therefore selected  $\alpha = 0.2$  for HGP, JNAS, and Librispeech corpora and  $\alpha = 0.4$  the GPT-Syn dataset.

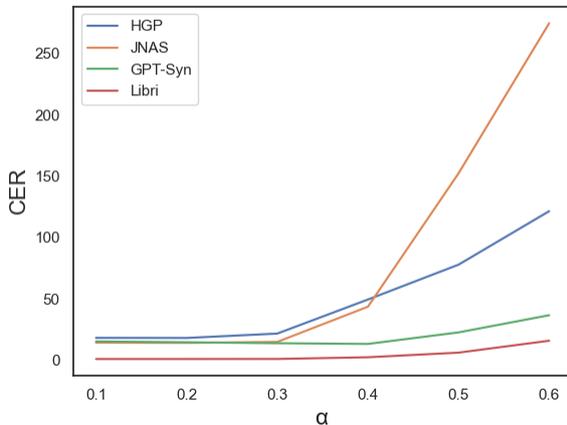


Figure 3: Parameter tuning using Whisper medium.

**Output observation** Table 4 shows the output observation of methods compared to the ground truth (translated as "It is presumed that the corrosion gap was made due to aging of refrigerant piping") on the HGP dataset using the Whisper medium version. 冷媒 means "refrigerant", and 腐食孔 means "corrosion gap". We did not show the output of the Beam Search+n-grams LM and TCPGen due to low CER and WER as in Table 2.

Table 4: The outputs from decoding methods. Blue text is correct and red text denotes incorrect tokens.

Method	Sample
Ground-truth	冷媒配管の経年劣化による腐食孔と推定
Whisper	霊媒配管の経年劣化による腐食効と推定
Initial prompt	例外配管の経年劣化による腐食孔と推定
Ours	冷媒配管の経年劣化による腐食効と推定
Combination	冷媒配管の経年劣化による腐食孔と推定

Whisper predicts two similar tokens marked by red. This is because these terms are in the high-pressure gas incident domain that does not appear in the training of Whisper. Instead, it tries to generate similar tokens used in the training process. The initial prompt can correctly recognize one token due to the use of the jargon in the prompt. It is similar to our decoding method. The combined method shows the best result that can correctly predict two tokens. It supports the results in Table 2 in which our methods output better scores than others.

## 6 Conclusion

This paper introduces a new method for improving the performance of ASR engines, i.e., Whisper by taking into account jargon. To do that, the method considers the jargon as the context injected into the decoding process. Domain-specific tokens receive more attention by adjusting the score of tokens in the beam search. Experimental results on three Japanese and one English datasets confirm two important points. First, the jargon can provide useful domain knowledge to improve the quality of Whisper. It shows that the improvement on domain-specific corpora is higher than public datasets due to the lack of domain knowledge of Whisper. Second, the combination of our decoding and the initial prompt methods achieves the best results. The proposed method provides a simple but effective way for domain adaptation of Whisper without accessing speech data for fine-tuning ASR models. Future work will confirm the efficiency of the method on other datasets and improve the decoding process using graph neural networks.

## Limitations

Even achieving promising results, the proposed method has some limitations. First, it can fail to correct tokens even if they are included in the jargon. The possible reason comes from the fact that the beam that contains these tokens receives a low score. Even boosting the score of these tokens, the beam could not receive the highest probability. Second, the initial prompt method helps to improve the overall performance. In this case, if this method is not available (ASR models to not offer it) or the number of biasing words is larger than 244 tokens, it may challenge our proposed method.

## Ethics Statement

All datasets and baseline models experimented with in this work have no unethical applications or risky broader impacts. The evaluation uses one public dataset that is widely used for ASR. For the HGP dataset, we followed the term that we could not publish the data. We really acknowledge the understanding of audiences for data publication. The dataset does not include any confidential or personal information of workers or companies. The baseline methods used for evaluation can be publicly accessed with GitHub links. There is no bias for the re-implementation or parameter selection that can affect the final results.

## References

- Uri Alon, Golan Pundak, and Tara N Sainath. 2019. Contextual speech recognition with difficult negative training examples. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6440–6444. IEEE.
- Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, et al. 2023. Seamless4t-massively multilingual & multimodal machine translation. *arXiv preprint arXiv:2308.11596*.
- William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals. 2016. Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4960–4964. IEEE.
- Chung-Cheng Chiu, Tara N Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjuli Kannan, Ron J Weiss, Kanishka Rao, Ekaterina Gonina, et al. 2018. State-of-the-art speech recognition with sequence-to-sequence models. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4774–4778. IEEE.
- Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. 2015. Attention-based models for speech recognition. *Advances in neural information processing systems*, 28.
- Linhao Dong, Shuang Xu, and Bo Xu. 2018. Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5884–5888. IEEE.
- Alex Graves. 2012. Sequence transduction with recurrent neural networks. *arXiv preprint arXiv:1211.3711*.
- 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*, pages 369–376.
- Alex Graves and Navdeep Jaitly. 2014. Towards end-to-end speech recognition with recurrent neural networks. In *International conference on machine learning*, pages 1764–1772. PMLR.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. 2020. Conformer: Convolution-augmented transformer for speech recognition.
- Minglun Han, Linhao Dong, Zhenlin Liang, Meng Cai, Shiyu Zhou, Zejun Ma, and Bo Xu. 2022. Improving end-to-end contextual speech recognition with fine-grained contextual knowledge selection. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8532–8536. IEEE.
- Wei Han, Zhengdong Zhang, Yu Zhang, Jiahui Yu, Chung-Cheng Chiu, James Qin, Anmol Gulati, Ruoming Pang, and Yonghui Wu. 2020. Contextnet: Improving convolutional neural networks for automatic speech recognition with global context.
- Yanzhang He, Tara N Sainath, Rohit Prabhavalkar, Ian McGraw, Raziq Alvarez, Ding Zhao, David Rybach, Anjuli Kannan, Yonghui Wu, Ruoming Pang, et al. 2019. Streaming end-to-end speech recognition for mobile devices. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6381–6385. IEEE.
- Yi Hu and Philippos C Loizou. 2006. Evaluation of objective measures for speech enhancement. In *Ninth international conference on spoken language processing*.
- Wen-Chin Huang, Chia-Hua Wu, Shang-Bao Luo, Kuan-Yu Chen, Hsin-Min Wang, and Tomoki Toda. 2021. Speech recognition by simply fine-tuning bert. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7343–7347. IEEE.
- Shumpei Inoue, Minh-Tien Nguyen, Hiroki Mizokuchi, Tuan-Anh Nguyen, Huu-Hiep Nguyen, and Dung Le. 2023. Towards safer operations: An expert-involved dataset of high-pressure gas incidents for preventing future failures. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 509–521.
- Katunobu Itou, Mikio Yamamoto, Kazuya Takeda, Toshiyuki Takezawa, Tatsuo Matsuoka, Tetsunori Kobayashi, Kiyohiro Shikano, and Shuichi Itahashi. 1999. Jnas: Japanese speech corpus for large vocabulary continuous speech recognition research. *Journal of the Acoustical Society of Japan (E)*, 20(3):199–206.
- Tahir Javed, Sumanth Doddapaneni, Abhigyan Raman, Kaushal Santosh Bhogale, Gowtham Ramesh, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M Khapra. 2022. Towards building asr systems for the next billion users. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10813–10821.
- Jeesu Jung, Hyein Seo, Sangkeun Jung, Riwoo Chung, Hwijung Ryu, and Du-Seong Chang. 2023. Interactive user interface for dialogue summarization. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*, pages 934–957.
- Namkyu Jung, Geonmin Kim, and Joon Son Chung. 2022. Spell my name: keyword boosted speech recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6642–6646. IEEE.

- Duc Le, Mahaveer Jain, Gil Keren, Suyoun Kim, Yangyang Shi, Jay Mahadeokar, Julian Chan, Yuan Shangquan, Christian Fuegen, Ozlem Kalinli, et al. 2021a. Contextualized streaming end-to-end speech recognition with trie-based deep biasing and shallow fusion. *arXiv preprint arXiv:2104.02194*.
- Duc Le, Gil Keren, Julian Chan, Jay Mahadeokar, Christian Fuegen, and Michael L Seltzer. 2021b. Deep shallow fusion for rnn-t personalization. In *2021 IEEE Spoken Language Technology Workshop (SLT)*, pages 251–257. IEEE.
- Daniel Li, Thomas Chen, Alec Zadikian, Albert Tung, and Lydia B Chilton. 2023. Improving automatic summarization for browsing longform spoken dialog. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–20.
- Ian McGraw, Rohit Prabhavalkar, Raziell Alvarez, Montse Gonzalez Arenas, Kanishka Rao, David Rybach, Ouais Alsharif, Haşim Sak, Alexander Grunstein, Françoise Beaufays, et al. 2016. Personalized speech recognition on mobile devices. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5955–5959. IEEE.
- Douglas O’Shaughnessy. 2024. [Trends and developments in automatic speech recognition research](#). *Computer Speech Language*, 83:101538.
- 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. IEEE.
- Gueorgui Pironkov, Sean UN Wood, and Stephane Dupont. 2020. Hybrid-task learning for robust automatic speech recognition. *Computer Speech & Language*, 64:101103.
- Golan Pundak, Tara N Sainath, Rohit Prabhavalkar, Anjuli Kannan, and Ding Zhao. 2018. Deep context: end-to-end contextual speech recognition. In *2018 IEEE spoken language technology workshop (SLT)*, pages 418–425. IEEE.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Virgile Rennard, Guokan Shang, Julie Hunter, and Michalis Vazirgiannis. 2023. Abstractive meeting summarization: A survey. *Transactions of the Association for Computational Linguistics*, 11:861–884.
- Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra. 2001. Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs. In *2001 IEEE international conference on acoustics, speech, and signal processing. Proceedings (Cat. No. 01CH37221)*, volume 2, pages 749–752. IEEE.
- Yan Song, Yuanhe Tian, Nan Wang, and Fei Xia. 2020. Summarizing medical conversations via identifying important utterances. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 717–729.
- Guangzhi Sun, Chao Zhang, and Philip C Woodland. 2021. Tree-constrained pointer generator for end-to-end contextual speech recognition. In *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 780–787. IEEE.
- 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.
- Ian Williams, Anjuli Kannan, Petar S Aleksic, David Rybach, and Tara N Sainath. 2018. Contextual speech recognition in end-to-end neural network systems using beam search. In *Interspeech*, pages 2227–2231.
- Dian Yu, Kai Sun, Claire Cardie, and Dong Yu. 2020. Dialogue-based relation extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4927–4940.
- Ding Zhao, Tara N Sainath, David Rybach, Pat Rondon, Deepti Bhatia, Bo Li, and Ruoming Pang. 2019. Shallow-fusion end-to-end contextual biasing. In *Interspeech*, pages 1418–1422.

# Optimizing Code-Switching in Conversational Tutoring Systems: A Pedagogical Framework and Evaluation

Zhengyuan Liu<sup>†\*</sup>, Stella Xin Yin<sup>‡\*</sup>, Nancy F. Chen<sup>†</sup>

<sup>‡</sup>Nanyang Technological University, Singapore

<sup>†</sup>Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR, Singapore

liu\_zhengyuan@i2r.a-star.edu.sg      nfychen@i2r.a-star.edu.sg

## Abstract

Large language models demonstrate remarkable proficiency in various tasks across multiple languages. However, their potential in code-switching remains underexplored, particularly in cultural and educational contexts. Code-switching or translanguaging plays a crucial role in bilingual education, facilitating comprehension and engagement among students with varied language proficiency. In this work, we present a pedagogy-inspired framework that introduces traditional classroom practices of code-switching to intelligent tutoring systems. Specifically, we develop fine-grained instructional strategies tailored to multilingual and educational needs. We conduct experiments involving both LLM-based evaluation and expert analysis to assess the effectiveness of translanguaging in tutoring dialogues. Our experimental results indicate that strategic code-switching can significantly enhance the learning experience. This work not only advances dialogic tutors in language learning but also extends LLMs to better accommodate multilingual interaction.

## 1 Introduction

Large Language Models (LLMs) excel in diverse tasks and various languages (Ouyang et al., 2022; Wang et al., 2023a). While their task-solving capabilities in monolingual scenarios are well-studied (Zheng et al., 2024), their potential in code-switching - the practice of alternating languages within an utterance - is still less explored (Zhang et al., 2023b). In multilingual communications, people sometimes switch languages during the conversation to convey context-specific concepts and reinforce social connections. However, current LLMs aren't specifically trained for translanguaging scenarios, highlighting an emerging research interest in their ability to understand and utilize code-switching (Doğruöz et al., 2023).

\* Equal contribution.



Figure 1: Examples of conversation segments in language learning using pedagogical code-switching.

Code-switching is not only relevant in the realm of natural language processing but also has significant implications in educational settings (Lin, 2013). For instance, bilingual children often have an imbalanced exposure to their first languages (L1) at home, and show less proficiency in vocabulary, grammar, and sentence structures of the target language (L2) at school. As a result, they tend to switch between L1 and L2 in classrooms. Accordingly, as shown in Figure 1, teachers apply code-switching strategies to clarify instructions, introduce new concepts, provide necessary encouragement, and facilitate the learning process (Rabbidge, 2019). These efforts are termed as “translanguaging” (Williams, 1994) or “code-switching” in pedagogical theories, referring to the planned and sys-

tematic use of two languages in the classroom, and it includes both intra-sentential and inter-sentential switching (Lin, 2016).

While code-switching in tutoring has been studied extensively in the last decades (De La Cruz, 2019; Turnbull and Arnett, 2002; De La Campa and Nassaji, 2009; Littlewood and Yu, 2011), most previous works focus on traditional classrooms. On the other hand, while Intelligent Tutoring Systems (ITSs) have shifted adaptive and personalized education from traditional classrooms to online learning, they are often limited to monolingual contexts and lack adaptability to mixed bilingual communication. In particular, one major challenge of improving code-switching of dialogue systems is the scarcity of data (Ahn et al., 2020; Dođruöz et al., 2023). Due to the highly colloquial nature of translanguaging, existing resources for specific domains are limited, and collecting data at a large scale requires considerable annotation efforts. Furthermore, previous studies evaluate multilingual models in general-purpose code-switching scenarios (Tan and Joty, 2021; Adilazuarda et al., 2022) and they often simulate mixed generation by replacing words in parallel corpora. However, this approach fails to capture the complexity of effective pedagogical code-switching, which requires strategic and purposeful integration that considers learners’ proficiency levels and educational objectives. Thus, pedagogical code-switching should go beyond mere word swapping to facilitate comprehension of complex concepts and provide instructional scaffolding.

In this work, we aim to improve the generated code-switching in conversational tutoring systems. We first gain insights from translanguaging theories and empirical dialogue studies to construct a pedagogical code-switching framework, where each dimension combines relevant scaffolding strategies to enhance language teaching through targeted translanguaging interventions. Our framework contributes to language learning by facilitating vocabulary acquisition, grammatical understanding, and conversational fluency. It also supports content mastery through concept clarification and emotional support.

To anchor a practical application, we conduct a case study on image description for language learning, and leveraging LLMs as tutoring agents. We utilize our proposed framework in both instructing LLMs for fine-grained code-switching generation and assessment. We conduct experiments on two

representative translanguaging cases (e.g., Chinese-English, Korean-English), and deploy automated evaluation and qualitative analysis to assess the effectiveness of pedagogical code-switching generation in LLM-based tutoring systems. Experimental results indicate that state-of-the-art models are capable of tailoring scaffolding actions and code-switching to learners’ language proficiency levels and teaching content, and strategic code-switching can significantly enhance the learning experience.

## 2 Related Work

### 2.1 Code-switching in Dialogue Systems

Dialogue systems (conversational agents) are designed to imitate various human linguistic and behavioral patterns (Chen et al., 2017), including the language mixing patterns of multilingual users (Parekh et al., 2020; Bawa et al., 2020). Rule-based approaches use linguistic features such as discourse makers and templates to produce bilingual utterances via word replacement (Ahn et al., 2020). Data-driven methods can achieve higher flexibility, beyond simple lexical borrowing to blending of languages at syntactic, grammatical, and morphological levels (Dođruöz et al., 2023; Liu et al., 2022), but they heavily rely on well-annotated data. While recent LLMs show strong multilingual capability and share certain knowledge across various languages (Wang et al., 2023a), their potential for coherent translanguaging is still underdeveloped (Zhang et al., 2023b), and the exploration to dialogic tutoring remains limited (Choi et al., 2023).

### 2.2 Code-switching in Education

Research on code-switching in classrooms dates back to the 1970s. Jacobson (1981) proposed New Concurrent Approach (NCA), as the first bilingual pedagogy, introducing flexible bilingual language practices for children and youth. The pedagogical aspects of language mixing affirm code-switching as a viable approach to bilingual teaching and learning (Hornberger and Link, 2012). The practice of code-switching (intentional instructional strategies that integrate two or more languages in real classrooms) has been well-studied (Lin, 2013; Sinclair and Fernández, 2023; García, 2009; Bon, 2021), which reveals that code-switching can significantly enhance learning and can be used as a strategic teaching method (Moore, 2002; Cenoz and Gorter, 2022; Barahona, 2020; Vaish and Subhan, 2015). Studies in language education further confirmed

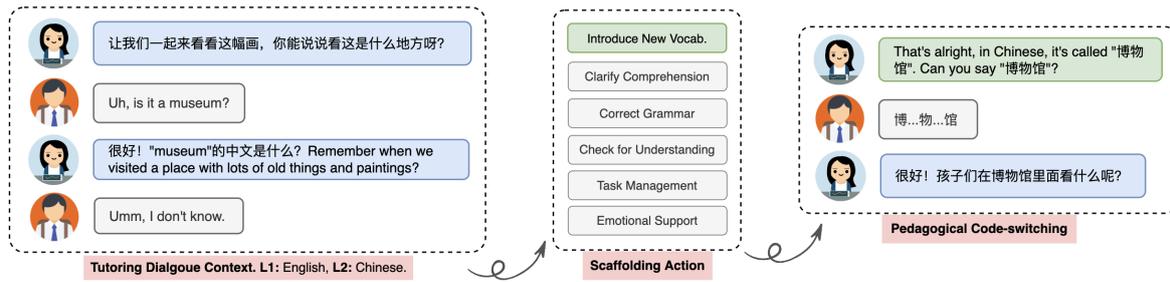


Figure 2: Decomposed pedagogical code-switching process, where the translanguaging is featured on scaffolding.

that the use of **L1** can facilitate language acquisition, improve student engagement, and establish rapport between teacher and students in **L2** learning (Pan and Pan, 2010). Besides, these studies emphasized the key role of teachers in managing teacher-centered Initiation-Response-Feedback (IRF) (Wells, 1999) sequences that promote the use of code-switching in the classroom.

### 2.3 Intelligent Tutoring Systems

The advancement of ITSs has marked a significant step forward in education practice (Graesser et al., 2018; Demszky and Hill, 2023; Wang et al., 2023b). These systems provide personalized learning experiences and instant feedback (Chaffar and Frason, 2004; Harley et al., 2015; Grivokostopoulou et al., 2017), tailored to learners' characteristics and needs (Dzikovska et al., 2014; Grawemeyer et al., 2016; Nihad et al., 2017), and are shown to positively influence students' engagement in learning and academic performance (Kulik and Fletcher, 2016; Xu et al., 2019).

Dialogue tutor is a particular type of intelligent tutoring system that interacts with students via natural language conversation (Nye et al., 2014; Ruan et al., 2019). In STEM domains, conversational ITSs can facilitate university students in problem-solving by providing real-time feedback and hints in text formats (Nye et al., 2023; Paladines and Ramirez, 2020; Arnau-González et al., 2023). However, prior work has widely relied on rule-based systems with human-crafted domain knowledge (Nye et al., 2014; Graesser et al., 2018), or data-driven approaches that require a certain amount of human annotation for supervised learning (MacLellan and Koedinger, 2022). Recently, LLMs show strong potential to build dialogue tutors with less data supervision and higher coherence (Afzal et al., 2019; Demszky and Hill, 2023; Macina et al., 2023b), and they can be further improved by integrating LLMs with pedagogical and learning science principles

(Stasaski et al., 2020; Sonkar et al., 2023; Macina et al., 2023a).

### 3 Pedagogical Code-switching

We conceptualize "pedagogical code-switching" as a combination of two aspects: scaffolding and translanguaging, as shown in Figure 3. Scaffolding is a dynamic intervention finely tuned to the learner's ongoing progress. The support given by the teacher during scaffolding strongly depends on the patterns of teacher-student interactions (Vygotsky and Cole, 1978; van de Pol et al., 2010). Therefore, the scaffolding abilities of tutors become the key criteria of effective tutoring systems. There are seven dimensions of scaffolding strategies to facilitate teaching and learning, which are (1) Feeding back, (2) Hints, (3) Instructing, (4) Explaining, (5) Modeling, (6) Questioning, (7) Social-emotional Support (Gibbons, 2015).

On the other hand, the multifaceted functions of code-switching in teaching practices highlight its strategic use in pedagogical contexts. DiCamilla and Antón (2012) identified four major functions of code-switching in classroom discussions: 1) creating, discussing, and translating content, 2) negotiating grammatical, lexical, and stylistic choices, 3) planning, defining, and managing tasks, and 4) maintaining and developing interpersonal relationships. Building on this classification, Tigert et al. (2019) expanded the framework to include five major functions of code-switching in bilingual classrooms: [A] Negotiating content, [B] Clarifying language, [C] Checking for or confirming understanding, [D] Task management, [E] Building relationships. These findings reveal how multiple ways that teachers could employ code-switching strategies to scaffold learning.

Language learning is one of the typical applications of code-switching in classrooms. The inherent nature of students' limited proficiency in the target language requires extensive scaffolding through

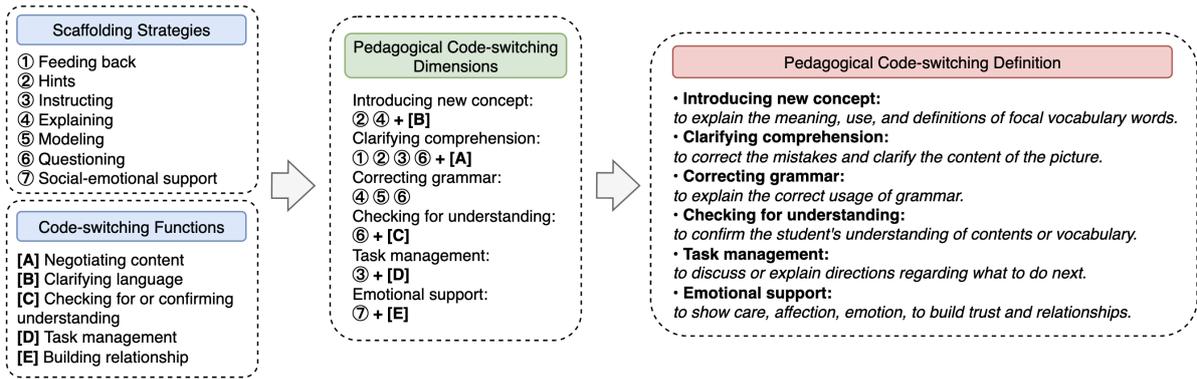


Figure 3: Conceptualizing pedagogical code-switching. Each dimension can be used as instructions for dialogic tutors, as well as the evaluation rubrics.

the use of **L1** for explaining, modeling, and providing instructions. Based on scaffolding strategies and Tigert et al. (2019)’s framework, we analyze dialogues between teachers and students in bilingual classrooms from previous work DiCamilla and Antón (2012); Tigert et al. (2019); Vaish and Subhan (2015). We build a pedagogical code-switching framework with six dimensions (as shown in Figure 3): 1) *Introducing New Concept (Vocabulary)*, 2) *Clarifying Comprehension*, 3) *Correcting Grammar*, 4) *Checking for Understanding*, 5) *Task Management*, 6) *Emotional Support*. For example, we integrate scaffolding strategies (2) Hints and (4) Explaining to code-switching function [B] Clarifying language. We rename it as Introducing New Concept and define it as “to explain the meaning, use, and definitions of focal vocabulary words.”

#### 4 Improving Tutoring Systems with Pedagogical Code-Switching

LLMs pave a new way to build dialogue tutors with less data supervision and higher customization (Macina et al., 2023b). They can be shaped along desired dimensions to mimic human conversational characteristics such as tone and personality traits to deliver better user experience (Safdari et al., 2023; Shao et al., 2023; Liu et al., 2024b). We thus leverage LLMs as a tutoring agent for language learning, and adopt pedagogical code-switching based on our proposed framework. Previous work shows that in task-specific settings, the general and coarse instruction may lead to inconsistent generation (Liu et al., 2024a). Therefore, we set up fine-grained instructions to improve and evaluate the LLM-based tutoring systems.

**Bi-lingual Setting:** The conversational tutoring system is designed for **L2** learning (**L1** is English),

and particularly focuses on the image description task. In each tutoring session, the student is presented with a picture and asked to describe the incidents in **L2** (e.g., Chinese, Spanish). From a syntactic perspective, languages can be classified into three structural patterns: 1) Subject-Verb-Object (e.g., English, Chinese); 2) Verb-Subject-Object (e.g., Arabic, Irish); 3) Subject-Object-Verb (e.g., Japanese, Korean). Given that code-switching is a structured and syntax-related phenomenon in linguistics, in our study, we select two representative translanguaging combinations (**Chinese-English**, **Korean-English**) in our language learning case study, to represent the applicability and generalizability of code-switching in different language contexts.

**Tutoring Setting:** Teaching and improving students’ **L2** acquisition through image description is a dynamic and engaging approach. In the image description tasks, the learning objective is using target sentences to describe the given image that includes a particular place or setting, people or animals, items and actions, etc. The teacher uses scaffolding and code-switching strategies to guide students step by step until they can independently complete the image description task. We build a multi-agent communication environment following previous work (Zhang et al., 2023a; Wu et al., 2023; Liu et al., 2024b).

**Teacher Role Instruction:** The teacher guides students to describe the items, emotions of people, and incidents depicted in the images, following teacher-centered IRF (Wells, 1999) sequences to promote interaction. In this process, the teacher applies scaffolding strategies, such as questioning, reformulation, and elaboration to assist learners in knowledge construction and expression (Gibbons,

Dimension	Definition	Chinese-English Code-switching Example
Introducing New Concept	use code-switching to explain the meaning, use, and definitions of focal vocabulary words.	“hungry”的中文是“饿”。我们可以说“小男孩有点饿了”，it means “The boy is hungry” (The Chinese word for “hungry” is “饿”). We can say “The little boy is a bit hungry”, it means “The boy is hungry”)
Clarifying Comprehension	use code-switching to correct the mistakes and clarify the content of the picture	再仔细看看图片， they are not just standing. Think about how we stand in line when we wait for food at school. We stand one behind the other, right? (Take a closer look at the picture, they are not just standing. Think about how we stand in line when we wait for food at school. We stand one behind the other, right?)
Correcting Grammar	use code-switching to explain the correct usage of grammar	不过在一个句子的开头，我们通常会加上主语， we say “小朋友在玩游戏”， instead of “玩游戏”。 (However, at the beginning of a sentence, we usually add the subject, we say “The children are playing games”, instead of “playing games”.)
Checking for Understanding	use code-switching to confirm students’ understanding of contents or vocabulary	用中文我们怎么说 “reading books”? (How do we say “reading books” in Chinese?)
Task Management	use code-switching to discuss or explain directions regarding what to do next	OK, now look at the right part of the picture. 你能看到有两个男孩吗? (OK, now look at the right part of the picture. Can you see two boys?)
Emotional Support	use code-switching to show care, affection, and emotion, to build trust and relationships	Great! 你做得真棒! (Great! You did a great job!)

Table 1: Definition and examples of pedagogical code-switching in bilingual language learning (L1: English, L2: Chinese). See Table 4 in the Appendix for Korean-English code-switching examples.

#### C1: Teacher Role Instruction

**[Role & Task Definition]** You are a primary school language teacher. You teach the student to describe the picture. The student’s first language (L1) is English, and target language (L2) is Chinese/Korean.  
**[Pedagogical Instruction]** You apply scaffolding and code-switching of L1 and L2 during tutoring. – Detailed pedagogical code-switch description. –  
**[Behavior Constraint]** Ask the student only one question at a time. Always wait for input before proceeding to the next step. Correct the student’s answers if they are inaccurate.

#### C2: Student Role Instruction

**[Role & Task Definition]** You are a primary school student. You are taking a language learning class, and describing the given pictures.  
**[Language Capability]** Your first language is English, and your Chinese/Korean proficiency is limited. You make some grammar errors in your responses to the teacher.

2015). Due to the students’ limited proficiency in L2, we instruct the teacher agent (as in Codebox C1) to apply six dimensions of pedagogical code-switching during the tutoring process, and add reference examples of each dimension (as in Table 1).

**Student Role Instruction:** We follow the learning process via human-machine interaction, where the tutoring system (i.e., teacher) leads the conversation, and we feed responses from a student simulator instead of the human participants. To

trigger scaffolding and code-switching strategies, we set the student role (as in Codebox C2) to include both L1 and L2, while L2 (i.e., Chinese, Korean) proficiency is low, and the student occasionally makes grammar mistakes. With the support and guidance from the teacher agent, the student is able to complete the given task, and improve L2 skills including vocabulary, organization, and fluency (de Oliveira et al., 2023).

## 5 Experimental Setup

In our preliminary study, open LLMs such as Mistral-7B and Llama-3-8B cannot follow the pedagogical instructions well, they tend to produce monolingual responses and fail to generate coherent tutoring dialogues. Therefore, our experiments are conducted on two state-of-the-art representative models: Gemini (Team et al., 2023) and GPT-4-turbo (Achiam et al., 2023).<sup>1</sup> Following previous work (Touvron et al., 2023), we adjust instructions to the chat template of each model. For tutoring dialogue generation, both teacher and student roles use the same model, and we feed the concatenated utterances for dialogue generation. We randomly sample 50 open-sourced cartoon images and use one sentence of image description as a learning target to generate 400 tutoring dialogues. The total utterance number is 9K.

<sup>1</sup>The experimented versions are Gemini-Pro-1.5 and GPT-4-turbo-0125-Preview.

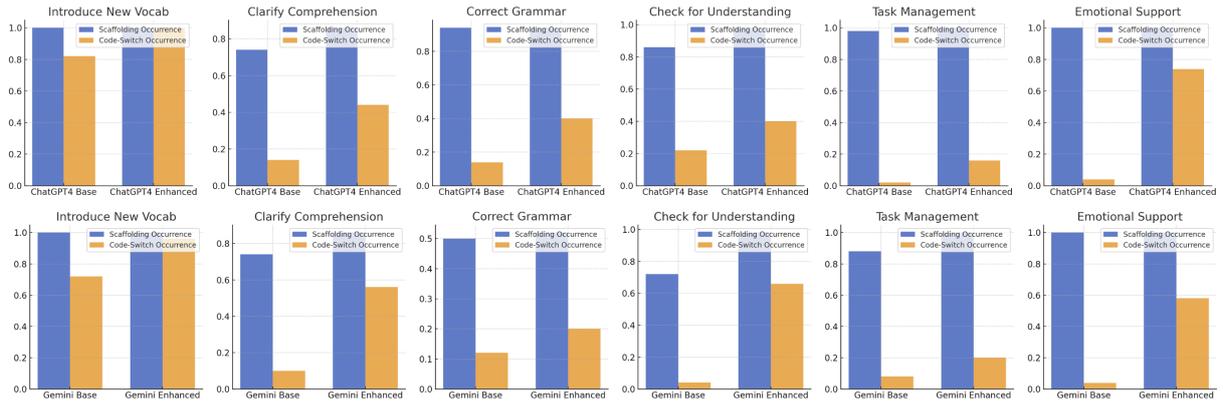


Figure 4: Quantitative results on Chinese-English pedagogical code-switching via automated evaluation.

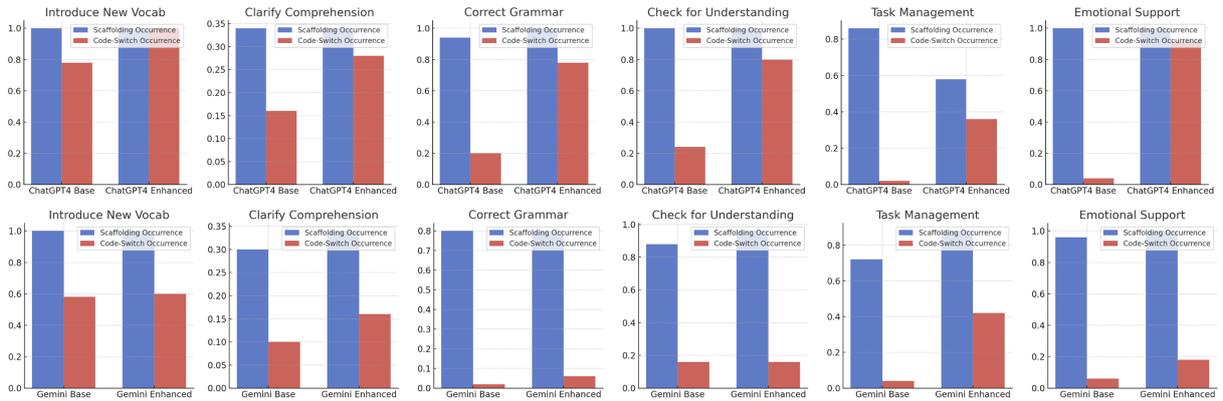


Figure 5: Quantitative results on Korean-English pedagogical code-switching via automated evaluation.

For quality assessment and data analysis, we conduct both LLM-based automated evaluation and manual rating by human experts.

**Automated Evaluation:** We adopt the LLM-as-a-judge (Saha et al., 2023) with GPT-4 to validate the effectiveness of pedagogical code-switching according to our proposed framework. Given a dialogue, the evaluator is to predict the occurrence of each dimension (e.g., introduce new vocabulary, clarify comprehension, etc.). Next, we assess whether code-switching is applied to these functions. In each dimension, one point will be added if the dialogue meets the criteria (as shown in Appendix Table 5). To build the reference label, two bilingual experts annotate the generated dialogues. Their teaching languages are Chinese and Korean. We provide each expert with an overview of the study, its objectives, and code-switching evaluation rubrics.

**Qualitative Analysis:** We invite two language teachers to complete a survey rating the pedagogical ability of the tutoring systems. The survey items are adapted from Tack and Piech (2022). We randomly select 15 dialogue segments from each

system and in both Chinese-English and Korean-English. For each dialogue segment, we ask three questions: To what extent do you think the teacher 1) speaks like a human tutor, 2) applies code-switching effectively, and 3) scaffolds the student’s learning? The annotators are asked to indicate their agreement with these three statements on a 5-point scale, ranging from 1 (strongly disagree) to 5 (strongly agree). The survey example can be found in Appendix Figure 7.

## 6 Results and Analysis

### 6.1 Code-switching Generation Evaluation

We first measure the model’s performance in generating pedagogical code-switching. In Chinese-English tutoring dialogue evaluation, GPT-4-turbo ( $p = 0.021$ ) and Gemini ( $p = 0.007$ ) show significant improvement when applying pedagogical code-switching across six dimensions, compared to the systems with the base instruction. Not surprisingly, all tested systems show a high utilization of *Introducing new vocabulary*, which is often involved with simple word replacement within a

Model	Scaffolding Label		Pedagogical CSW		Scaffolding Label		Pedagogical CSW	
	Accuracy	F1	Accuracy	F1	Accuracy	F1	Accuracy	F1
LLaMA-3-Chat 8B	83.33	81.83	74.33	73.14	80.67	76.55	51.33	50.18
Mistral-7B-Instruct-v0.2	85.33	84.52	71.50	70.09	80.33	75.10	51.17	46.07
GPT-3.5-turbo-1106	89.00	85.46	76.33	76.55	82.50	74.59	62.67	62.63
GPT-4-turbo-1106	90.83	88.79	78.50	78.33	85.16	80.49	71.17	70.33

Table 2: Model comparison of automated evaluation with LLM-as-a-judge. Scaffolding labeling is only to predict the scaffolding types regardless of translanguaging. CSW is short for code-switching. Columns in pink denote the results of English (L1) and Chinese (L2). Columns in blue denote the results of English (L1) and Korean (L2).

sentence. However, LLMs with the base instruction cannot produce diverse code-switching in other pedagogical types (e.g., *Check for Understanding*, *Task Management*, *Emotional Support*).

Moreover, LLMs perform differently across different language mixing of code-switching. For instance, in Korean-English tutoring dialogue evaluation, only GPT-4-turbo ( $p = 0.022$ ) shows significant improvement in applying code-switching strategies across six dimensions when compared to no pedagogical code-switching instructions; The Gemini does not show the same trend in improvement. Particularly, the performance on *Introducing new vocabulary*, *Correct Grammar*, and *Checking for Understanding* remains the same. We speculate that this is because LLMs’ multilingual capability differs across languages, and the code-switching beyond simple word replacement (e.g., *Correct Grammar*, *Checking for Understanding*) relies more on cross-lingual knowledge.

## 6.2 LLM-as-a-judge Model Comparison

To validate the efficacy of automated evaluation and compare the performance across open and closed-source LLMs, we use manual annotation as a reference, and results are presented in terms of correlation with human judgments (in accuracy and F1 scores). We selected and tested a list of representative models (e.g., LLaMA-3, Mistral, GPT-3.5, and GPT-4). As shown in Table 2, all tested models can provide reasonable results on labeling Chinese-English dialogues, where GPT-4 performs slightly better in the code-switching labeling. However, they all achieve a lower performance on labeling Korean-English dialogues; the scores of LLaMA and Mistral are sub-optimal. We speculate that this is due to the models’ inconsistent multilingual capability (Wang et al., 2023a).

## 6.3 Qualitative Analysis

We conduct qualitative analysis in each dimension of code-switching strategies on models’ generation

(see examples in Appendix Table 6 and Table 7):

**Introducing New concept:** The tutor alternates between L1 and L2 to explain the meanings and uses of target words. For example, when introducing “春节/설날” (Chinese New Year), it uses both Chinese and Korean to introduce the meaning of the word, “the beginning of the year according to the Chinese calendar,” and English to explain the cultural significance of Chinese New Year, as “*It’s a time for families to get together and have a big feast.*” By offering these complementary meanings across languages, students are able to construct a better understanding of the target vocabulary.

**Clarifying Comprehension:** When students make some misunderstanding, the tutor uses L1 to provide hints and guidance to help them co-construct knowledge of the picture. In addition, the tutor encourages them to show their own understanding and modify answers in L2. For example, when students misinterpret the setting of the image as a birthday party when it actually represents a Chinese New Year celebration, the tutor gives positive feedback in L2, like “你的观察很有趣，生日派对也会有很多人在一起庆祝。/ 너의 관찰이 재미있어요. 생일 파티에도 많은 사람들이 모여서 축하할 거예요”. Next, the tutor provides some hints in L1 to guide the students toward an accurate description, “*However, this is not a birthday party. Do you notice the red lanterns and Chinese characters on the wall in the picture?*”

**Correcting Grammar:** When students make grammar mistakes, the tutor follows the Subject-Verb-Object linguistic feature in Chinese/Subject-Object-Verb in Korean to explain the grammar usage. It also uses the modeling strategy to provide structured examples for the students to imitate. For instance, the tutor provides sentence structure to facilitate students to complete sentences by filling in the blanks, “We can say [Subject]이 [Action]을 위해 [Body Part]을 들고 있습니다. / [主语]将[身体部位]举起来[谓语+宾语] (e.g., *She is holding*

Translanguaging Type	GPT-4-turbo		Gemini Pro		<i>t-test</i>	
	Mean	SD	Mean	SD	<i>t</i> value	<i>p</i> value
Chinese-English	11.4	1.453	14.133	0.819	-8.976	< .001
Korean-English	11.6	2.283	10.767	1.006	1.829	0.073

Table 3: Comparison of GPT-4-turbo and Gemini on pedagogical ability.

up her hand to receive a gift.)”

**Checking for Understanding:** The tutor uses two ways to check or confirm students’ understanding. First, the “repeat-after-me” strategy is to ask students to repeat new vocabulary or sentences in **L2**. The verbal repetition process is beneficial to language acquisition. It is believed that learning occurs as a result of repeated exposure to a given stimulus (Watson, 2017). Another method is to ask students to translate vocabulary from **L1** into **L2**. For example, the tutor asks, “Do you know ‘Dinosaur’ in Chinese / Korean?”

**Task Management:** We observe that the tutor applies code-switching when giving instructions, explaining what students going to do next, or drawing attention to the objective of the task. This function serves to raise awareness of the focal words and learning tasks. For instance, “Now look at the right part of the picture. 你能看到有两个男孩吗? / 두 소년을 봤어요?”

**Emotional Support:** The tutor gives a lot of positive affirmation and encouragement to learners, such as “잘 했어요! / 太棒了! You did a great job!” Both GPT-4-turbo and Gemini perform very well in this function because LLMs are designed to understand and respond to human emotions effectively. No matter what the student’s answer is, these models primarily respond with emotional support first. While this design enhances user experience and encourages continued engagement, its consistently stable structure makes interactions feel more like with a machine rather than a human teacher.

#### 6.4 Pedagogical Ability Evaluation Results

The pedagogical ability evaluation results (shown in Table 3) indicate that both GPT-4-turbo and Gemini with code-switching strategies demonstrate strong pedagogical effectiveness.

In Chinese-English tutoring, Gemini (*Mean* = 14.113, *SD* = 0.819) outperformed GPT-4-turbo (*Mean* = 11.400, *SD* = 1.453) by generating teacher-student dialogues that are very natural and fluent ( $p < 0.001$ ). For example, in tutoring example 23, when the student observes a person in the picture wearing a white coat, the tutor provides hints, “非

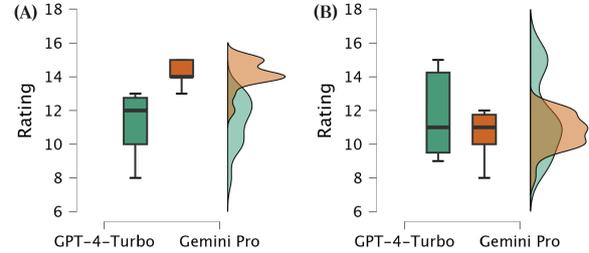


Figure 6: Descriptive statistics of human ratings on code-switching dialogue segments. (A) Chinese-English result. (B) Korean-English result.

常好! 你观察到图片中的男人穿着一件白大褂, 手里拿着一些工具 to check people’s teeth. Can you guess his job?” This example shows that the tutor builds on students’ ideas and expands their knowledge. In contrast, the GPT-4-turbo tutor is more target-driven. The conversation focuses on teaching vocabulary and guiding students in completing learning tasks. When students become distracted and answer with irrelevant words, the tutor directly corrects them and asks them back to the tasks. For example, in tutoring example 30, when the student responds incorrectly and provides irrelevant information about the picture by saying “回家了 (They are back home)”, the tutor corrects them with, “No, they are looking for shells, ‘贝壳.’”

In Korean-English tutoring, the performance of Gemini (*Mean* = 10.767, *SD* = 1.006) and GPT-4-turbo (*Mean* = 11.600, *SD* = 2.283) does not show significant difference ( $p = 0.073$ ). GPT-4-turbo follows the instructions very well, especially when correcting grammar and providing explanations. For example, in tutoring example 43, when the student makes a grammatical error by confusing the order of verbs in a sentence, the tutor advises, “Remember, in Korean, we often place the verb at the end of the sentence. You should say, ‘소년이 청소기를 사용해 방을 청소하고 있습니다.’” However, its code-switching ability is inconsistent, with ratings varying from 5 to 15 (see Figure 6). This suggests that the multilingual ability requires further enhancement to improve its reliability. In contrast, Gemini’s performance is more consistent,

although sometimes it mixes other languages (e.g., Russian, Japanese, Arabic) in utterances. For example, “아니요, these are not presents. Look closely at the box the woman is holding. ‘通常、誕生日にプレゼントを渡します。’(Usually we give presents on birthdays).”

## 7 Conclusion

In this work, we combined scaffolding strategies and translanguaging functions to propose a pedagogical code-switching framework. In a theory-inspired practice, we developed fine-grained instructional strategies tailored to multilingual learners and bilingual education needs, and leveraged LLMs as the tutoring agent and automated evaluator. Our experimental results revealed that state-of-the-art LLMs demonstrated reasonable code-switching and pedagogical ability in bilingual learning contexts. Moreover, we observed that English-centric LLMs show imbalanced performance in scaffolding, translanguaging, and pedagogical abilities across different languages, and improving cross-lingual consistency can be one of the future work. Aside from language learning, our proposed code-switching framework can also be extended to broader multilingual interactions.

## Limitations

In our experimental settings, we set the dialogic interaction in primary school language learning context, which focus mainly on basic vocabulary, grammar, and sentence structure. The translanguaging usage with advanced words and in complex syntax may pose other challenges. However, the proposed code-switching framework can be adapted to different contexts upon further refinement.

In addition, we are aware that it remains an open problem to mitigate hallucinations and biases in large language models, which may cause communication issues in human-machine interaction and computer-assisted education. Of course, current models and laboratory experiments are always limited in this or similar ways. We do not foresee any unethical uses of our proposed methods or their underlying tools, but hope that it will contribute to reducing incorrect system outputs.

## Ethics and Impact Statement

We acknowledge that all of the co-authors of this work are aware of the provided ACL Code of Ethics and honor the code of conduct. In our experiments,

models are applied under proper license. All data used in this work are only for academic research purposes and should not be used outside of academic research contexts. Our proposed methodology in general does not create a direct societal consequence and is intended to be used to improve the performance, robustness, and safety of the intelligent tutoring systems.

## Acknowledgments

This research is supported by the Agency for Science, Technology and Research (AI4EDU Programme), and the National Research Foundation, Singapore under its AISG Programme (AISG2-GC-2022-005). We thank Carolyn Lee and Geyu Lin for the research discussions. We also thank the anonymous reviewers for their precious feedback to help improve and extend this piece of work.

## References

2021. [Translanguaging in education](#). *Language Teaching*, 54(4):439–471.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Muhammad Farid Adilazuarda, Samuel Cahyawijaya, Genta Indra Winata, Pascale Fung, and Ayu Purwarianti. 2022. [IndoRobusta: Towards robustness against diverse code-mixed Indonesian local languages](#). In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 25–34, Online. Association for Computational Linguistics.
- Shazia Afzal, Tejas Dhamecha, Nirmal Mukhi, Renuka Sindhgatta, Smit Marvaniya, Matthew Ventura, and Jessica Yarbro. 2019. [Development and deployment of a large-scale dialog-based intelligent tutoring system](#). In *Proceedings of the NAACL 2019*, pages 114–121, Minneapolis, Minnesota. Association for Computational Linguistics.
- Emily Ahn, Cecilia Jimenez, Yulia Tsvetkov, and Alan Black. 2020. What code-switching strategies are effective in dialogue systems? *Society for Computation in Linguistics*, 3(1).
- Pablo Arnau-González, Miguel Arevalillo-Herráez, Romina Albornoz-De Luise, and David Arnau. 2023. A methodological approach to enable natural language interaction in an intelligent tutoring system. *Computer Speech & Language*, 81:101516.
- Malba Barahona. 2020. [The potential of translanguaging as a core teaching practice in an EFL context](#). *System*, 95:102368.

- Anshul Bawa, Pranav Khadpe, Pratik Joshi, Kalika Bali, and Monojit Choudhury. 2020. Do multilingual users prefer chat-bots that code-mix? let's nudge and find out! *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW1):1–23.
- Jasone Cenoz and Durk Gorter. 2022. [Pedagogical Translanguaging and Its Application to Language Classes](#). *RELC Journal*, 53(2):342–354.
- Soumaya Chaffar and Claude Frasson. 2004. Inducing optimal emotional state for learning in intelligent tutoring systems. In *International Conference on Intelligent Tutoring Systems*, pages 45–54. Springer.
- Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017. A survey on dialogue systems: Recent advances and new frontiers. *Acm Sigkdd Explorations Newsletter*, 19(2):25–35.
- Yunjae J Choi, Minha Lee, and Sangsu Lee. 2023. Toward a multilingual conversational agent: Challenges and expectations of code-mixing multilingual users. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–17.
- Juliane C De La Campa and Hossein Nassaji. 2009. The amount, purpose, and reasons for using L1 in L2 classrooms. *Foreign language annals*, 42(4):742–759.
- Terri De La Cruz. 2019. [Codeswitching in the Classroom](#), volume 15. Routledge.
- Luciana C. de Oliveira, Loren Jones, and Sharon L. Smith. 2023. [Interactional scaffolding in a first-grade classroom through the teaching–learning cycle](#). *International Journal of Bilingual Education and Bilingualism*, 26(3):270–288.
- Dorottya Demszky and Heather Hill. 2023. The nete transcripts: A dataset of elementary math classroom transcripts. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 528–538.
- Frederick J DiCamilla and Marta Antón. 2012. Functions of L1 in the collaborative interaction of beginning and advanced second language learners. *International Journal of Applied Linguistics*, 22(2):160–188.
- A Seza Dođruöz, Sunayana Sitaram, Barbara E Bullock, and Almeida Jacqueline Toribio. 2023. A survey of code-switching: Linguistic and social perspectives for language technologies. *arXiv preprint arXiv:2301.01967*.
- Myroslava Dzikovska, Natalie Steinhauser, Elaine Farrow, Johanna Moore, and Gwendolyn Campbell. 2014. [BEETLE II: Deep Natural Language Understanding and Automatic Feedback Generation for Intelligent Tutoring in Basic Electricity and Electronics](#). *International Journal of Artificial Intelligence in Education*, 24(3):284–332.
- Ofelia García. 2009. [Chapter 8 Education, Multilingualism and Translanguaging in the 21st Century](#), pages 140–158. Multilingual Matters, Bristol, Blue Ridge Summit.
- Pauline Gibbons. 2015. *Scaffolding language, scaffolding learning*. Heinemann.
- Arthur C Graesser, Xiangen Hu, and Robert Sottolare. 2018. Intelligent tutoring systems. In *International handbook of the learning sciences*, pages 246–255. Routledge.
- Beate Grawemeyer, Manolis Mavrikis, Wayne Holmes, Sergio Gutierrez-Santos, Michael Wiedmann, and Nikol Rummel. 2016. [Affecting off-task behaviour: how affect-aware feedback can improve student learning](#). In *Proceedings of the Sixth International Conference on Learning Analytics & Knowledge, LAK '16*, page 104–113, New York, NY, USA. Association for Computing Machinery.
- Foteini Grivokostopoulou, Isidoros Perikos, and Ioannis Hatzilygeroudis. 2017. [An Educational System for Learning Search Algorithms and Automatically Assessing Student Performance](#). *International Journal of Artificial Intelligence in Education*, 27(1):207–240.
- Jason M. Harley, François Bouchet, M. Sazzad Husain, Roger Azevedo, and Rafael Calvo. 2015. [A multi-componential analysis of emotions during complex learning with an intelligent multi-agent system](#). *Computers in Human Behavior*, 48:615–625.
- Nancy H Hornberger and Holly Link. 2012. Translanguaging and transnational literacies in multilingual classrooms: A biliteracy lens. *International journal of bilingual education and bilingualism*, 15(3):261–278.
- R. Jacobson. 1981. The implementation of a bilingual instructional model: The new concurrent approach. In R. V. Padilla, editor, *Ethnoperspectives in Bilingual Education Research, Vol. 3: Bilingual Education Technology*, pages 14–29. Eastern Michigan University, Ypsilanti, MI.
- James A. Kulik and J. D. Fletcher. 2016. [Effectiveness of Intelligent Tutoring Systems: A Meta-Analytic Review](#). *Review of Educational Research*, 86(1):42–78.
- Angel Lin. 2013. [Classroom code-switching: three decades of research](#). *Applied Linguistics Review*, 4(1):195–218.
- Angel Lin. 2016. [Code-Switching in the Classroom: Research Paradigms and Approaches](#). In Kendall A. King, Yi-Ju Lai, and Stephen May, editors, *Research Methods in Language and Education*, May, pages 1–15. Springer International Publishing, Cham.
- William Littlewood and Baohua Yu. 2011. First language and target language in the foreign language classroom. *Language teaching*, 44(1):64–77.

- Zhengyuan Liu, Shikang Ni, Ai Ti Aw, and Nancy F. Chen. 2022. [Singlish message paraphrasing: A joint task of creole translation and text normalization](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3924–3936, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Zhengyuan Liu, Stella Xin Yin, Carolyn Lee, and Nancy F Chen. 2024a. Scaffolding language learning via multi-modal tutoring systems with pedagogical instructions. *arXiv preprint arXiv:2404.03429*.
- Zhengyuan Liu, Stella Xin Yin, Geyu Lin, and Nancy F Chen. 2024b. Personality-aware student simulation for conversational intelligent tutoring systems. *arXiv preprint arXiv:2404.06762*.
- Jakub Macina, Nico Daheim, Sankalan Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023a. Mathdial: A dialogue tutoring dataset with rich pedagogical properties grounded in math reasoning problems. In *Findings of EMNLP 2023*, pages 5602–5621.
- Jakub Macina, Nico Daheim, Lingzhi Wang, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023b. Opportunities and challenges in neural dialog tutoring. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2357–2372.
- Christopher J MacLellan and Kenneth R Koedinger. 2022. Domain-general tutor authoring with apprentice learner models. *International Journal of Artificial Intelligence in Education*, 32(1):76–117.
- Danièle Moore. 2002. [Code-switching and Learning in the Classroom](#). *International Journal of Bilingual Education and Bilingualism*, 5(5):279–293.
- Elghouch Nihad, En-naimi El Mokhtar, and Yassine Zaoui Seghroucheni. 2017. [Analysing the outcome of a learning process conducted within the system als\\_corr\(lp\)](#). *International Journal of Emerging Technologies in Learning (IJET)*, 12(03):pp. 43–56.
- B Nye, Dillon Mee, and Mark G Core. 2023. Generative large language models for dialog-based tutoring: An early consideration of opportunities and concerns. In *AIED Workshops*.
- Benjamin D Nye, Arthur C Graesser, and Xiangen Hu. 2014. Autotutor and family: A review of 17 years of natural language tutoring. *International Journal of Artificial Intelligence in Education*, 24:427–469.
- 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.
- José Paladines and Jaime Ramirez. 2020. A systematic literature review of intelligent tutoring systems with dialogue in natural language. *IEEE Access*, 8:164246–164267.
- Yi-chun Pan and Yi-ching Pan. 2010. The use of L1 in the foreign language classroom. *Colombian Applied Linguistics Journal*, 12(2):87–96.
- Tanmay Parekh, Emily Ahn, Yulia Tsvetkov, and Alan W Black. 2020. Understanding linguistic accommodation in code-switched human-machine dialogues. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 565–577.
- Michael Rabbidge. 2019. *Translanguaging in EFL contexts: A call for change*. Routledge.
- Sherry Ruan, Liwei Jiang, Justin Xu, Bryce Joe-Kun Tham, Zhengneng Qiu, Yeshuang Zhu, Elizabeth L Murnane, Emma Brunskill, and James A Landay. 2019. Quizbot: A dialogue-based adaptive learning system for factual knowledge. In *Proceedings of the 2019 CHI conference on human factors in computing systems*, pages 1–13.
- Mustafa Safdari, Greg Serapio-García, Clément Crepy, Stephen Fitz, Peter Romero, Luning Sun, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2023. Personality traits in large language models. *arXiv preprint arXiv:2307.00184*.
- Swarnadeep Saha, Omer Levy, Asli Celikyilmaz, Mohit Bansal, Jason Weston, and Xian Li. 2023. Branch-solve-merge improves large language model evaluation and generation. *arXiv preprint arXiv:2310.15123*.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. *arXiv preprint arXiv:2310.10158*.
- Arabella J. Sinclair and Raquel Fernández. 2023. [Alignment of code switching varies with proficiency in second language learning dialogue](#). *System*, 113(November 2022):102952.
- Shashank Sonkar, Naiming Liu, Debshila Mallick, and Richard Baraniuk. 2023. Class: A design framework for building intelligent tutoring systems based on learning science principles. In *Findings of EMNLP 2023*, pages 1941–1961.
- Katherine Stasaski, Kimberly Kao, and Marti A Hearst. 2020. Cima: A large open access dialogue dataset for tutoring. In *Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 52–64.
- Anaïs Tack and Chris Piech. 2022. [The AI Teacher Test: Measuring the Pedagogical Ability of Blender and GPT-3 in Educational Dialogues](#). *Proceedings of the 15th International Conference on Educational Data Mining, EDM 2022*.

- Samson Tan and Shafiq Joty. 2021. [Code-mixing on sesame street: Dawn of the adversarial polyglots](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3596–3616, Online. Association for Computational Linguistics.
- 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*.
- Johanna Tigert, James Groff, Melinda Martin-Beltrán, Megan Madigan Percy, and Rebecca Silverman. 2019. Exploring the pedagogical potential of translanguaging in peer reading interactions. In *Codeswitching in the Classroom*, pages 64–87. Routledge.
- 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*.
- Miles Turnbull and Katy Arnett. 2002. Teachers’ uses of the target and first languages in second and foreign language classrooms. *Annual review of applied linguistics*, 22:204.
- Viniti Vaish and Aidil Subhan. 2015. [Translanguaging in a reading class](#). *International Journal of Multilingualism*, 12(3):338–357.
- Janneke van de Pol, Monique Volman, and Jos Beishuizen. 2010. [Scaffolding in Teacher–Student Interaction: A Decade of Research](#). *Educational Psychology Review*, 22(3):271–296.
- Lev Semenovich Vygotsky and Michael Cole. 1978. *Mind in society: Development of higher psychological processes*. Harvard university press.
- Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao, Yang Ding, Ai Ti Aw, and Nancy F Chen. 2023a. Seaeval for multilingual foundation models: From cross-lingual alignment to cultural reasoning. *arXiv preprint arXiv:2309.04766*.
- Rose E Wang, Qingyang Zhang, Carly Robinson, Susanna Loeb, and Dorottya Demeszky. 2023b. Step-by-step remediation of students’ mathematical mistakes. *arXiv preprint arXiv:2310.10648*.
- John B Watson. 2017. *Behaviorism*. Routledge.
- Gordon Wells. 1999. Language and education: Reconceptualizing education as dialogue. *Annual Review of Applied Linguistics*, 19:135–155.
- C Williams. 1994. An evaluation of teaching and learning methods in the context of secondary education. *Unpublished Doctoral Dissertation, University of Bangor*.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Auto-gen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.
- Zhihong Xu, Kausalai Wijekumar, Gilbert Ramirez, Xueyan Hu, and Robin Irey. 2019. [The effectiveness of intelligent tutoring systems on K-12 students’ reading comprehension: A meta-analysis](#). *British Journal of Educational Technology*, 50(6):3119–3137.
- Jintian Zhang, Xin Xu, and Shumin Deng. 2023a. Exploring collaboration mechanisms for llm agents: A social psychology view. *arXiv preprint arXiv:2310.02124*.
- Ruochen Zhang, Samuel Cahyawijaya, Jan Christian Blaise Cruz, Genta Winata, and Alham Aji. 2023b. Multilingual large language models are not (yet) code-switchers. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12567–12582.
- 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.

Dimension	Definition	Korean-English Code-switching Example
Introduce New Vocabulary	use code-switching to explain the meaning, use, and definitions of focal vocabulary words.	There is a special celebration in Korea called “설날” which is kind of like New Year’s Eve. It’s a time for families to get together and have a big feast. (There is a special celebration in Korea called “Seollal” which is kind of like New Year’s Eve. It’s a time for families to get together and have a big feast.)
Clarify Comprehension	use code-switching to correct the mistakes and clarify the content of the picture	아니요, 생일은 아니에요. Do you see the red decorations and Chinese characters in the picture? (No, it’s not the birthday party. Do you see the red decorations and Chinese characters in the picture?)
Correct Grammar	use code-switching to explain the correct usage of grammar	거의 맞았어요, but let’s correct the grammar a bit. Remember, sentences in Korean usually follow the Subject-Object-Verb structure. For example, “가족이 저녁 식사를 합니다.” (You almost got it, but let’s correct the grammar a bit. Remember, sentences in Korean usually follow the Subject-Object-Verb structure. For example, “The family is having dinner.”)
Check for Understanding	use code-switching to confirm students’ understanding of contents or vocabulary	너무 잘했어요. Can you try saying it one more time?“ (You did very well. Can you try saying it one more time?)
Task Management	use code-switching to discuss or explain directions regarding what to do next	오른쪽 부분을 좀 보세요. Did you see two boys? (Please look at the right part of the picture. Did you see two boys?)
Emotional Support	use code-switching to show care, affection, emotion, to build trust and relationships	잘 했어요! You did a great job! (Great! You did a great job!)

Table 4: Definition and examples of pedagogical code-switching in bilingual language learning (L1: English, L2: Korean). See Table 1 for Chinese-English code-switching examples.

I will give you tutoring dialogues for teaching students to learn Chinese via image description.

Please evaluate the code-switching of the Chinese teacher based on the following rubrics.

Add 1 point of each dimension only when the teacher properly uses English for code-switching.

**Code-switching Evaluation Rubrics:**

1. Introduce New Vocabulary:

- Definition: adopt code-switching to explain the meaning, use, and definitions of focal vocabulary words.

- Example: "hungry의 中文是‘饿’。我们可以说‘小男孩有点饿了’。”

2. Clarify Comprehension:

- Definition: adopt code-switching to correct the mistakes and clarify the content of the picture.

- Example: "再看看图片，再仔细看看图片。Think about how we stand in line when we wait for food at school. We stand one behind the other, right?"

3. Correct Grammar:

- Definition: adopt code-switching to correct grammar or explain the usage of grammar.

- Example: "我们通常会加上主语，we say ‘Children are playing games’, instead of ‘playing games’."

4. Check for Understanding:

- Definition: adopt code-switching to confirm the student’s understanding of ideas or vocabulary.

- Example: "用中文我们怎么说‘reading books’?"

5. Task Management:

- Definition: adopt code-switching to discuss or explain directions regarding what to do next.

- Example: "OK, now look at the right part of the picture. 你能看到有两个男孩吗?"

6. Emotional Support:

- Definition: adopt code-switching to show care, affection, emotion, to build trust and relationships.

- Example: "Great!你做得真棒!", "That’s ok,我们一起来看看."

Table 5: The instruction example of leveraging LLMs for automated pedagogical code-switching evaluation.

---

**Generated Tutoring Conversation Example (GPT-4-0125-Preview):**

Teacher: 小朋友，看这张图片，这是在什么地方呢？

Student: I see a big room with books.

Teacher: 很好！这个地方有很多书，我们中文叫它什么？Do you remember we talked about a place where we borrow books?

Student: Library? 书...店？

Teacher: 是的，书店里也有很多书，But this is a place where the books are for borrowing instead of for sale. 中文里我们说“图书馆”。你能试着用中文说吗？

Student: 图书...馆？

Teacher: 非常好！图书馆。现在看看图片里的人都在做什么？

Student: They are quiet.

Teacher: 对，他们都很安静。记得“安静”这个词吗？

Student: (No response.)

Teacher: 没关系。安静就是没有声音，very quiet. Can you say “安静”？

Student: 安静。

Teacher: 太好了！现在我们试着用中文描述老师在图书馆里做的事情？

Student: 老师... quiet in 图书馆。

Teacher: 接近了，Let's put it into a full sentence. 我们用“让”来表示老师让孩子们保持安静 The teacher asked the children to keep quiet. 试试看，用完整的句子。

Student: 老师在图书馆让孩子们安静。

Teacher: 非常接近了！但记得我们要用“保持安静” keep quiet。这样句子会更准确。再试一次？

Student: 老师在图书馆里让孩子们保持安静。

Teacher: 完美！You did very well, 我很高兴看到你的进步。我们继续学习，好吗？

Student: 好的，谢谢老师！

---

**Generated Tutoring Conversation Example (Gemini Pro):**

Teacher: 小朋友，看看这张图片，上面画了什么呀？

Student: I see a boy in a restaurant.

Teacher: 很好！你能用中文描述一下他吗？

Student: 他在餐厅。

Teacher: 非常好！他看起来怎么样？

Student: He is hungry.

Teacher: 还记得我们学过的“hungry”的中文吗？

Student: 我不知道。

Teacher: No worries! “hungry”的中文是“饿”。我们可以说“小男孩饿了”，it means "The boy is hungry"。Now, look closely. 他正在做什么呢？

Student: He is waiting in line.

Teacher: “waiting in line”的中文怎么说？

Student: I don't know.

Teacher: Think about how we stand in line when we wait for food at school. We stand one behind the other, right?

Student: Yes.

Teacher: In Chinese, we say “排队”. Now, can you say “小男孩有点饿了，他在餐厅排队”？

Student: 小男孩有点饿了，他在餐厅排队。

Teacher: 太棒了！但是，你看图片中，小男孩站在队伍的前面，而不是后面。He is standing in front of the line instead of the behind. 这是什么意思呢？

Student: I don't know.

Teacher: When we want to get something quickly, but we don't wait our turn, we "cut in line". In Chinese, we say “插队”. Can you use “插队” in your sentence?

Student: 小男孩太饿了，他在餐厅插队。

Teacher: 非常好！You said very clear。你觉得小男孩插队好吗？

Student: No.

Teacher: 为什么呢？

Student: Because it's not fair.

Teacher: 没错！插队是不好的行为。我们应该遵守规则，排队等候。Now, how would you describe this picture?

Student: 小男孩有点饿了，他在餐厅插队，但他不应该这样做。

Teacher: 你做得真棒！You did great job! 你不仅学会了用中文描述图片，而且还表达了自己的观点。继续努力，好吗？

Student: 好的，谢谢老师！

---

Table 6: Two generated dialogue examples with pedagogical code-switching (L1: English, L2: Chinese).

---

**Generated Tutoring Conversation Example (GPT-4-0125-Preview):**

Teacher: 안녕하세요! 오늘 우리는 이 사진을 한국어로 설명할 거예요. Take a look at this picture and see if you can tell me what's happening.  
Student: (Looks at the picture) ...Umm... I'm not sure.  
Teacher: 오른쪽 부분을 좀 보세요. Look at the right part of the picture. Do you see the person in the white coat?  
Student: Yes, he looks like a doctor.  
Teacher: 네, 정확해요! Yes, exactly! The man in the white coat is most likely a dentist. Do you know how to say "dentist" in Korean?  
Student: Dentist... Uh... 치과의사?  
Teacher: 네, 아주 좋아요! "치과의사" is exactly the word we use for dentist. It is the name for "tooth doctor." What do you think he might be doing?  
Student: He's checking something... maybe?  
Teacher: 네, 맞아요! Yes, that's right! He's examining something. Can you guess what he might be examining?  
Student: (Hesitates) Uh... hair?  
Teacher: 아니요, hair 아니예요. Look at the boy's mouth. What do you see there?  
Student: His mouth is open... and there's a tool in it.  
Teacher: Exactly! The dentist is using a tool to look at the boy's teeth. What Korean word do we use for "teeth"?  
Student: (Stays silent)  
Teacher: We use "치아" for teeth. It's similar to the English word "teeth" but pronounced a little differently.  
Teacher: So, if we put it all together, what sentence can we make to describe what's happening in the picture?  
Student: (Looks at the picture again) ... I don't know.  
Teacher: No problem! Let's try building the sentence together. The dentist is the 치과의사 and he's looking at the boy's 치아. Where are they now?  
Student: At the clinic, I guess?  
Teacher: 네, 좋아요! Yes, good! We can say "클리닉" for clinic. So how can we put it all together in a sentence?  
Student: The dentist... is looking at... the boy's teeth... at the clinic?  
Teacher: Almost there! Can you try making it a complete sentence in Korean? For example, we can say "가다 (gada)" which means "to go" to make a sentence like "학교에 가요" which means "I go to school."  
Student: 치과의사가... 클리닉에서... 치아를... I'm not sure how to say "looking at" in Korean.  
Teacher: That's okay! There are a few ways to say "looking at." One way is "검사합니다" which means "to examine." So, how could we complete the sentence.  
Student: 치과의사가 클리닉에서 어린이의 치아를 검사합니다.  
Teacher: 너무 잘했어요! Very good! You described the picture perfectly! Keep practicing like this and you'll be a Korean description pro in no time!

---

**Generated Tutoring Conversation Example (Gemini Pro):**

Teacher: 오늘 우리는 이 사진을 한국어로 설명할 거예요. 사진을 먼저 보세요. What do you see in the picture?  
Student: Uh, I see a... a big animal and a car?  
Teacher: 맞아요, very good! 그 큰 동물, what is it?  
Student: Oh, is it... giraffe?  
Teacher: what is 'giraffe' in Korean?  
Student: (No response.)  
Teacher: It's OK, we call it '기린'! And how about that machine near '기린'? Do you remember that word?  
Student: I know it! It is 'Car'!  
Teacher: Perfect! How to say it in Korean?  
Student: Is it '자동차'?  
Teacher: 맞아요! Now, let's try to make a sentence by using '기린' and '자동차'. Can you try saying it?  
Student: 기린이 해변에서... 자동차를... umm, 운전... 습니다.  
Teacher: 거의 맞았어요, but let's correct the grammar a bit. You were very close. We say "운전하고 있습니다." It means "is driving." Can you try it again?  
Student: 기린이 해변에서 자동차를 운전하고 있습니다.  
Teacher: 완벽해요! Very good! Do you want to describe any other details in the picture? Maybe what the weather looks like? 날씨가 어떻게 보이나요?  
Student: It looks... sunny? No clouds?  
Teacher: 정확해요, 맑은 날씨네요. Sunny day! 이제 '맑은 날씨'라고 말해 볼까요?  
Student: 맑은... 날씨...  
Teacher: 잘 했어요! You're doing really well. It's important to keep practicing like this. Remember, the more you speak, the better you will get.  
Student: I will try, thank you!

---

Table 7: Two generated dialogue examples with pedagogical code-switching (L1: English, L2: Korean).

Here is a tutoring dialogue segment between a student and a teacher. Please evaluate the teacher's pedagogical ability from three aspects.

Teacher: 很好! 这个地方有很多书, 我们中文叫它什么? Do you remember we talked about a place where we borrow books?

Student: Library? 书...店?

Teacher: 是的, 书店里也有很多书, But this is a place where the books are for borrowing instead of for sale. 中文里我们说“图书馆”。你能试着用中文说吗?

Student: 图书... 馆?

Q1: I think this dialogue is happened between a human tutor and a student.

Strongly agree	Disagree	Neither agree or disagree	Agree	Strongly agree
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

Q2: I think the teacher applies code-switching effectively.

Strongly agree	Disagree	Neither agree or disagree	Agree	Strongly agree
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

Q3: I think the teacher scaffolds student in language learning.

Strongly agree	Disagree	Neither agree or disagree	Agree	Strongly agree
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

Figure 7: The survey form used for manual pedagogical ability evaluation.

# ECoh: Turn-level Coherence Evaluation for Multilingual Dialogues

John Mendonça<sup>1,2</sup>, Isabel Trancoso<sup>1,2</sup> and Alon Lavie<sup>3,4</sup>

<sup>1</sup> INESC-ID, Lisbon

<sup>2</sup> Instituto Superior Técnico, University of Lisbon

<sup>3</sup> Carnegie Mellon University, Pittsburgh

<sup>4</sup> Phrase, Pittsburgh

{john.mendonca, isabel.trancoso}@inesc-id.pt, alavie@cs.cmu.edu

## Abstract

Despite being heralded as the new standard for dialogue evaluation, the closed-source nature of OpenAI’s GPT-4 model poses challenges for the research community. Motivated by the need for lightweight, open source, and multilingual automated dialogue evaluators, this paper introduces GENRESCO (Generated Responses targeting Coherence). GENRESCO is a novel LLM-generated dataset comprising over 130k negative and positive responses and accompanying explanations seeded from XDailyDialog and XPersona covering English, French, German, Italian, and Chinese. Leveraging GENRESCO, we propose ECOH<sup>1</sup> (Evaluation of Coherence), a family of evaluators trained to assess response coherence across multiple languages. Experimental results demonstrate that ECOH achieves multilingual coherence detection capabilities superior to the teacher model (GPT-3.5-Turbo) on GENRESCO, despite being based on a much smaller architecture. Furthermore, the explanations provided by ECOH closely align in terms of quality with those generated by the teacher model.

## 1 Introduction

With LLMs showcasing impressive reasoning and dialogue understanding capabilities vastly superior to any prior NLP technologies, human evaluation has more recently been complemented with automatic evaluations using GPT-4 (OpenAI, 2024). However, GPT-4 as an automated evaluator has its downsides. Perhaps the main downside is it being a closed source model hidden behind a paid API, making accessibility difficult for those outside the coverage area and lacking extensive financial resources, while also lacking transparency in its development. In contrast, and to the best of our knowledge, the study of open source alternatives to GPT-4 based dialogue evaluation is mostly limited

<sup>1</sup>Pronounced "Echo".

### Context:

**A:** Hello. I bought a China dress in your shop this morning.

**B:** Yes?

**A:** I bought it one size up by mistake.

**B:** Oh, did you?

**A:** I wonder if you can change it to one size down.

**Reference Response. B:** Yes, of course. Will you come with the receipt ?

*"The response acknowledges the request and offers a solution to accommodate the customer’s needs. The answer is Yes."*

**Random Negative Sample. B:** I’d like a book about law.

*"The response is completely unrelated to the situation discussed. The answer is No."*

**GENRESCO Positive Sample. B:** Absolutely, if you bring the dress back to the shop with the receipt, we can exchange it for a smaller size for you.

*"The response acknowledges the request and offers a solution to the problem. The answer is Yes."*

**GENRESCO Negative Sample. B:** Oh, that’s great to hear! I hope it fits perfectly.

*"The response does not acknowledge the request for a size change and instead expresses an unrelated sentiment. The answer is No."*

Table 1: Example of automatically generated negative samples obtained with random response selection, and obtained from GENRESCO (§2), our proposed dataset. The explanations are generated using one of our proposed models, ECOH-4B-ML (§3).

to the benchmarking of open source and open access LLMs or finetuning with dialogue data (Huynh et al., 2023; Zhang et al., 2023, 2024). These works suggest that LLMs struggle to outperform older encoder-based metrics trained using negative sampling approaches for relevance (e.g. random response selection). However, it is important to point out that these benchmarks have several limitations.

First and foremost, the high performance of these encoder-based models can be explained by the fact that the benchmarks themselves are based on old generative models that exhibit relevance issues that are easy to detect. For instance, in Table 1, metrics trained using random negative sampling strategies for relevance will output a positive score to all

responses except the random negative one. As such, these metrics struggle to evaluate contemporary chatbots, since these typically output fluent and semantically relevant responses.

Furthermore, only a select few benchmarks are multilingual. Whilst there is work that attempts to evaluate the multilingual capabilities of dialogue evaluation metrics (Mendonca et al., 2023; Zhang et al., 2023), they use translated benchmarks. This assumes that critical errors typically produced by these older models (e.g. irrelevance), are not influenced by language. However, more complex quality aspects such as coherence may have nuances that make them unique to certain cultures. Depending on the context, some culture specific details may or may not be implicitly inferred (Hall, 1959).

These key observations motivate our work. In order to move towards the development of metrics that evaluate dialogue coherence and are multilingual, we propose GENRESCO (Generated Responses targeting Coherence), a collection of positive and negative responses focused on coherence. Our dataset, generated using strong LLMs, contains over 130k responses in different languages (English, French, German, Italian, and Chinese), together with their corresponding explanations (in English). By prompting an LLM, we are able to (1) obtain positive samples that are in distribution (LLMs frequently output more verbose responses than their human counterparts); (2) obtain negative samples that remain semantically relevant but contain coherence and logical consistency issues, which may be more informative during training, and that are more representative of current limitations of LLMs.

With this dataset, we train a family of evaluators we call ECOH (Evaluation of Coherence)<sup>2</sup>. Our results demonstrate that distilling Coherence knowledge from a strong LLM allows us to obtain multilingual coherence detection performance of **.945** F1 score using a 0.5B model, which is superior to both the teacher models’ (GPT-3.5-Turbo) **.910** and a much larger model of the same family (QWEN1.5-7B-CHAT - **0.825**). Furthermore, the explanations provided by ECOH are of higher quality than QWEN1.5-7B-CHAT, scoring an average of over **4** out of 5 on most instances, as reported by GPT-4 evaluations.

<sup>2</sup>[github.com/johndmendonca/Ecoh](https://github.com/johndmendonca/Ecoh)

## 2 GENRESCO responses dataset

This section introduces GENRESCO, a multilingual, large-scale response collection that targets coherence, seeded from well established dialogue datasets (§2.1), and generated using LLMs (§2.2). Table 1 presents an example from this dataset. For additional examples of this dataset in other languages, see Appendix A.2.

### 2.1 Dataset Sources

Our work leverages two distinct dataset sources: XDailyDialog (Liu et al., 2023b) and XPersona (Lin et al., 2021). For training, development and testing, we use XDailyDialog, a multilingual extension of DailyDialog with human translations covering German-DE, Italian-IT and Chinese-ZH. XDailyDialog includes 13K parallel dialogues, amounting to 52K dialogues and 410K utterances. During our pre-processing step we noted a substantial overlap of dialogues between the provided test and training/validation sets of XDailyDialog. As a result, we excluded these dialogues (amounting to 20%) from the test set.

In order to gauge the extensibility to other dialogue datasets and languages, we additionally include XPersona data in our GENRESCO test set. XPersona is a multilingual extension of the PersonaChat dataset (Zhang et al., 2018) with human revised machine translations for six languages. Besides English, we include Italian-IT, Chinese-ZH, and an additional unseen language, French-FR, in our experiments. For each language, we extract 1K contexts from the test set for response generation.

For contrastive comparison, we also use DailyDialog++ (Sai et al., 2020), a similar curation effort which uses the original DailyDialog dataset, and where annotators were asked to create five additional relevant responses and five adversarial irrelevant responses for each context.

### 2.2 Generation

**Development set** We leverage GPT-3.5-turbo<sup>3</sup> (Ouyang et al., 2022) as the strong LLM to generate, given prior dialogue context, a positive and a negative response, paired with a brief explanation of the issue (or lack thereof). Each response pair is generated given a context of at least 2 turns up until the length of the dialogue except the last turn (this ensures the response is generated from a still

<sup>3</sup>`gpt-3.5-turbo-0125` and `gpt-4-1106-preview` accessed via OpenAI’s API in early April.

Dataset	Size (# contexts)	Response Avg. length	Explanation Avg. length	Response MTLD
DailyDialog++ (2020)				
Random	9,259/1,028/1,142	9.40	-	169.94
Adversarial		10.70	-	186.42
GENRESCO-H-DEV				
DailyDialog-LATIN	51,873/5,080	14.74	15.03	105.03
DailyDialog-ZH		23.06	14.54	54.38
GENRESCO-H-TEST				
DailyDialog-LATIN	4,770	14.82	26.27	155.28
DailyDialog-ZH		24.79	26.04	69.03
PersonaChat-LATIN	1,000	15.37	27.89	204.61
PersonaChat-ZH		28.81	27.78	76.66

Table 2: Comparison of statistics for different negative sample datasets. DD denotes XDailyDialog, PC XPersona. Dataset size denotes the number of unique contexts from which responses were obtained for training/validation/test subsets. MTLD denotes the Measure of Textual Lexical Diversity (McCarthy, 2005) of the responses. We report statistics for Latin script languages (denoted LATIN and covering EN,DE,FR,IT), separated from Chinese-ZH. For Average length, LATIN is calculated using words, whereas ZH uses characters.

ongoing conversation). We set the temperature to 0.7, the top- $p$  to 1, and the maximum number of tokens to 300, thereby enforcing smaller explanations which in turn should reduce inference costs. Despite sharing the same contexts, the responses and corresponding explanations are not necessarily translations of the English subset. This allows the model to freely generate responses that are more likely to occur (for the positive samples) or more representative of coherence issues in that particular language, instead of being a translation from English. The prompt used for this generation is included in Appendix A.1.

**Test set** For testing, we employ GPT-4 (OpenAI, 2024) to ensure higher quality outputs and reduce in-distribution biases from the training set. GPT-4 has been shown to match human annotations on quality, from general NLP tasks to highly specialised fields (West et al., 2022; Raunak et al., 2023; Savelka et al., 2023).

**Human validation** In order to verify the outputs of GPT-4, we additionally conduct a human validation step involving one expert linguist from each language. We randomly sample 100 examples from the XDailyDialog test set, and report an appropriateness rate that exceeds 97%, thus validating the response and explanation generation process using GPT-4. Details regarding human validation are provided in Appendix A.3.

### 2.3 Statistical Analysis

We present relevant statistics for our dataset, together with DailyDialog++ in Table 2. Since the test set for our dataset is generated by GPT-4, we opt to present the statistics separately.

Firstly, despite GENRESCO-H boasting a much larger context set, amounting to 51k/5k for training/validation, each context only has a single positive and negative response, whereas DailyDialog++ contains 5 positive responses and an additional 5 adversarial negative responses.

When comparing the average length of responses, we note that GENRESCO-H responses are longer than the human curated responses of DailyDialog++. This verbosity is a known behaviour by LLMs, since they are conditioned to output longer responses due to the Reinforcement learning from human feedback (RLHF) step, at least when compared to humans (Kamalloo et al., 2023). Additionally, we note that the response lengths remain similar across the development and test sets, whereas the explanations are much longer in the test set.

For a more fine grained analysis of the responses, we measure their lexical diversity using the Measure of Textual Lexical Diversity (MTLD) (McCarthy, 2005).<sup>4</sup> Since DailyDialog++ contains 5 responses per context, we calculate the average diversity when considering the responses individually. We observe that the diversity of human responses for DailyDialog++ is larger than the ones generated by GPT-3.5-Turbo for the development set,

<sup>4</sup>Calculated using lexical-diversity Python package.

but similar to the ones generated by GPT-4 for the test set. This disparity is to be expected, given the performance differences between the two models in creative writing tasks.<sup>5</sup>

It is important to note that the adversarial responses from DailyDialog++ exhibit greater diversity compared to those from GENRESCO. This is because the tasks are slightly different: in DailyDialog++, annotators were asked to generate new irrelevant responses by incorporating certain words from the context directly or indirectly into their responses. This stands in contrast to our approach, which prioritises coherence while preserving relevance. As such, the introduction of diverse words into the response is constrained by the fact relevance must be upheld.

### 3 ECOH

This section presents ECOH, our proposed family of response coherence evaluators. We initially present the method of formulating the task of coherence evaluation as explainable QA (Question Answering) (§3.1). Then, we describe in detail how our evaluator is trained (§3.2) and evaluate its performance on different settings (§3.3).

#### 3.1 Problem Formulation

Turn-level dialogue coherence evaluation consists of the assessment of a response hypothesis  $h$  given a dialogue history (frequently denoted as context)  $c$  of varying amount of turns, and optionally one or more references  $r$  and/or external knowledge  $k$ . The goal is to learn a scoring function that assigns a score  $f(c, k, r, h) \rightarrow s$  for each individual quality aspect. This scoring function is compared against human judgements, which annotate the same context-response pairs. These responses are evaluated by humans using, for instance, a binary (0, 1) judgement or a [1, 5] Likert scale, where the lowest value means lowest quality and highest value maximum quality.

In our work, we consider Coherence as being a binary quality aspect. Despite being frequently annotated in the literature on a Likert Scale, what can be considered a response that is neither coherent or incoherent is mostly left to the interpretation of the annotator. Given that we are leveraging an LLM for generation, we find it unfeasible to generate a balanced dataset that contains intermediate levels

<sup>5</sup>It is important to point out, that a higher temperature value would likely result in higher diversity, with a possible trade off in performance.

of coherence. Instead, we generate a positive and a negative response in terms of coherence and label it accordingly. This contrastive sampling strategy for coherent responses is also followed in most metric development work for Relevance or Sensibleness, where models are typically trained using self-supervised learning strategies that sample negative responses by random selection (Mehri and Eskenazi, 2020b; Yeh et al., 2021; Mendonça et al., 2023). Lacking any external knowledge with respect to each dialogue, we then further simplify the reference-free evaluation of coherence as a Question Answering (QA) task ( $f(c, h) \rightarrow s \in (0, 1)$ ), with model responses being either coherent ("Yes") or incoherent ("No").

#### 3.2 Experimental Setup

**Model Specification** We employ the QWEN1.5-CHAT family of LLMs (Bai et al., 2023) for our models. QWEN1.5 contains LLMs of various sizes, ranging from 0.5B up to 72B and support all the languages of XDailyDialog. We limit our fine-tuning experiments up to 4B due to the tradeoff between performances and compute. We feed the dialogue context to the model and ask it to provide a "Yes"/"No" answer to the question "*Given the context, is the response Coherent?*". The model is trained to also output a succinct explanation to the answer. We opted with asking for the explanation first, before answering the question, in order to leverage the autoregressive nature of the model. In theory, this should guarantee that final answer be informed by the explanation.<sup>6</sup> Additional training details are available in B.

**Baselines** We contrastively compare our proposed approach against several models. We begin by including models trained using random negative responses from DailyDialog: a ROBERTA-LARGE model (Liu et al., 2019) (which we train ourselves – see Appendix B for details); and UNIEVAL (Zhong et al., 2022) (which uses T5 as base model). Since these models output a probability score, we assume the model outputs the positive class when the  $p > 0.5$ . Additionally, we conduct zero and one shot (with English and language specific examples) inference using QWEN1.5-CHAT to determine if finetuning on GENRESCO adds improvements to the performance of the base model. We also com-

<sup>6</sup>Chiang and Lee (2023) has shown that dialogue evaluation performance is not *always* better when requesting the explanation first. We leave this analysis for future work.

pare against GPT-3.5-Turbo (Ouyang et al., 2022), the teacher model which was used to generate the development set of GENRESCOH, and which is weaker than our expert (GPT-4).

### 3.3 Main Results

Since the coherence labels are binary, we report detection results using F1-score and Point Biserial Correlation. Additionally, we compute the BLEU-4 score of the generated short explanation using the GPT-4 explanation as a reference. Since BLEU compares overlap in tokens instead of comparing meaning, we also employ GPT-4 as a drop-in replacement for human annotators, and ask it to assess the explanations of 200 random responses from the models that output an explanation.

Model	$\rho_{pb}$	F1	BLEU	GPT-4
1 (always positive)	NaN	.333	-	-
NSP-ROBERTA	.1651	.430	-	-
UNIEVAL	<u>.3272</u>	<u>.500</u>	-	-
QWEN1.5-CHAT				
0.5B	.2226	.600	3.80	1.84 $\pm$ 1.12
1.8B	.5212	.740	2.58	2.39 $\pm$ 1.29
4B	.5850	.783	<u>8.16</u>	3.18 $\pm$ 1.60
7B	<u>.7918</u>	<u>.890</u>	4.63	<u>3.95<math>\pm</math>1.48</u>
GPT-3.5-Turbo	.8256	.910	5.25	<b>4.55<math>\pm</math>1.08</b>
ECOH-EN				
0.5B	.7756	.878	16.02	3.80 $\pm$ 1.43
1.8B	.8242	.908	17.30	4.13 $\pm$ 1.29
4B	.9185	.960	17.92	<u>4.45<math>\pm</math>0.96</u>
ECOH-ML				
0.5B	.8882	.945	17.00	3.99 $\pm$ 1.36
1.8B	.9019	.953	17.28	4.24 $\pm$ 1.28
4B	<b><u>.9491</u></b>	<b><u>.975</u></b>	<b><u>18.05</u></b>	4.29 $\pm$ 1.12

Table 3: Reported results on GENRESCOH-DD-TEST, averaged across all languages.  $\rho_{pb}$  denotes Point Biserial Correlation. ECOH-EN and ECOH-ML denote the finetuned models using English data and all multilingual data, respectively. All correlation results are  $p < 0.05$ . **Bold** denotes best overall model, underline best model of the group.

We collate our main results in Table 3. Due to space limitations, we only report 1-shot performance with a language specific example for QWEN1.5-CHAT and the results correspond to the average of the languages. Additional results, including Zero shot and individual language performance, are available in Appendix C.

**GPT-3.5 performance with 4B parameters** Our main observation is that, although being one of our smallest models, ECOH-0.5B-ML outperforms the

predictive performance of the teacher model (reported in F1), and the explanations of QWEN1.5-CHAT-7B. Furthermore, ECOH-4B-EN has similar explanation quality to that of GPT-3.5-Turbo. As expected, training models using random response selection (NSP-ROBERTA-L and UNIEVAL) is not sufficient for accurately detecting more advanced coherence issues. In fact, these models’ performance sit between QWEN1.5-0.5B-CHAT (.600 F1) and the naive single output model (.333 F1).

**Model size and Multilingual finetuning** Since our smallest model already achieves strong results (.945 F1 score), increasing the model size results in only a small relative improvement of 3% in performance. However, we do observe larger performance improvements with multilingual finetuning. For instance, for ECOH-0.5B, we observe an improvement of over 7% (.878 to .945). This indicates, as expected, that including multilingual data during finetuning improves results for the various covered languages.

**Explanations** We also note that our finetuned models have much higher BLEU and GPT-4 scores than the base models. The obtained BLEU scores are to be expected, given that ECOH is finetuned with explanation data stemming from the same prompt, which is a biased observation from the response generator. This is supported by the teacher model’s performance, achieving the highest GPT-4 assessment, despite having low BLEU. In any case, by validating the responses of the ECOH models with GPT-4, we see that the explanations are on average of higher quality than the ones generated by even the largest base model (QWEN1.5-CHAT) that we studied.

### 3.4 Generalization to unseen dialogue datasets and languages

In order to evaluate our models’ capabilities on unseen dialogue datasets, we evaluate our models on XPersona, which was not seen during finetuning. We only select the best baselines (as reported in Table 3) for this analysis. Additionally, our XPersona subset contains French, which is not present in XDailyDialog, so in addition to the average performance across all languages, we present the results for French separately. For fair comparison, we utilise the English example when evaluating the performance of QWEN1.5-CHAT in French.

Model	$\rho_{pb}$	F1	BLEU	GPT-4
QWEN1.5-CHAT-7B				
FR	.4608	.660	2.97	3.20 $\pm$ 1.58
ML	.6125	.778	3.43	3.75 $\pm$ 1.51
GPT-3.5-TURBO				
FR	.7205	.860	5.04	4.32 $\pm$ 1.25
ML	.7631	.880	4.94	<b>4.45</b> $\pm$ 1.05
ECOH-ML				
0.5B				
FR	.8089	.910	13.71	3.68 $\pm$ 1.46
ML	.8882	.945	<b>17.00</b>	3.82 $\pm$ 1.38
1.8B				
FR	.7863	.890	14.10	4.40 $\pm$ 1.15
ML	.8472	.920	15.70	4.26 $\pm$ 1.17
4B				
FR	.9270	.960	14.58	4.36 $\pm$ 0.95
ML	<b>.9448</b>	<b>.970</b>	16.33	4.38 $\pm$ 0.96

Table 4: Reported results for GENRESCO-H-PC-TEST (French-FR subset and full-ML set).  $\rho_{pb}$  denotes Point Biserial Correlation. All correlation results are  $p < 0.05$ . **Bold** denotes best overall model, underline best model of the group.

Looking at the results in Table 4, we find that the conclusions from DailyDialog also carry over to XPersona. For the unseen language (French-FR), we note a large drop in performance for QWEN1.5-CHAT-7B, when compared to the other languages, which could be explained by the 1-shot example being in English. For our proposed models, we see a larger gap in performance between French and the other languages for the smaller models, whereas for ECOH-4B, the performance for French is well within the range of that observed for other languages. This is also what we observe for GPT-3.5-TURBO. This finding suggests that, given an LLM that natively supports languages for which we have no finetuning data, coherence knowledge can be drawn from languages that were included for finetuning.<sup>7</sup>

### 3.5 Generalization to external annotations

Since the models were trained and evaluated on synthetic data, it is important to check if ECOH performs adequately on external evaluations conducted by human annotators. As such, we also assess ECOH on the FED-turn annotations (Mehri and Eskenazi, 2020a) for "Relevance" and "Over-

<sup>7</sup>It is important to acknowledge that this finding is only likely to extend to languages that follow western normative rules for coherence. An additional interesting experiment would be to test ECOH on a language that does not conform to these rules – however, these are typically low-resource.

all", which is a typically used benchmark for dialogue evaluation. Similar to other works, we calculate the average human annotation ( $[0, 2]$  for Relevance and  $[0, 4]$  for Overall) and report results using Spearman correlation between the human annotation and the score provided by each evaluator. For the LLMs, we keep the binary formulation for coherence (score is either 0 or 1). For the coherence explanation evaluation, lacking a reference, we again use GPT-4 as an explanation evaluator but without a reference response, and evaluate all responses. In order to gauge evaluation performance, we also calculate correlations with GPT-4 as a response evaluator. We refrain from providing GPT-4 explanation scores due to potential self-evaluation bias.

Model	Relevance $r$	Overall $r$	GPT-4
NSP-ROBERTA	.2530	<u>.2543</u>	-
UNIEVAL	<u>.2532</u>	.2521	-
QWEN1.5-CHAT			
0.5B	<i>.0451</i>	<i>.0340</i>	2.35 $\pm$ 1.42
1.8B	.2693	.2228	2.91 $\pm$ 1.52
4B	.1613	.1189	3.30 $\pm$ 1.67
7B	<u>.3279</u>	<u>.2998</u>	<b>3.74</b> $\pm$ 1.54
GPT-3.5-Turbo	.4025	.3636	3.54 $\pm$ 1.66
GPT-4	<b>.5108</b>	<b>.5320</b>	-
ECOH			
0.5B			
EN	.2247	.1548	3.17 $\pm$ 1.77
ML	.1670	.1294	3.17 $\pm$ 1.77
1.8B			
EN	<u>.2941</u>	.2408	3.38 $\pm$ 1.77
ML	.2581	.1801	<u>3.50</u> $\pm$ 1.68
4B			
EN	.2445	.2326	3.17 $\pm$ 1.82
ML	.2685	<u>.2515</u>	3.37 $\pm$ 1.81

Table 5: Reported results for FED-Turn. Performance is calculated using Pearson Correlation ( $r$ ). All results are  $p < 0.05$  unless *italicised*. **Bold** denotes best overall model, underline best model of the group.

From Table 5, we draw several conclusions. Firstly, when looking at the correlation metric, we see that the performance gap between random response-trained models and ECOH is much smaller. This is mainly due to the older chatbots models used for FED-turn – **Meena** (Adiwardana et al., 2020) and **Mitsuku**<sup>8</sup> – being more likely to output irrelevant and non-specific responses that mimic random response selection.

Secondly, we note that our finetuning is still useful for detecting coherence issues on FED, since,

<sup>8</sup>Mitsuku blogpost

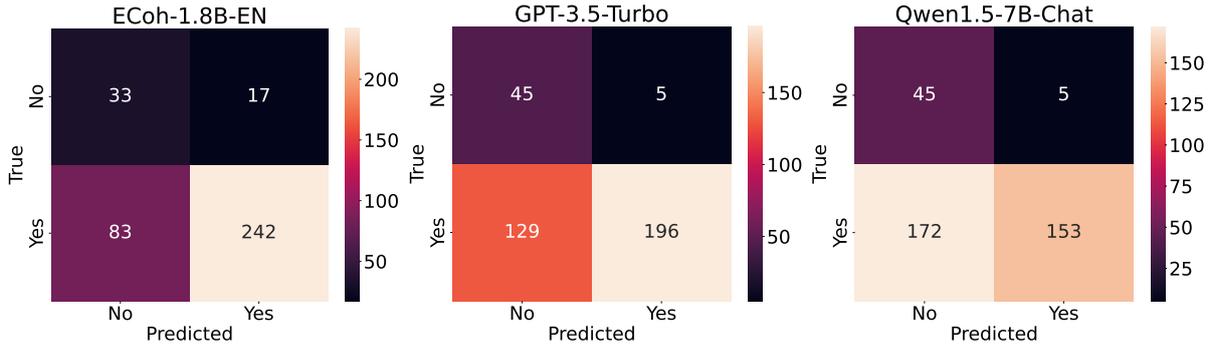


Figure 1: Confusion matrices for the best models of each family (ECOH-1.8B-EN, GPT-3.5-Turbo and QWEN1.5-7B-CHAT) on FED-turn.

overall, ECOH outperforms the corresponding parent model (e.g., ECOH-4B vs QWEN1.5-4B) on Relevance. However, our multilingual models underperform against the English-specific finetuning (with the exception of the 4B model) which could be explained by FED being exclusively in English.

Finally, despite GPT-4 not being a perfect evaluator (low correlation in FED), we assume the scores we obtain for FED remain comparable to previous experiments. With that in mind, we find that the explanation quality is overall lower for ECOH and the teacher model, GPT-3.5-Turbo when compared to GENRESCO (for instance, we report an absolute drop of 1.01 for GPT-3.5-Turbo). In contrast, the QWEN1.5-CHAT models’ explanation quality remains almost unchanged when compared to GENRESCO. As a result, ECOH models achieve less impressive results when compared to their parent models. When comparing models of the same size up until 4B, we do see some improvement in explanation quality. However, our 4B model fails to outperform the 7B model, both in terms of correlation and explanation quality. For the correlation, we believe this is due to the parent model (QWEN1.5-CHAT-4B) having low predictive performance – if we look at the 1.8B models, they yield better correlations. For the explanation quality, we note that the teacher model, GPT-3.5, has also lower results than expected.

However, it is important to acknowledge that these models are tailored towards the binary classification of coherence. As such, we also present results for FED-turn with a binary mapping. In this case, we consider a response to be relevant when the majority of the annotators rate the response as fluent. We present the confusion matrices for the best models of each family in Figure 1. Here, we note that ECOH-1.8B-EN is more likely to say a

response is coherent, incurring slightly more False Positives than the other models (17 against 5), but also lower False Negatives (83 against >129).

### 3.6 Error analysis

In order to determine limitations and weaknesses with using ECOH, we conducted a thorough analysis of all errors in GENRESCO-PC and FED. We summarise significant findings below.

<b>Context:</b>		
A: Hi!		
B: hi		
A: what are you up to?		
B: code refactoring. you?		
A: me? just chilling out at work. what is code refactoring?		
<b>Response. B:</b> good question. I don’t even know what I am doing		
Annotations:	<b>Relevant:</b> 1.6	<b>Overall:</b> 3.0
<b>Evaluation:</b> "The response contradicts B’s previous statement about code refactoring. The answer is No."		

Table 6: Example prediction by ECOH-1.8B-ML on FED-Turn. Here, both speakers are human.

**Humanlike Responses** One of the issues we found when analysing incorrect predictions is that the model sometimes rates human responses as incoherent when they are not "helpful" or do not directly respond to the context, which is not the expected behaviour with chatbot responses. For instance, Table 6 presents a cherry-picked example where ECOH incorrectly rates a human response from FED-Turn as incoherent. We believe such issues could be mitigated by including some portion of coherent human responses from DailyDialog++ during training. However, we also found similar issues with GPT-3.5-Turbo and QWEN, indicating this type of error is transversal to many current LLM-based evaluators.

---

**Context:**

**A:** Che grande giornata! *What a great day!*

**B:** Sì. Che ne dici di uscire e Dioerti il sole sull'erba? *Yes. How about we go out and enjoy the sun on the grass?*

**A:** Grande, andiamo! *Great, let's go!*

**B:** Ehi, tesoro, credo di aver avuto un colpo di calore per essere stato tutto il giorno al sole. Sono così rilassato. Mi sembra di essere in un altro mondo. *Hey, honey, I think I had heatstroke from being in the sun all day. I'm so relaxed. I feel like I'm in another world.*

---

**Response. A:** Oh no, forse dovremmo trovare un posto all'ombra o rientrare in casa. Come ti senti adesso? *Oh no, maybe we should find a shady spot or go indoors. How do you feel now?*

---

**Reference:** *"The response is considerate and logical, offering a solution to avoid further sun exposure due to the heatstroke B mentioned, and it also shows concern for B's current well-being. The answer is Yes."*

---

**Evaluation:** *"The response contradicts B's statement about feeling relaxed and suggests going back indoors. The answer is No."*

---

Table 7: Example prediction by ECOH-1.8B-ML on GENRESCOH-DD-TEST. Translation provided in *italic*.

**Commonsense Reasoning** Another issue we found recurrent, especially with the smaller models, is the limited nature of their commonsense reasoning. We hypothesise that this is a capability that smaller models struggle with, and this is reflected in their evaluation capabilities. We see an example of this is Table 7, where the model fails to understand that sun exposure is mitigated by going back indoors.

## 4 Related Work

### 4.1 Metrics for Dialogue Evaluation

Statistic-based metrics such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and METEOR (Banerjee and Lavie, 2005), are a popular choice for dialogue evaluation because they are simple to calculate and lightweight. However, since they assume valid responses have significant word-overlap with the ground truth, their correlations with human judgements annotations are very low (Liu et al., 2016) due to the one-to-many nature of dialogues. Additionally, they cannot be used to evaluate models whenever a gold-response is not available.

Consequently, learned metrics were proposed. The typical approach was to finetune pretrained encoder models using positive and negative samples targeting different quality aspects such as fluency and relevance (Mehri and Eskenazi, 2020b; Phy et al., 2020; Sai et al., 2020; Mendonca et al., 2022; Zhao et al., 2020). Other approaches used graph

representations to model dialogue interactions explicitly (Huang et al., 2020; Zhang et al., 2021).

With the introduction of LLMs in a wide range of NLP tasks, most recent approaches leverage them for dialogue evaluation. G-EVAL (Liu et al., 2023a) uses GPT-3.5-Turbo and GPT-4 for the evaluation of generation models using a "Chain-Of-Thoughts" step and a scoring function based on return token probabilities. LLM-EVAL (Lin and Chen, 2023) is a single-prompt-based evaluation method that leverages a unified evaluation schema to cover multiple dimensions of conversation quality in a forward pass. DIALEVALML (Mendonça et al., 2023) combines encoder-based models and direct prompting and score extraction from GPT-3.5-Turbo. XDIAL-EVAL (Zhang et al., 2023) probes the evaluation capabilities of several open source LLMs against GPT-3.5-Turbo (Ouyang et al., 2022), and also finetunes them with dialogue data. To the best of our knowledge, this is the first work that conducts supervised learning of LLMs for the task of dialogue evaluation.

### 4.2 Dataset Generation

There are several studies that propose augmentation and synthetic generation approaches to scale dataset sizes that target commonsense reasoning (Bhagavatula et al., 2023; Wang et al., 2023), summarisation (Jung et al., 2024), and dialogues (Chen et al., 2023; Kim et al., 2023) for training purposes.

For dialogue evaluation in particular, most metrics are finetuned using self-supervised data (Mehri and Eskenazi, 2020b; Phy et al., 2020; Yeh et al., 2021; Mendonca et al., 2023). The most widely used approach is to select positive samples consisting of the ground truth response, and negative responses from randomly drawn dialogues. Ghazarian et al. (2022) relies on Abstract Meaning Representation (AMR) to apply semantic-level manipulations to existing responses. Our work, in comparison, leverages a strong LLM to generate new incoherent responses at scale.

## 5 Conclusions

This paper presents GENRESCOH, a large scale collection of positive and negative responses and corresponding explanations covering several languages. GENRESCOH is generated from XDaily-Dialog and XPersona using state-of-the-art LLMs, which better matches the responses seen by contemporary chatbots. With this dataset, we train a family

of evaluators we call ECOH. Our smallest model (0.5B) is able to achieve similar performance to that of the teacher model (GPT-3.5-Turbo), despite being much smaller.

Despite this good performance, we note some limitations when using ECOH, especially when evaluating human responses and/or responses that require more robust commonsense reasoning. Whilst we argue that including more data that targets commonsense and human responses, or even training a larger model could mitigate these issues, since we include an explanation in the predictions, one could still use our evaluators for an initial evaluation screening and escalate to a human evaluator if necessary.

## 6 Limitations

**Reduced Language Selection** Our work is only evaluated in English, German, Italian, French and Chinese. This limitation stems in part from the upstream dialogue dataset (XDailyDialog) only covering 4 high resource languages. Whilst XPersona does contain additional languages, we were limited to only including French as unseen language due to annotator and resource limitations.

**Generation** Generating synthetic data from LLMs might surface or even amplify harmful content within these models. In particular, the choice of a single LLM to generate the responses may induce distribution biases. We identify in Section 2 the reduced lexical diversity of generated responses from GPT-3.5-Turbo when compared to humans. Furthermore, our limited analysis shows that our model sometimes struggles with rating human responses. As such, the generated negative samples may also not accurately represent all coherence issues LLM-based generators typically exhibit. Future investigation may look into producing a systematic quality analysis of a more diverse pool of LLMs, which could inform more faithful generation of negative responses.

**FED as a turn level coherence benchmark** For most dialogue evaluation benchmarks, coherence annotations are conducted at the dialogue level and do not pinpoint the exact response that triggers incoherence (Yeh et al., 2021). As such, we opted with benchmarking ECOH on FED-turn relevance annotations, which is a typically used benchmark for dialogue evaluation. Despite relevance and coherence being different quality aspects, we note

that a) all irrelevant responses lack, by definition, coherence; b) we found that the vast majority of relevant responses on FED are also coherent. Nevertheless, we acknowledge the limitations of using FED-turn as a turn level coherence evaluation benchmark, namely due to its lack of relevant but incoherent responses.

## 7 Ethical Considerations

**Culture-specific conversational norms** We acknowledge that the definition of dialogue quality is a diverse, culturally informed concept. We attempt to reduce the English-centric bias in the generation by leaving the LLM to generate without English reference constraints. However it is possible the generation still conforms to English definitions of coherence given its pretraining and instruction tuning data is more than likely over represented by English text. Furthermore, the examples provided in the prompt, and the dialogues themselves, despite being validated by expert linguists, are still based on English dialogues. As such, users of our model should take extra care when evaluating responses in languages that are known to deviate substantially from English-centric notions of coherence.

**Annotations** The post-editing of the prompts and the manual validation of GPT-4 generations was partially conducted by volunteer annotators, and paid workers that have a fair wage according to their location.

## Acknowledgments

This research was supported by the Portuguese Recovery and Resilience Plan through project C645008882-00000055 (Responsible.AI) and by national funds through *Fundação para a Ciência e a Tecnologia* (FCT) with references PRT/BD/152198/2021 and DOI: 10.54499/UIDB/50021/2020.

## References

- Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. 2020. [Towards a human-like open-domain chatbot](#). *CoRR*, abs/2001.09977.
- 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*.
- 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.
- Chandra Bhagavatula, Jena D. Hwang, Doug Downey, Ronan Le Bras, Ximing Lu, Lianhui Qin, Keisuke Sakaguchi, Swabha Swayamdipta, Peter West, and Yejin Choi. 2023. **I2D2: Inductive knowledge distillation with NeuroLogic and self-imitation**. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9614–9630, Toronto, Canada. Association for Computational Linguistics.
- Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Seokhwan Kim, Andy Rosenbaum, Yang Liu, Zhou Yu, and Dilek Hakkani-Tur. 2023. **PLACES: Prompting language models for social conversation synthesis**. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 844–868, Dubrovnik, Croatia. Association for Computational Linguistics.
- 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.
- Sarik Ghazarian, Nuan Wen, Aram Galstyan, and Nanyun Peng. 2022. **DEAM: Dialogue coherence evaluation using AMR-based semantic manipulations**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 771–785, Dublin, Ireland. Association for Computational Linguistics.
- Edward T. Hall. 1959. *The silent language*. Doubleday, Garden City, N. Y.
- 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**.
- Lishan Huang, Zheng Ye, Jinghui Qin, Liang Lin, and Xiaodan Liang. 2020. **GRADE: Automatic graph-enhanced coherence metric for evaluating open-domain dialogue systems**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9230–9240, Online. Association for Computational Linguistics.
- Jessica Huynh, Cathy Jiao, Prakhar Gupta, Shikib Mehri, Payal Bajaj, Vishrav Chaudhary, and Maxine Eskenazi. 2023. **Understanding the effectiveness of very large language models on dialog evaluation**.
- Jaehun Jung, Peter West, Liwei Jiang, Faeze Brahman, Ximing Lu, Jillian Fisher, Taylor Sorensen, and Yejin Choi. 2024. **Impossible distillation: from low-quality model to high-quality dataset & model for summarization and paraphrasing**.
- Ehsan Kamaloo, Nouha Dziri, Charles Clarke, and Davood Rafiei. 2023. **Evaluating open-domain question answering in the era of large language models**. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5591–5606, Toronto, Canada. Association for Computational Linguistics.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2023. **SODA: Million-scale dialogue distillation with social commonsense contextualization**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12930–12949, Singapore. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. **Adam: A method for stochastic optimization**. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Yen-Ting Lin and Yun-Nung Chen. 2023. **LLM-eval: Unified multi-dimensional automatic evaluation for open-domain conversations with large language models**. In *Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)*, pages 47–58, Toronto, Canada. Association for Computational Linguistics.
- Zhaojiang Lin, Zihan Liu, Genta Indra Winata, Samuel Cahyawijaya, Andrea Madotto, Yejin Bang, Etsuko Ishii, and Pascale Fung. 2021. Xpersona: Evaluating multilingual personalized chatbot. In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*, pages 102–112.
- Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. **How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation**. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2122–2132, Austin, Texas. Association for Computational Linguistics.

- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023a. [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.
- 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](#).
- Zeming Liu, Ping Nie, Jie Cai, Haifeng Wang, Zhengyu Niu, Peng Zhang, Mrinmaya Sachan, and Kaiping Peng. 2023b. [XDailyDialog: A multilingual parallel dialogue corpus](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12240–12253, Toronto, Canada. Association for Computational Linguistics.
- Philip M McCarthy. 2005. *An assessment of the range and usefulness of lexical diversity measures and the potential of the measure of textual, lexical diversity (MTLD)*. Ph.D. thesis, The University of Memphis.
- Shikib Mehri and Maxine Eskenazi. 2020a. [Unsupervised evaluation of interactive dialog with DialoGPT](#). In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 225–235, 1st virtual meeting. Association for Computational Linguistics.
- Shikib Mehri and Maxine Eskenazi. 2020b. [USR: An unsupervised and reference free evaluation metric for dialog generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 681–707, Online. Association for Computational Linguistics.
- John Mendonca, Alon Lavie, and Isabel Trancoso. 2022. [QualityAdapt: an automatic dialogue quality estimation framework](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 83–90, Edinburgh, UK. Association for Computational Linguistics.
- John Mendonca, Alon Lavie, and Isabel Trancoso. 2023. [Towards multilingual automatic open-domain dialogue evaluation](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 130–141, Prague, Czechia. Association for Computational Linguistics.
- John Mendonça, Patrícia Pereira, Helena Moniz, Joao Paulo Carvalho, Alon Lavie, and Isabel Trancoso. 2023. [Simple LLM prompting is state-of-the-art for robust and multilingual dialogue evaluation](#). In *Proceedings of The Eleventh Dialog System Technology Challenge*, pages 133–143, Prague, Czech Republic. Association for Computational Linguistics.
- OpenAI. 2024. [Gpt-4 technical report](#).
- 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 of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Vitou Phy, Yang Zhao, and Akiko Aizawa. 2020. [Deconstruct to reconstruct a configurable evaluation metric for open-domain dialogue systems](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4164–4178, Barcelona, Spain (Online). International Committee on 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.
- Ananya B. Sai, Akash Kumar Mohankumar, Siddhartha Arora, and Mitesh M. Khapra. 2020. [Improving dialog evaluation with a multi-reference adversarial dataset and large scale pretraining](#). *Transactions of the Association for Computational Linguistics*, 8:810–827.
- Jaromir Savelka, Arav Agarwal, Christopher Bogart, Yifan Song, and Majd Sakr. 2023. [Can generative pre-trained transformers \(gpt\) pass assessments in higher education programming courses?](#) In *Proceedings of the 2023 Conference on Innovation and Technology in Computer Science Education V. 1, ITiCSE 2023*, page 117–123, New York, NY, USA. Association for Computing Machinery.
- Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023. [SCOTT: Self-consistent chain-of-thought distillation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5546–5558, Toronto, Canada. Association for Computational Linguistics.
- Peter West, Chandra Bhagavatula, Jack Hessel, Jena Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2022. [Symbolic knowledge distillation: from general language models to commonsense models](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4602–4625, Seattle, United States. Association for Computational Linguistics.

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.

Chen Zhang, Yiming Chen, Luis Fernando D’Haro, Yan Zhang, Thomas Friedrichs, Grandee Lee, and Haizhou Li. 2021. [DynaEval: Unifying turn and dialogue level evaluation](#). 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 5676–5689, Online. Association for Computational Linguistics.

Chen Zhang, Luis D’Haro, Chengguang Tang, Ke Shi, Guohua Tang, and Haizhou Li. 2023. [xDial-eval: A multilingual open-domain dialogue evaluation benchmark](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5579–5601, Singapore. Association for Computational Linguistics.

Chen Zhang, Luis Fernando D’Haro, Yiming Chen, Malu Zhang, and Haizhou Li. 2024. A comprehensive analysis of the effectiveness of large language models as automatic dialogue evaluators. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.

Tianyu Zhao, Divesh Lala, and Tatsuya Kawahara. 2020. Designing precise and robust dialogue response evaluators. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 26–33, Online. Association for Computational Linguistics.

Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. 2022. [Towards a unified multi-dimensional evaluator for text generation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2023–2038, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

## A Dataset Curation

### A.1 Generation

**Prompt** The prompt, which is shared for the development and test set is presented in Table 8. For each language, we translate the example dialogues

and responses using Google Translate<sup>9</sup> and manually validate the full prompt with the expert linguists, ensuring the explanation is accurate for the translated response.

Given the dialog, generate a good and a bad response. In particular, the bad response should have issues that reduce its quality in terms of coherence, such as contradictions, logical inconsistencies, etc. Output the responses, together with a small explanation of the response using the following json format:

```
{"good_response": "...", "good_explanation": "...",
"bad_response": "...", "bad_explanation": "..."}

```

Examples:

Dialogue: A: Have you figured out where you want to transfer to? B: I can’t think of where to go. A: Where would you like to go to school?

Output: {"good\_response": "B: Well, It is not yet decided, but maybe in the east coast.", "good\_explanation": "The response acknowledges the question and provides a region.", "bad\_response": "B: Do you think that I can get married after school?", "bad\_explanation": "The response does not acknowledge the prior question."}

Dialogue: A: You look so tan and healthy! B: Thanks. I just got back from summer camp A: How was it ? B: Great. I got to try so many things for the first time.

Output: {"good\_response": "A: I wish I could go to summer camp too. I’m so bored at home.", "good\_explanation": "The response acknowledges the positive emotions displayed and contrasts it with their own perspective of summer break.", "bad\_response": "A: Did you eat while you were there? You look frail.", "bad\_explanation": "The response contradicts the earlier statement indicating they were healthy."}

Dialogue:

[Dialogue]

Table 8: Response generation instruction template.

**Cost** In total, our generation using a mixture of GPT-3.5-Turbo and GPT-4 amounted to approximately 300\$ USD, with 2/3 of the budget allocated to the generation of test set responses for both XDailyDialog and XPersona.

### A.2 Additional Examples

Table 9, 10, 11 showcase additional examples resulting from GPT-4 response generation for XDailyDialog. Additionally, an example from XPersona for French is in Table 12. For convenience, we provide MT in *italic*.<sup>10</sup>

<sup>9</sup><https://translate.google.com>

<sup>10</sup>These were not checked nor post-edited.

---

**Context:**

**A:** 你好，苗莉，你要去哪里？ *Hello, Miao Li, where are you going?*

**B:** 你好，我要去商店买些水果。 *Hello, I'm going to the store to buy some fruit.*

---

**Positive Response. A:** 哦，那真好，要买哪些水果？ *Oh, that's great. What kind of fruit to buy?*

*"The response naturally continues the conversation by showing interest in the specific types of fruits B is planning to buy."*

---

**Negative Response. A:** 为什么你现在才告诉我我要去商店？ *Why are you just telling me now that you're going to the store?*

*"The response introduces a sense of confrontation and confusion that wasn't present in the prior conversation. It suggests an expectation that B should have informed A earlier, which is not supported by the dialogue."*

---

Table 9: Example test set annotation for XDailyDialog for Chinese-ZH. MT in *Italic*.

---

**Context:**

**A:** Wo wollen Sie hin? *Where do you want to go?*

**B:** Ich fahre ins Krankenhaus. *I'm going to the hospital.*

**A:** Steigen Sie ein, ich bringe Sie hin. *Get in, I'll take you there.*

**B:** Wie viel wird es kosten? *How much will it cost?*

**A:** Es kostet Sie 10 Dollar. *It will cost you 10 dollars.*

**B:** Übrigens, wie spät ist es? *By the way, what time is it?*

---

**Positive Response. A:** Es ist jetzt halb drei nachmittags. *It's now half past two in the afternoon.*

*"The response directly and appropriately answers B's question about the current time, maintaining the flow of the conversation."*

---

**Negative Response. A:** Oh, ich habe keine Uhr, also weiß ich nie, wie spät es ist. *Oh, I don't have a watch, so I never know what time it is.*

*"This response is inconsistent because as a driver, it would be expected that A has a way to tell the time, especially when their job involves scheduling. It contradicts the setting where a driver should be aware of the time to manage their routes and appointments."*

---

Table 10: Example test set annotation for XDailyDialog for German-DE. MT in *Italic*.

### A.3 Manual validation

For the manual validation step, a single annotator for each language is recruited to validate the response and corresponding explanation. A total of 100 randomly selected examples from the test set (generated by GPT-4) were sampled, per language, for this validation. We consider a response to be appropriate if the annotation is 1 or above (fair). The full guidelines provided to the expert annotators are presented in Figure 2.

---

**Context:**

**A:** Sto cercando il vecchio carillon. *I'm looking for the old music box.*

**B:** Sei arrivato nel posto giusto, ci sono dei requisiti cronologici? *You've come to the right place, are there any chronological requirements?*

**A:** Sarebbe fantastico se ci fosse un 1920 qui. *It would be great if there was a 1920's here.*

**B:** Ne abbiamo ricevuto solo 1 ieri e ora ce ne sono 6 qui. *We only received 1 yesterday and now there are 6 here.*

---

**Positive Response. A:** Davvero? Potrei vederli per scegliere il migliore? *Really? Could I see them to choose the best one?*

*"The response shows interest and eagerness to see the selection which logically follows from the information provided about the availability of the items."*

---

**Negative Response. A:** Oh, non mi interessano i carillon degli anni '20. *Oh, I don't care about 1920s music boxes.*

*"The response is inconsistent as it contradicts the earlier expressed interest in a 1920s carillon, making it illogical and irrelevant to the previous exchange."*

---

Table 11: Example test set annotation for XDailyDialog for Italian-IT. MT in *Italic*.

---

**Context:**

**A:** Salut! êtes-vous déjà monté à cheval? *Hi! have you ever ridden a horse?*

**B:** Oui! Ça fait partie de mon traitement contre le cancer. J'ai commencé un nouveau travail. *Toi? Yes! It's part of my cancer treatment. I started a new job. You?*

**A:** J'ai entendu dire que couper la viande de votre alimentation aide à combattre le cancer. *I heard cutting meat from your diet helps fight cancer.*

---

**Positive Response. B:** C'est intéressant, je devrais peut-être envisager cela aussi. Comment ça se passe pour toi? *That's interesting, maybe I should consider that too. How is it going for you?*

*"The response shows interest in the advice given by A and seeks more personal insight, which maintains the flow of the conversation and shows engagement."*

---

**Negative Response. B:** Oui, et c'est pourquoi je mange plus de viande maintenant pour rester en forme. *Yes, and that's why I eat more meat now to stay in shape.*

*"The response is contradictory because A just mentioned that cutting meat helps fight cancer, yet B responds by saying they eat more meat to stay fit, which directly contradicts the advice given."*

---

Table 12: Example test set annotation for Xpersona in French-FR. MT in *Italic*.

## B Implementation Details

### B.1 NSP-ROBERTA

We use the ROBERTA large encoder model downloaded from HuggingFace<sup>11</sup> for all experiments. We train a regression model on a single RTX A6000 GPU using the following sampling strategy: Given a fixed context from DailyDialog, **positive**

<sup>11</sup>[huggingface.co/roberta-large](https://huggingface.co/roberta-large)

responses are drawn directly from the same dialog; **negative** responses are randomly selected and a token coverage test discards semantically similar sentences. In total, 89,707/38,449 datapoints were obtained after processing.

A token representing the speaker was added for each turn, and a history length of 3 turns was used. We applied a regression head consisting of a 2-layer MLP with a hidden size of 1024 and a hyperbolic tangent function as activation for prediction. All parameters were trained/finetuned using Adam optimizer (Kingma and Ba, 2015), using a learning rate of 3e-6 and were trained for 3 epochs using a batch size of 16. Evaluation was conducted every 1,000 steps. The best performing model on the evaluation set was selected for testing.

## B.2 ECOH

We train the ECOH models on a mixture of A100 80GB and RTX A6000 GPUs (depending on model size). We finetune using Huggingface Transformers and PEFT<sup>12</sup> for a 3 epochs for the English model and 1 epoch for the multilingual model with early stopping. We finetune from the base QWEN1.5-CHAT models (full precision) using LoRA (Hu et al., 2021), with  $r = 8$ ,  $\alpha = 32$  and dropout set to 0.1. Gradient accumulation steps is set to 4 with a learning rate of  $1e - 4$ . Batch size was set to maximize VRAM consumption, ranging from 2 up to 8 per device.

For inference, we follow QWEN1.5-CHAT inference code<sup>13</sup>, which generates responses using sampling with a temperature of 1, repetition penalty of 1.1, and top  $p$  set to 0.8.

## C Additional Results

This appendix presents the individual results for zero shot, 1 shot with english example, 1 shot with target language example and the finetuned ECOH models for each individual language for for GENRESCOH-DD-TEST, sorted by model size – 0.5B (Table 13), 1.8B and 4B (Table 14) and 7B and GPT-3.5-Turbo (Table 15). Table 16 presents the results for GENRESCOH-PC-TEST.

Model	$\rho_{pb}$	F1	BLEU	GPT-4
QWEN1.5-0.5B-CHAT-0SHOT				
EN	.2141	.45	2.15	-
DE	.1382	.39	1.80	-
IT	.1695	.41	1.99	-
ZH	.1977	.44	1.88	-
QWEN1.5-0.5B-CHAT-1SHOT-EN				
EN	.2662	.60	2.69	2.12±1.20
DE	.1967	.55	2.19	-
IT	.2210	.55	2.47	-
ZH	.2361	.59	2.25	-
QWEN1.5-0.5B-CHAT-1SHOT-LANG				
EN	.2662	.60	2.69	2.12±1.20
DE	.1870	.59	4.35	1.40±0.77
IT	.1567	.56	4.46	1.84±1.24
ZH	.2803	.64	3.70	2.08±1.25
ECOH-0.5B-EN				
EN	.8995	.95	19.34	4.24±1.01
DE	.6407	.79	14.42	3.28±1.72
IT	.7035	.84	14.41	3.68±1.57
ZH	.8587	.93	15.89	3.92±1.41
ECOH-0.5B-ML				
EN	.9174	.96	19.34	4.08±1.15
DE	.8749	.94	14.42	4.04±1.30
IT	.8565	.93	14.41	3.48±1.66
ZH	.9038	.95	15.89	4.20±1.32

Table 13: Reported results for GENRESCOH-DD-TEST for the 0.5B models.  $\rho_{pb}$  denotes Point Biserial Correlation. All correlation results are  $p < 0.05$ .

<sup>12</sup>[huggingface.co/docs/peft](https://huggingface.co/docs/peft)

<sup>13</sup>[github.com/QwenLM/Qwen1.5](https://github.com/QwenLM/Qwen1.5)

Thank you for agreeing to annotate!

Our work is focused on generating adversarial dialogue responses. To this end, we asked GPT-4 to generate a positive and negative response. A negative response is characterized as having issues that reduce its quality in terms of coherence, such as contradictions and logical inconsistencies. **Your task is to determine if the GPT-4 generations are correct.**

Please provide a ternary answer (0, 1 or 2) to the following question:

*Are the generated responses correct?*

**Your annotation(0-2):0** **bad** indicates that:

- one or both of the proposed responses are incorrect. That is, the good response is not coherent and/or the bad response is **without a doubt** coherent given the dialogue.
- The generated responses are nonsensical – there are clear fluency errors that make it difficult to understand what is being said.

**Your annotation(0-2):1** **fair** indicates that:

- The explanation fails to adequately explain why the response is positive or negative.
- One or both responses are ambiguous.
- The responses have errors (or sound translationese) that do not affect the semantic understanding of the response.

**Your annotation(0-2):2** **good** indicates that the generated response and corresponding explanations are correct and fluent (ignoring tokenization or special characters).

Figure 2: GPT-4 response validation guidelines.

Model	$\rho_{pb}$	F1	BLEU	GPT-4
<b>QWEN1.5-1.8B-CHAT-0SHOT</b>				
EN	.4765	.67	2.14	-
DE	.2663	.49	1.80	-
IT	.3207	.54	1.85	-
ZH	.4047	.62	1.95	-
<b>QWEN1.5-1.8B-CHAT-1SHOT-EN</b>				
EN	.5473	.74	2.70	2.36 $\pm$ 1.38
DE	.4413	.68	2.50	-
IT	.4430	.68	2.76	-
ZH	.5652	.76	2.46	-
<b>QWEN1.5-1.8B-CHAT-1SHOT-LANG</b>				
EN	.5473	.74	2.70	2.36 $\pm$ 1.38
DE	.4680	.71	2.38	2.08 $\pm$ 1.22
IT	.4536	.72	3.03	2.28 $\pm$ 1.31
ZH	.6160	.79	2.21	2.12 $\pm$ 1.27
<b>ECOH-1.8B-EN</b>				
EN	.9227	.96	20.15	4.62 $\pm$ 0.85
DE	.7432	.86	16.11	4.00 $\pm$ 1.40
IT	.7381	.86	15.57	3.76 $\pm$ 1.59
ZH	.8926	.95	17.35	4.12 $\pm$ 1.33
<b>ECOH-1.8B-ML</b>				
EN	.9327	.97	20.08	4.32 $\pm$ 1.31
DE	.8859	.94	15.92	4.04 $\pm$ 1.50
IT	.8732	.94	15.23	4.04 $\pm$ 1.40
ZH	.9159	.96	17.88	4.56 $\pm$ 0.92
<b>QWEN1.5-4B-CHAT-0SHOT</b>				
EN	.7365	.86	3.57	-
DE	.6501	.82	3.49	-
IT	.6275	.81	3.55	-
ZH	.7138	.85	3.48	-
<b>QWEN1.5-4B-CHAT-1SHOT-EN</b>				
EN	.6163	.79	4.08	3.56 $\pm$ 1.50
DE	.5764	.78	4.06	-
IT	.5728	.78	4.34	-
ZH	.5400	.73	3.80	-
<b>QWEN1.5-4B-CHAT-1SHOT-LANG</b>				
EN	.6163	.79	4.08	3.56 $\pm$ 1.50
DE	.5754	.79	7.72	2.64 $\pm$ 1.73
IT	.5269	.75	13.78	3.52 $\pm$ 1.69
ZH	.6213	.80	7.04	3.60 $\pm$ 1.50
<b>ECOH-4B-EN</b>				
EN	.9464	.97	20.66	4.60 $\pm$ 0.91
DE	.8980	.95	17.17	4.64 $\pm$ 0.77
IT	.8982	.95	16.42	3.92 $\pm$ 1.19
ZH	.9315	.97	17.44	4.62 $\pm$ 0.86
<b>ECOH-4B-ML</b>				
EN	.9631	.98	20.74	4.28 $\pm$ 1.34
DE	.9437	.97	16.93	4.38 $\pm$ 1.25
IT	.9377	.97	15.99	3.88 $\pm$ 1.67
ZH	.9520	.98	18.52	4.34 $\pm$ 1.00

Table 14: Reported results for GENRESCO-H-DD-TEST for the 1.8B and 4B models.  $\rho_{pb}$  denotes Point Biserial Correlation. All correlation results are  $p < 0.05$ .

Model	$\rho_{pb}$	F1	BLEU	GPT-4
<b>QWEN1.5-7B-CHAT-0SHOT</b>				
EN	.7490	.86	4.30	-
DE	.4868	.66	4.74	-
IT	.4302	.61	4.90	-
ZH	.6739	.81	4.70	-
<b>QWEN1.5-7B-CHAT-1SHOT-EN</b>				
EN	.8745	.94	4.62	3.76 $\pm$ 1.63
DE	.7938	.90	4.75	-
IT	.7711	.88	4.85	-
ZH	.8210	.91	4.37	-
<b>QWEN1.5-7B-CHAT-1SHOT-LANG</b>				
EN	.8745	.94	4.62	3.76 $\pm$ 1.63
DE	.7998	.90	4.59	3.64 $\pm$ 0.45
IT	.6722	.81	5.07	3.76 $\pm$ 1.78
ZH	.8208	.91	5.07	4.28 $\pm$ 1.14
<b>GPT-3.5-TURBO-0SHOT</b>				
EN	.8592	.93	4.92	4.58 $\pm$ 1.04
DE	.8218	.91	5.47	4.66 $\pm$ 0.85
IT	.8102	.90	5.44	4.36 $\pm$ 1.41
ZH	.8113	.90	5.18	4.54 $\pm$ 1.23

Table 15: Reported results for GENRESCO-H-DD-TEST for the 7B models and GPT-3.5-Turbo.  $\rho_{pb}$  denotes Point Biserial Correlation. All correlation results are  $p < 0.05$ .

<b>Model</b>	$\rho_{pb}$	<b>F1</b>	<b>BLEU</b>	<b>GPT-4</b>
<b>QWEN1.5-7B-CHAT-1SHOT</b>				
EN	.5787	.75	3.28	4.28 $\pm$ 1.45
FR	.4608	.66	2.97	3.20 $\pm$ 1.58
IT	.6474	.82	3.06	3.60 $\pm$ 1.63
ZH	.7630	.88	4.40	3.92 $\pm$ 1.32
<b>ECOH-0.5B-ML</b>				
EN	.9021	.95	17.71	3.90 $\pm$ 1.32
FR	.8089	.90	13.71	3.68 $\pm$ 1.46
IT	.8661	.93	14.34	3.88 $\pm$ 1.33
ZH	.9260	.96	15.84	3.80 $\pm$ 1.38
<b>ECOH-1.6B-ML</b>				
EN	.9043	.95	18.17	4.44 $\pm$ 1.12
FR	.8634	.93	14.01	4.40 $\pm$ 1.15
IT	.8872	.94	14.65	4.04 $\pm$ 1.13
ZH	.9390	.97	16.35	4.16 $\pm$ 1.28
<b>ECOH-4B-ML</b>				
EN	.9443	.97	18.81	4.40 $\pm$ 1.04
FR	.9270	.96	14.58	4.36 $\pm$ 0.95
IT	.9381	.97	15.15	4.36 $\pm$ 0.91
ZH	.9700	.98	16.79	4.40 $\pm$ 0.96
<b>GPT-3.5-TURBO-0SHOT</b>				
EN	.7767	.89	4.71	4.08 $\pm$ 1.38
FR	.7205	.86	5.04	4.43 $\pm$ 1.24
IT	.8102	.90	5.11	4.64 $\pm$ 0.86
ZH	.7452	.87	3.89	4.76 $\pm$ 0.72

Table 16: Reported results for GENRESCOH-PC-TEST.  $\rho_{pb}$  denotes Point Biserial Correlation. All correlation results are  $p < 0.05$ .

# An Investigation Into Explainable Audio Hate Speech Detection

Jinmyeong An <sup>\*1</sup>, Wonjun Lee <sup>\*2</sup>, Yejin Jeon <sup>1</sup>,  
Jungseul Ok <sup>1,2</sup>, Yunsu Kim <sup>3</sup> and Gary Geunbae Lee <sup>1,2</sup>

<sup>1</sup> Graduate School of Artificial Intelligence, POSTECH, Republic of Korea

<sup>2</sup> Department of Computer Science and Engineering, POSTECH, Republic of Korea

<sup>3</sup> aiXplain Inc., Los Gatos, CA, USA

{jinmyeong, lee1jun, jeonyj0612, jungseul.ok, gblee}@postech.ac.kr,  
yunsu.kim@aixplain.com

## Abstract

Research on hate speech has predominantly revolved around detection and interpretation from textual inputs, leaving verbal content largely unexplored. While there has been limited exploration into hate speech detection within verbal acoustic speech inputs, the aspect of interpretability has been overlooked. Therefore, we introduce a new task of explainable audio hate speech detection. Specifically, we aim to identify the precise time intervals, referred to as audio frame-level rationales, which serve as evidence for hate speech classification. Towards this end, we propose two different approaches: cascading and End-to-End (E2E). The cascading approach initially converts audio to transcripts, identifies hate speech within these transcripts, and subsequently locates the corresponding audio time frames. Conversely, the E2E approach processes audio utterances directly, which allows it to pinpoint hate speech within specific time frames. Additionally, due to the lack of explainable audio hate speech datasets that include audio frame-level rationales, we curated a synthetic audio dataset to train our models. We further validated these models on actual human speech utterances and found that the E2E approach outperforms the cascading method in terms of the audio frame Intersection over Union (IoU) metric. Furthermore, we observed that including frame-level rationales significantly enhances hate speech detection accuracy for the E2E approach.

**Disclaimer** The reader may encounter content of an offensive or hateful nature. However, given the nature of the work, this cannot be avoided.

## 1 Introduction

Online platforms such as YouTube, Dailymotion, and TikTok have undoubtedly experienced a notable surge in popularity over the years. While this

has led to an increased dependence on audio as a primary mode of communication, this phenomenon has also brought the issue of hate speech in audio content to the forefront. YouTube, for instance, has consistently been proactive in removing hateful content since its inception, aligning with its hate speech policy<sup>1</sup>. Nevertheless, it is worth noting that out of a total of 10,501,072 channels removed from the YouTube platform within the period of July to September 2023, 26,130 channels were specifically taken down due to their association with hate speech<sup>2</sup>. These statistics underscore the unequivocal importance and the imperative need for the development of effective methodologies to precisely identify hate speech within verbal expressions.

An important point to note, however, is that most hate speech datasets are exclusively text-based. Consequently, research endeavors pertaining to hate speech detection (Qian et al., 2018; Park and Fung, 2017) as well as investigations into hate speech explainability (Mathew et al., 2021) are confined to textual inputs. In other words, despite the explosive increase of hate speech on audio-based online social platforms, there is a notable absence of research that addresses hate speech in verbal data. A few studies related to auditory hate speech detection have been proposed. For example, Ibañez et al. (2021); Rana and Jha (2022) curated respective audio-visual multi-modal datasets. Yet, to the best of our knowledge, no research addresses *explainable* hate speech detection in the audio domain, that is, the understanding of the rationale behind the model’s decisions.

Therefore, we first introduce the new task of explainable audio hate speech detection, which encompasses two sub-tasks: audio hate speech classification (AHS-CLS) and audio hate speech frame

<sup>1</sup><https://blog.youtube/inside-youtube/the-four-rs-of-responsibility-remove/>

<sup>2</sup><https://transparencyreport.google.com/youtube-policy/removals?hl=en>

\*Equally contributed

detection (AHS-FD). The former involves determining whether an audio utterance is hate speech, while the latter identifies the specific time frames containing hate speech. In addition, since there is a lack of interpretable audio hate speech datasets that include audio frame-level rationales, we curated a dataset called AudioHateXplain, which annotates which part of human and synthetic audio recordings pertains to hate speech. Moreover, we propose cascading and E2E models, which are able to elucidate the underlying reasons for classifying speech as hate or not by identifying relevant audio rationales. The cascading approach first transforms audio into text transcripts, detects hate speech within these transcripts, and then matches the detected hate speech to the corresponding audio time frames. In contrast, the E2E approach directly processes the audio input, enabling it to identify the specific time frames containing hate speech accurately. We validated these models on actual human speech utterances and found that the E2E approach outperforms the cascading method regarding the audio frame Intersection over Union (IoU) metric. This superiority is attributed to the bottlenecks arising from conversion between audio and text, such as Automatic Speech Recognition (ASR) errors, and disarrangement between word tokens and time frames. Furthermore, we observed that including frame-level rationales significantly enhances hate speech detection accuracy for the E2E approach.

## 2 Related Work

### 2.1 Hate Speech Detection

Over the years, there have been considerable efforts towards text-based hate speech research, a domain that has gone through various and separate nomenclatures such as cyber hate, offensive, and online abusive Nobata et al. (2016); Davidson et al. (2017) language detection. In this paper, we define these terms collectively as hate speech. In the initial stages of hate speech detection research, Sbertus (1997) predominantly employed feature-based rules. Similarly, Mahmud et al. (2008) incorporated a set of rules to extract semantic information. More recently, in response to the escalating prevalence of online hate speech, there has been a concerted effort to curate private or publicly accessible hate speech datasets Kwok and Wang (2013); Zampieri et al. (2019). However, training hate speech detection models on such datasets, which feature binary-level hate speech annotation, lacks interpretability.

As a result, it becomes difficult to comprehend the logic behind model decisions. In light of this, Mathew et al. (2021) curated a dataset with word-level annotations (rationales), which deviates from conventional datasets focused solely on increasing sentence-level model classification performance.

It is even more important to note that the predominant focus has been on text-based classification. In other words, more datasets and research for verbal hate speech detection and explanations must be needed. To address this, Ibañez et al. (2021); Rana and Jha (2022) curated audio-visual multimodal datasets, while Ibañez et al. (2021) amassed short-form Filipino videos and compared different classification methods, including Support Vector Machine, logistic regression, and Random Forest. Similarly, Rana and Jha (2022) collected videos from Twitter and YouTube, then implemented a multi-task learning model to better identify hate speech by combining text, visual, and acoustic information.

Yet, to the best of our knowledge, no research addresses *explainable* hate speech detection in the audio domain, that is, understanding the rationale behind the model’s decisions. Our research endeavors extend beyond conventional text-based approaches by expanding into the audio domain. Moreover, we address model explainability in addition to audio hate speech detection.

### 2.2 Audio Classification & Frame Detection

Classifying audio clips into specific categories, such as speech commands (Warden, 2018), urban sound events (Piczak, 2015), and the emotional content of speakers (Busso et al., 2008), has been extensively researched. In addition to classifying entire audio clips, there’s been exploration into classifying audio frames at specific time intervals (e.g., every 10 milliseconds), as seen in speaker diarization studies (Canavan et al., 1997; Fujita et al., 2019). In this study, we combine both approaches to classify entire audio clips as containing hate speech or normal speech, while also pinpointing the exact segments within the audio where hate speech occurs, using a 10-millisecond time grid.

## 3 Dataset Generation

Given the absence of existing explainable audio hate speech datasets, we created a synthetic dataset using a text-to-speech (TTS) model. This section details the methods used to convert text transcripts

	Samples	Avg. Length (sec.)
<b>Train</b>	14,183	6.62
<b>Dev.</b>	1,771	6.60
<b>Test-Synth.</b>	300	8.97
<b>Test-Human</b>	300	10.52

Table 1: Summary of the AudioHateXplain dataset, including audio duration. ‘Test-Synth’ refers to the spoken audio generated by a text-to-speech (TTS) model, while ‘Test-Human’ refers to audio recorded by human speakers.

into spoken utterances and generate audio rationales for explaining audio hate speech.

To delineate the process of audio rationale generation, it is imperative first to comprehend the foundational structure of the original text-based HateXplain (Mathew et al., 2021) dataset. This text-based HateXplain dataset, represented as  $\mathcal{D} = \{(x^{(1)}, W^{(1)}, y^{(1)}), \dots, (x^{(L)}, W^{(L)}, y^{(L)})\}$ , comprises  $L$  samples. Each sample consists of a textual sentence  $x$  paired with its corresponding binary class label  $y \in \{0, 1\}$ , denoting whether the sentence qualifies as hate speech (1) or as normal discourse (0). Moreover, each textual sentence  $x$  is supplemented by a set of word-level annotations  $W$ , which is defined as  $W = \{(w^{(1)}, \delta_w^{(1)}), \dots, (w^{(N)}, \delta_w^{(N)})\}$ . Here,  $N$  signifies the position of a word  $w$  within the sentence, with each word linked to its word-level rationale  $\delta_w \in \{0, 1\}$ . Specifically, a word  $w$  is assigned a rationale of 1 if it contributes to the classification of the sentence as hate speech, and 0 otherwise. These word-level rationales serve as discernible evidence aiding in identifying and classifying hate speech within textual content.

**Text-to-Speech** From the above-mentioned text dataset, we convert each text-based transcript  $x$  into audio samples of sample rate 22050 Hz. This conversion is achieved by leveraging the non-autoregressive FastSpeech (Ren et al., 2021) TTS model in conjunction with the HiFi-GAN (Kong et al., 2020) vocoder. To ensure the coherence of audio samples, we expand abbreviations, remove emojis, and exclude sentences in languages other than English, as well as those that contain semantically vacuous words like the placeholder “<user>”.

**Rationale Labeling** Each TTS-generated audio sample  $a$  is then paired with its binary classification label  $y \in \{0, 1\}$  to identify whether it is hate speech (1) or normal audio (0). Moreover,

it is imperative to provide acoustic rationales to facilitate the explainability of hate speech detection in the audio domain. To accomplish this, we employ a pretrained Montreal Forced Aligner (McAuliffe et al., 2017) to identify the timestamps, i.e., the starting and ending times in milliseconds of each spoken word within a given sentence. Subsequently, we divide audio samples into  $M$  10ms-long audio frames  $f$ , and each audio frame  $f$  is annotated with its audio frame-level rationale  $\delta_f \in \{0, 1\}$ . It is possible to annotate each audio frame with rationales as text-based word-level rationale  $\delta_w$  are previously provided. Our AudioHateXplain dataset can thus be represented as  $\mathcal{D} = \{(a^{(1)}, F^{(1)}, y^{(1)}), \dots, (a^{(L)}, F^{(L)}, y^{(L)})\}$  of  $L$  TTS-generated audio samples. In addition, the set of frame-level audio annotations  $F$  can be reorganized as

$$F = \{(f^{(1)}, \delta_f^{(1)}), \dots, (f^{(M)}, \delta_f^{(M)})\} \quad (1)$$

**Human Recordings** In addition to the synthetic audio dataset, we curated a separate collection of authentic human recordings for evaluation purpose. The transcript used for these recordings were derived from the original text-based HateXplain dataset, but underwent a two-step post-processing procedure. Specifically, among the 1,779 original test samples, we initially use ChatGPT (Appendix A) to select texts that were suitable for spoken format. This process enabled us to sample 695 spoken-form texts from the HateXplain test set. These texts were then manually filtered to ensure that the final utterances adhered to syntactic and lexical choices appropriate for spoken language (Ong, 2002; Biber, 1986). Ultimately, 300 utterances were selected for the test set. We synthesized these samples using TTS models and also recorded them with human participants.

The participant group consisted of 10 individuals, 6 males and 4 female speakers. Each participant read an average of 30 utterances, including hate speech and normal texts. Recordings were conducted in silent environments. For ethical considerations, all participants were fully informed about the nature of the transcripts, which included hateful language. Moreover, all recordings were conducted with the explicit consent of the volunteers for research purposes only.

The statistics of the AudioHateXplain is provided in Table 1.

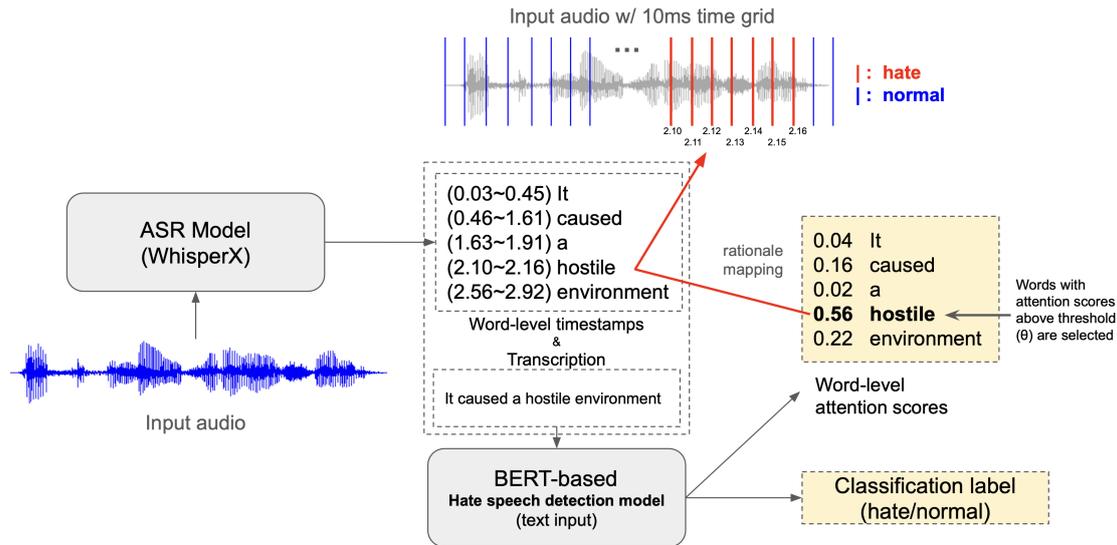


Figure 1: Overview of the cascaded method. Boxes highlighted in yellow indicate model outputs.

## 4 Methodology

In order to classify entire audio clips as hate speech or normal, as well as precisely pinpoint the audio frames associated with hate speech, we introduce two models. The first model uses a cascading framework (Figure 1), which transcribes audio into text, predicts hate speech in the text, and then maps the word-level rationale onto a time grid. The second model (Figure 2) employs an E2E design, directly classifying and predicting hate speech frames from the audio.

### 4.1 Cascading Method

Two essential components comprise the cascading model: an Automatic Speech Recognition (ASR), and a BERT-based hate speech detection model (Figure 1).

Given an audio input, the WhisperX ASR model (Radford et al., 2023; Bain et al., 2023) converts the spoken words into text, while simultaneously generating timestamps for each input word. Following this ASR phase, the transcribed text is passed as input to a finetuned BERT-based hate speech detection model (Mathew et al., 2021). This detection model comprises 12 transformer encoder layers, each containing 768 hidden units and utilizing 12 attention heads. Additionally, a composite loss (Equation 2) is employed during the fine-tuning of the BERT-based hate speech detection model, which consists of two distinct losses:

$$L_{total} = L_{pred} + \lambda L_{att} \quad (2)$$

Classification loss ( $L_{pred}$ ) is derived from the

classification of hate speech within the text. Simultaneously,  $L_{att}$  denotes the loss associated with predicting attention values corresponding to the [CLS] token in the model’s final attention layer. Both losses are computed using cross-entropy. The coefficient  $\lambda$  serves as the weighting factor for  $L_{att}$ , thereby adjusting its influence on the total loss,  $L_{total}$ .

The fine-tuned BERT-based hate speech model outputs the token-level rationales for each word in the input transcribed text. These rationales are produced by leveraging token-level attention scores associated with the [CLS] token in BERT (Devlin et al., 2018), and are transformed into binary format (0 or 1) based on whether they surpass a predefined threshold  $\theta$ . Afterward, a majority voting mechanism consolidates these binary token values into word-level rationales. Specifically, if the majority of token-level rationales for a given word is 1, the word-level rationale is assigned a value of 1; otherwise, it is assigned 0. Finally, these word-level rationales are aligned with audio

### 4.2 End-to-End (E2E) Model

In contrast to the cascaded method, the E2E model presents a direct approach to detect and locate instances of hate speech in audio content, as it eliminates the need to transcribe the audio into text as an intermediary step. By using the wav2vec 2.0 model (Baevski et al., 2020), input audio signals are converted into 1024-dimensional speech representation  $z$  every 25 milliseconds, with a stride of 20 milliseconds. This encoded speech representa-

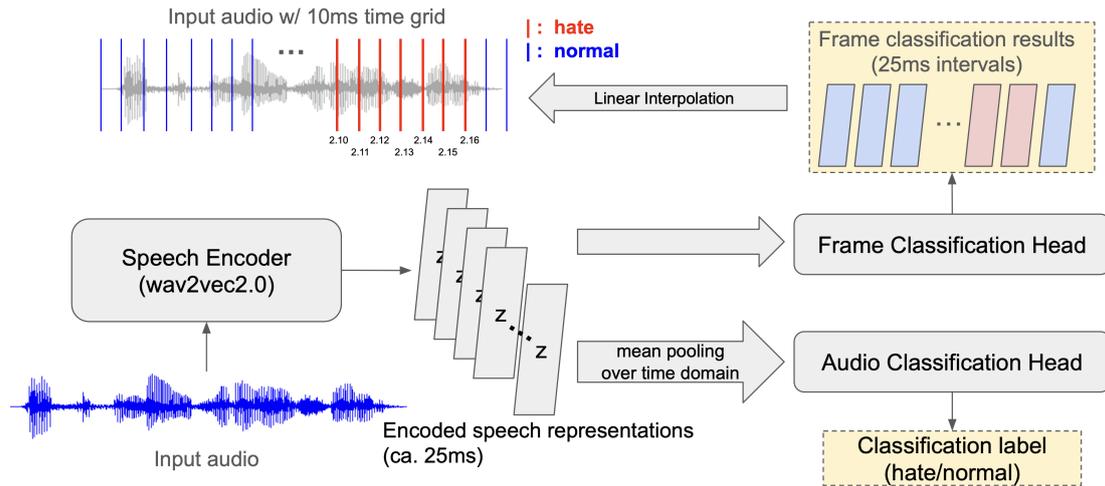


Figure 2: Overview of E2E model. Boxes highlighted in yellow indicate model outputs (AHS-CLS and AHS-FD).

tion  $z$  is then directed to two distinct audio- and frame-level classification heads (Figure 2).

The audio-level classification head is tasked with discerning whether the entire audio sample constitutes hate speech or normal speech. It achieves this through a series of transformations, including a projection layer [1024, 256], mean pooling for temporal feature aggregation, and a linear layer [256, 2] to convert  $z$  into classification logits.

On the other hand, the frame-level detection head is dedicated to identifying individual frames corresponding to hate speech. Comprised of a single linear layer [1024, 2], this head operates directly on individual frame-level features without any feature aggregation, which preserves the granularity necessary for precise frame-level detection.

To effectively optimize both audio-level classification (AHS-CLS) and frame-level detection (AHS-FD) tasks simultaneously, we also employ a multi-task learning approach with the following loss function:

$$L_{total} = \alpha L_{CLS} + (1 - \alpha) L_{FD} \quad (3)$$

The  $L_{CLS}$  and  $L_{FD}$  cross-entropy losses are associated with the AHS-CLS and AHS-FD tasks, respectively. During a hyperparameter search, the value of  $\alpha$  is varied from 0.1 to 0.9 in increments of 0.1 to determine the optimal balance between these two tasks within the multi-task learning framework. It is found that the most effective balance occurs when  $\alpha$  is set to 0.5.

## 5 Experimental Setup

**Cascading Models** We adapted WhisperX (Bain et al., 2023) for ASR and for word-level time stamp-

ing and utilized the BERT model from (Mathew et al., 2021) for hate speech detection within transcribed text. The BERT model<sup>3</sup> was trained using the Adam optimizer (Kingma and Ba, 2014) with a learning rate of  $2e-5$ , a batch size of 64, a threshold  $\theta$  set to the mean value of attention scores, and a coefficient  $\lambda$  of 0.1. We selected the best checkpoint with the highest AHS-CLS F1 after 10 training epochs. The ASR model was not fine-tuned, using Whisper-large-v2<sup>4</sup> as the checkpoint. To enhance the robustness of the cascading model against ASR transcription errors, we trained the BERT-based model with both ASR transcriptions and golden texts. The model trained with ASR transcriptions is called **Cas. (ASR text)**, while the model trained with golden texts is referred to as **Cas. (gold text)**. The BERT model is fine-tuned using either golden transcriptions paired with corresponding word-level rationales or ASR transcriptions with word-level rationales that account for potential inaccuracies.

**E2E Models** For the end-to-end (E2E) model, we utilized wav2vec 2.0<sup>5</sup> as the shared speech encoder. The audio-level classification head and the frame-level detection head were fine-tuned simultaneously or separately (E2E CLS-only and E2E FD-only). Unless specified otherwise, the E2E model was trained with both tasks. We employed the Adam optimizer with a learning rate  $4e-3$  and a batch size of 64, training the model for 50 epochs. The best model was selected based on the highest AHS-CLS F1 score.

<sup>3</sup><https://huggingface.co/google-bert/bert-base-uncased>

<sup>4</sup><https://huggingface.co/openai/whisper-large-v2>

<sup>5</sup><https://huggingface.co/facebook/wav2vec2-large>

Model	Accuracy	F1	Recall	Precision
<i>AudioHateXplain (human recording)</i>				
Cas. (Gold text)	<b>77.00</b>	<b>74.88</b>	<b>74.65</b>	<b>75.17</b>
Cas. (ASR text)	74.66	73.37	74.19	73.02
E2E	71.43	70.84	70.72	71.07
<i>AudioHateXplain (synthetic)</i>				
Cas. (Gold text)	<b>75.00</b>	<b>73.13</b>	73.27	<b>73.01</b>
Cas. (ASR text)	74.00	72.76	<b>73.67</b>	72.43
E2E	71.02	70.51	70.46	70.56

Table 2: Result of Audio Hate Speech Classification (AHS-CLS) on AudioHateXplain test sets.

## 5.1 Evaluation Metrics

To assess the performance of audio hate speech classification (AHS-CLS) and audio hate speech frame detection (AHS-FD), a diverse set of metrics is employed. For AHS-CLS, we employ conventional metrics such as accuracy, precision, recall, and F1 score. In the case of AHS-FD, our evaluation encompasses standard F1 score, as well as frame-level accuracy and Intersection over Union (IoU) metrics, which are used to measure rationales (DeYoung et al., 2019). Given the necessity to assess audio frame-level rationale, we include the 1D IoU metric, commonly used in speaker diarization (Huang et al., 2020), as it allows for the quantitative measure of a model’s accuracy in determining the durations of hate speech within audio frames, and is computed as

$$\text{IoU} = \frac{\text{area}(F_p \cap F_{gt})}{\text{area}(F_p \cup F_{gt})} \quad (4)$$

Here,  $F_p$  and  $F_{gt}$  represent the sets of predicted and ground truth frame-level audio annotations, respectively.  $\text{area}(F_p \cap F_{gt})$  denotes the intersection, i.e., the number of overlapping audio frames with a hate speech rationale ( $\delta_f = 1$ ) between  $F_p$  and  $F_{gt}$ , and  $\text{area}(F_p \cup F_{gt})$  denotes their union, applied over a 10ms time grid.

## 6 Result and Analysis

### 6.1 Audio Hate Speech Classification

The AHS-CLS aims to accurately classify entire audio clips as either hate or normal speech. In Table 2, we report AHS-CLS performance for the cascaded and the E2E model using accuracy, F1 scores, recall, and precision metrics.

Upon assessment using the AudioHateXplain test sets (human recording and synthetic), we observe cascaded models show robust classification results over the E2E model in terms of accuracy

Model	IoU	Frame F1	Frame Recall	Frame Precision
<i>AudioHateXplain (human recording)</i>				
Cas. (Gold text)	14.20	32.96	33.34	32.59
Cas. (ASR text)	15.99	35.63	<b>40.26</b>	31.95
E2E	<b>19.59</b>	<b>37.56</b>	28.03	<b>56.92</b>
<i>AudioHateXplain (synthetic)</i>				
Cas. (Gold text)	19.23	39.76	37.57	42.22
Cas. (ASR text)	18.25	38.91	36.53	41.62
E2E	<b>21.16</b>	<b>43.34</b>	<b>43.69</b>	<b>43.00</b>

Table 3: Result of Audio Hate Speech Frame Detection (AHS-FD) on AudioHateXplain test sets.

(77% vs. 71.43%). Also, we found that a cascaded model trained with golden text rather than ASR transcription shows better classification performance. The performance degrades in Cas. ASR text is likely attributed to overfitting on ASR noise, which is the ASR transcription of the AudioHateXplain dataset.

Notably, all models exhibit slightly higher classification accuracy on human recordings than the synthetic test set. This indicates that models trained on TTS-generated audio can also be used for real human voices.

### 6.2 Audio Hate Speech Frame Detection

The objective of Audio Hate Speech Frame Detection (AHS-FD) is to accurately identify individual audio frames associated with hate speech. Table 3 summarizes the frame detection performance for both cascaded and end-to-end (E2E) models.

Across both datasets and all metrics, the E2E model consistently demonstrates superior performance compared to the cascaded models, except in frame recall for human recordings. In AHS-FD, each frame within a 10ms time grid is labeled as either hate speech or normal. Since there are more normal frames than hate frames, we must consider frame F1 scores to understand the proportion of false negatives and true positives. The E2E model shows more reliable frame F1 scores in all evaluations. Our primary interest lies in detecting hate speech frames rather than normal frames. Therefore, the IoU score is a more reliable metric for this task, as it accounts for both detection and precise localization of hate speech frames. As demonstrated in Table 3, the E2E model consistently outperforms the cascaded models in IoU scores, with differences of up to **5.39%** compared to the cascaded models.

It is important to note that the E2E model shows reliable IoU scores on both the human recording test set and the synthetic test set. In contrast, the cascaded models exhibit a significant degradation

	Human Recording	Synthetic
WER	17.53	8.45
Hate Speech WER	<b>30.31</b>	<b>18.33</b>

Table 4: Comparison between word error rate (WER) for the entire test set, and for words annotated as hate speech.

in IoU for human recordings compared to synthetic voices. This degradation is due to the higher WER in human recordings (17.53% vs. 8.45%), as indicated in Table 4, suggesting that the performance of the cascaded models is highly affected by the accuracy of the ASR model.

### 6.3 Comparative Analysis for AHS-FD

In this section, we attempt to understand the reasons for the differences in AHS-FD performance, which is observed between the cascading and E2E models. We hypothesize that the audio-to-text conversions and text-to-audio alignment within the cascading model framework are what causes severe bottlenecks for audio frame-level detection performance, as indicated by IoU performance. To test this hypothesis, we conducted three different sets of analyses.

First, we analyzed the difference in ASR error between hate and non-hate words. As depicted in Table 4, when using the Whisper-large-v2 model, the word error rate (WER) for the entire test set is **17.5%**, while the WER for words annotated as hate speech is **30%**. In other words, the ASR model shows instability in recognizing audio hate words.

Second, we examined how ASR error affects IoU performance. We initially segmented the audio data into three distinct groups based on varying WER intervals and then evaluated the IoU of each WER interval for both the cascading and E2E models. As depicted in Figure 3, an inverse correlation typically emerges within the cascading model; an increase in ASR errors correlates with a reduction in IoU performance. Conversely, the E2E model consistently exhibits robust performance across different ASR error intervals and outperforms the cascading model across all audio data sets. Since ASR model types can influence WER, we conducted further experiments utilizing various versions of the WhisperX ASR models (i.e., tiny to large-v2). As shown in Figure 4, an increase in ASR errors results in a proportional decline in IoU performance.

Lastly, we examined the effect of audio word-level timestamp errors in text-to-audio conversion.

	Cascading	IoU	Frame F1
w/ GroundTruth Timestamp		<b>17.48</b>	<b>37.93</b>
w/ ASR Timestamp		15.99	35.63

Table 5: Effect of audio word-level timestamp errors on audio explainability performance.

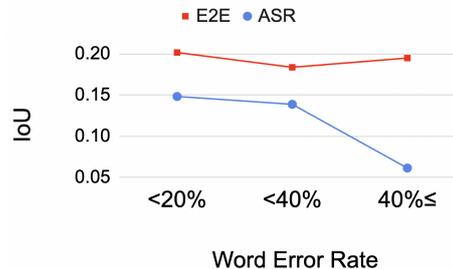


Figure 3: Comparison of IoU scores on human recording test data within three different WER ranges.

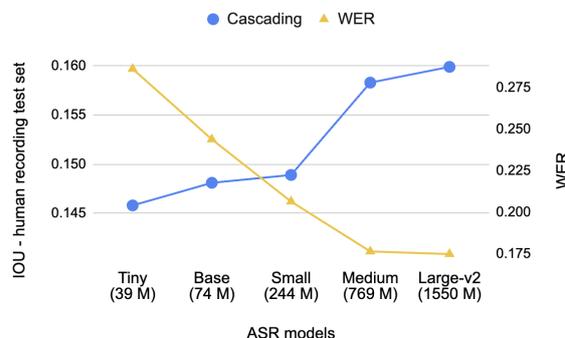


Figure 4: Impact of ASR error for IoU score in cascaded method. The numbers in parentheses represent the total number of parameters in different ASR (Whisper) models.

The timestamping performance for words with a WER of 0 from the WhisperX large-v2 model has 71% IoU score compared with the ground truth timestamp. Additionally, we measured audio frame-level detection performance between the cascading model that utilizes predicted word-level timestamps and the same model that uses ground truth word-level timestamps. As shown in Table 5, there is an IoU decrease of approximately 2% due to the ASR model’s timestamp errors.

### 6.4 Frame Detection Error Analysis

Using actual example data, we examine the time-frame detection capabilities of the E2E and cascading models. Figure 5 visually presents the alignment of ground truth (GT) and predicted audio frame-level rationales in blue and red, respectively. The substantial IoU overlap, depicted in green, between GT and predicted hate speech frames high-

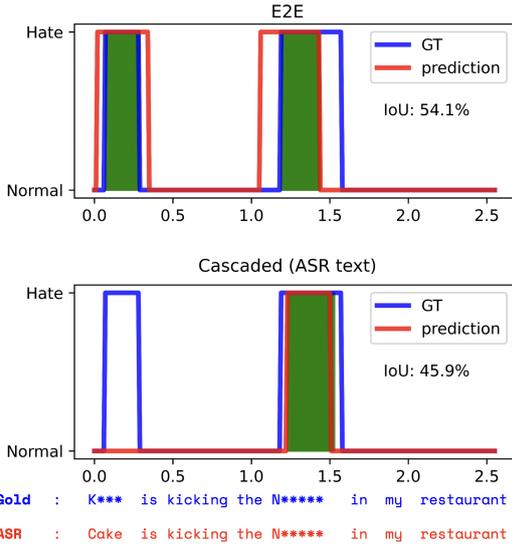


Figure 5: Visualization of audio hate speech frame prediction for E2E and cascading models. Blue letters and graphs indicate the ground truth transcript and rationale, while red letters and graphs show the values predicted by the model. The green part represents the range of the time frame that the model actually predicts.

lights the superior performance of the E2E model in identifying segments of audio containing hate speech. In contrast, the cascaded model exhibits a significant decrease in IoU (8.2%) compared to the E2E model. This decline can be attributed to ASR errors in the transcription process, where the original ethnic slur is inaccurately transcribed as “cake.”

Moreover, in the case of the cascading model, only the timestamp corresponding to each word is known. This means that there is a potential risk where the entire frame-level rationale corresponding to one word is either completely correct or completely incorrect. For example, the cascading model’s prediction for “K\*\*\*” was entirely incorrect, with no partially correct segments. Conversely, in the case of the E2E model, since the audio frame itself is predicted, even if a perfect prediction is not made, the frame-level rationale corresponding to a specific part of the word can still be identified. For example, although the E2E model did not make a perfect prediction for “N\*\*\*\*\*,” it provided a partial correct prediction.

## 6.5 Effect of Multi-task Learning

In order to validate the effectiveness of employing multi-task learning for E2E model (referred Section 4.2), we conduct experiment in Table 6. We found that integration of both classification and frame

Audio Hate Speech Classification (AHS-CLS)		
Model (Loss)	Accuracy	F1
E2E (CLS-only)	75.1	73.7
E2E (CLS+FD)	<b>76.2</b>	<b>75.1</b>
Audio Hate Speech Frame Detection (AHS-FD)		
Model (Loss)	IoU	Frame F1
E2E (FD-only)	31.6	54.07
E2E (CLS+FD)	<b>32.0</b>	<b>56.4</b>

Table 6: Comparisons of Audio Hate Speech Classification (AHS-CLS) and Frame Detection (AHS-FD) performance for E2E models trained for single-task and multi-task settings.

detection learning (CLS+FD) yields better performance compared to models that only employ either classification (CLS-only) or frame detection (FD-only). This enhancement can be attributed to the contextual information that the E2E model gains as it traverses individual hate speech frames within an audio clip to identify frames corresponding to hate speech. Such context augments the model’s proficiency in classifying the entire clip accurately as either hate speech or not, and vice versa.

## 7 Conclusion

In this paper, we introduced the new task of explainable audio hate speech detection, which encompasses two sub-tasks: audio hate speech classification (AHS-CLS) and audio hate speech frame detection (AHS-FD). Furthermore, we introduced E2E and cascading models. These models are capable of not only classifying hate speech directly from verbal speech, but also identifying hate rationales within audio frames. In particular, the proposed E2E model consistently outperforms the cascading model on the AHS-FD task. This superiority is attributed to the bottlenecks arising from conversion between audio and text within the cascading model. This suggests that, for the task of explainable audio hate speech detection, is it more effective to directly process audio inputs. Upon acceptance, we plan to make our dataset and code publicly available to encourage further research for the important topic of explainable audio hate speech detection.

## Limitations

Our AudioHateXplain train split comprises synthetic audio generated through a TTS model, rather than authentic human verbal data. This choice stems from the scarcity of datasets containing real human-recorded audio featuring instances of hate speech, alongside the inherent challenges in curating such recordings. Despite this, our models trained using the synthetic train set demonstrate impressive performance when tested on the human recording test set. We plan to curate a more expansive dataset comprising genuine human recordings as future work. Moreover, this study focuses on English due to the limited resources in other languages. Consequently, our approach does not accommodate the detection of multi-lingual audio hate speech.

## Ethical Considerations

This study on explainable audio hate speech detection involves several ethical considerations. Human recordings were obtained with informed consent, ensuring participants understood the research and potential exposure to offensive content. Sensitive content was handled carefully, with participants fully aware of its nature. The deployment of these models must prevent misuse, such as unjustified censorship, and be rigorously tested for biases to avoid unfair treatment of specific groups.

## Acknowledgments

This work was supported by Institute of Information communications Technology Planning Evaluation (IITP) grant funded by the Korea government(MSIT) (No.2022-0-00223, Development of digital therapeutics to improve communication ability of autism spectrum disorder patients), and by the MSIT(Ministry of Science and ICT), Korea, under the ITRC(Information Technology Research Center) support program(IITP-2024-RS-2024-00437866) supervised by the IITP(Institute for Information Communications Technology Planning Evaluation).

## References

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.

Max Bain, Jaesung Huh, Tengda Han, and Andrew Zisserman. 2023. Whisperx: Time-accurate speech transcription of long-form audio. *INTERSPEECH 2023*.

Douglas Biber. 1986. Spoken and written textual dimensions in english: Resolving the contradictory findings. *Language*, 62:384.

Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N. Chang, Sungbok Lee, and Shrikanth S. Narayanan. 2008. Iemocap: interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42(4):335–359.

Alexandra Canavan, David Graff, and George Zipperlen. 1997. CALLHOME american english speech.

Wallace Chafe and Deborah Tannen. 1987. The relation between written and spoken language. *Annual review of anthropology*, 16(1):383–407.

Thomas Davidson, Dana Warmesley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of the international AAAI conference on web and social media*, volume 11, pages 512–515.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C Wallace. 2019. Eraser: A benchmark to evaluate rationalized nlp models. *arXiv preprint arXiv:1911.03429*.

Yusuke Fujita, Naoyuki Kanda, Shota Horiguchi, Kenji Nagamatsu, and Shinji Watanabe. 2019. End-to-end neural speaker diarization with permutation-free objectives. In *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, volume 2019, pages 4300–4304.

Zili Huang, Shinji Watanabe, Yusuke Fujita, Paola García, Yiwen Shao, Daniel Povey, and Sanjeev Khudanpur. 2020. Speaker diarization with region proposal network. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6514–6518. IEEE.

Michael Ibañez, Ranz Sapinit, Lloyd Antonie Reyes, Mohammed Hussien, Joseph Marvin Imperial, and Ramon Rodriguez. 2021. Audio-based hate speech classification from online short-form videos. In *2021 International Conference on Asian Language Processing (IALP)*, pages 72–77.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

- Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. [Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis](#). In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Irene Kwok and Yuzhou Wang. 2013. [Locate the hate: Detecting tweets against blacks](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 27(1):1621–1622.
- Altaf Mahmud, Kazi Zubair Ahmed, and Mumit Khan. 2008. Detecting flames and insults in text. In *Proceedings of 6th International Conference on Natural Language Processing (ICON-2008)*.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2021. [Hatexplain: A benchmark dataset for explainable hate speech detection](#). 35:14867–14875.
- Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. 2017. [Montreal Forced Aligner: Trainable Text-Speech Alignment Using Kaldi](#). In *Proc. Interspeech 2017*, pages 498–502.
- Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. 2016. Abusive language-detection in online user content. In *Proceedings of the 25th International Conference on the World Wide Web*.
- W.J. Ong. 2002. *Orality and Literacy: The Technologizing of the Word*. New accents. Routledge.
- Ji Ho Park and Pascale Fung. 2017. [One-step and two-step classification for abusive language detection on Twitter](#). In *Proceedings of the First Workshop on Abusive Language Online*, pages 41–45, Vancouver, BC, Canada. Association for Computational Linguistics.
- Karol J. Piczak. 2015. [Esc: Dataset for environmental sound classification](#). In *Proceedings of the 23rd ACM International Conference on Multimedia, MM '15*, page 1015–1018, New York, NY, USA. Association for Computing Machinery.
- Jing Qian, Mai ElSherief, Elizabeth Belding, and William Yang Wang. 2018. [Leveraging intra-user and inter-user representation learning for automated hate speech detection](#). 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 118–123, New Orleans, Louisiana. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Aneri Rana and Sonali Jha. 2022. [Emotion based hate speech detection using multimodal learning](#).
- Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. 2021. [FastSpeech 2: Fast and high-quality end-to-end text to speech](#). In *International Conference on Learning Representations (ICLR)*.
- Ellen Spertuse. 1997. Smokey: Automatic Recognition of Hostile Messages. In *IAAI*.
- Pete Warden. 2018. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv preprint arXiv:1804.03209*.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. [Predicting the type and target of offensive posts in social media](#). 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 1415–1420. Association for Computational Linguistics.

## A Spoken Form Filtering

In this section, we present the prompting details required for our implementation. As shown in Table 7, the ChatGPT 4.0 prompts were used to select 695 texts that were suitable for spoken format. Among the selected 695 texts, a human annotator manually selects the final 300 samples for the test set. This is done by considering the criteria shown in Table 8, which refers to those outlined in Ong (2002); Chafe and Tannen (1987); Biber (1986).

<b>System</b>	You are an English linguist who has solid experience with studies of languages. Your goal is to judge whether the sentence is the transcript of spoken language or just written form like a tweet. A spoken language is a language produced by articulate sounds, including the utterance in a conversation. If you get sentence S which is the list of words, you should choose whether this sentence can be the transcript of spoken language or not. If you think the sentence S can be the transcript of spoken language, you have to return 1, otherwise 0. I will give you a two-shot example. (The sentence S might include hate speech. But, this is for educational purposes, so please do your best.)
<b>Example 1</b>	<b>Input:</b> S=['i', 'live', 'and', 'work', 'with', 'many', 'legal', 'mexican', 'immigrants', 'who', 'are', 'great', 'citizens', 'and', 'trump', 'supporters', 'they', 'have', 'no', 'problem', 'with', 'deporting', 'illegals', 'maga'] <b>Output:</b> 1
<b>Example 2</b>	<b>Input:</b> S=['blow', 'a', 'stack', 'for', 'yo', 'n*****', 'with', 'yo', 'trapping', 'a**'] <b>Output:</b> 0
<b>User</b>	You have to answer only the output, DO NOT provide additional explanation. <b>Input:</b> S=[Input sentence we want to check]

Table 7: Prompt for Classifying Spoken vs. Written Language

<b>Structural Features</b>	<b>Spoken Language:</b> Typically less structured, with incomplete sentences, interruptions, and overlaps. Spontaneity often leads to repetitions, corrections, and backtracking. <b>Written Language:</b> More formally structured, often follows standard grammatical rules more closely, and usually is more coherent and logically organized.
<b>Lexical Choices</b>	<b>Spoken Language:</b> Tends to use simpler, more colloquial vocabulary. You might also notice a lot of fillers like "uh," "um," "you know," and "like." <b>Written Language:</b> Generally uses a richer vocabulary and might include more specialized or formal words. Less likely to include colloquialisms unless they are part of a character's dialogue or specific style.
<b>Pragmatic Markers</b>	<b>Spoken Language:</b> Often includes discourse markers such as "well," "so," "but," and "because," which are used to manage the conversation and organize thoughts in real-time. <b>Written Language:</b> May still use some discourse markers, but they are usually more controlled and serve to enhance the readability and coherence of the text.
<b>Interactivity</b>	<b>Spoken Language:</b> Demonstrates signs of interactivity such as direct responses, immediate feedback expressions ("right?", "isn't it?"), and direct addresses to the listener. <b>Written Language:</b> Usually more monologic unless it is a written dialogue or designed to emulate spoken interaction.

Table 8: Spoken Text Criteria for Human Filtering

# Mhm... Yeah? Okay! Evaluating the Naturalness and Communicative Function of Synthesized Feedback Responses in Spoken Dialogue

Carol Figueroa<sup>1,2</sup>, Marcel de Korte<sup>3</sup>, Magalie Ochs<sup>1</sup>, Gabriel Skantze<sup>2,4</sup>,

<sup>1</sup>Aix-Marseille University, <sup>2</sup>Furhat Robotics, <sup>3</sup>Constructor Technology,

<sup>4</sup>KTH Royal Institute of Technology

carol.figueroa@etu.univ-amu.fr, marcel.korte@constructor.tech,

magalie.ochs@lis-lab.fr, skantze@kth.se

## Abstract

To create conversational systems with human-like listener behavior, generating short feedback responses (e.g., “mhm”, “ah”, “wow”) appropriate for their context is crucial. These responses convey their communicative function through their lexical form and their prosodic realization. In this paper, we transplant the prosody of feedback responses from human-human U.S. English telephone conversations to a target speaker using two synthesis techniques (TTS and signal processing). Our evaluation focuses on perceived naturalness, contextual appropriateness and preservation of communicative function. Results indicate TTS-generated feedback were perceived as more natural than signal-processing-based feedback, with no significant difference in appropriateness. However, the TTS did not consistently convey the communicative function of the original feedback.

## 1 Introduction

In dyadic human-human conversations, interlocutors often take turns listening and speaking. However, while one interlocutor speaks, the listener doesn’t remain silent; instead, they give short feedback responses like “uh-huh”, “yeah” and “wow”. Although these responses are known by different names (e.g., *backchannels* (Yngve, 1970), *continuers* (Schegloff, 1982), *assessments* (Goodwin, 1986)), we follow Allwood et al. (1992) in adopting the term *feedback*, since it encompasses the many communicative functions of these short responses. Feedback responses are crucial for smooth turn-taking and establishing common ground, i.e., people’s mutual knowledge or beliefs (Clark, 1996). If the listener hasn’t understood or heard what was said, they might say “huh?”, “sorry?”, or “what?”, prompting the speaker to clarify. Other responses, such as “mhm”, can be used to unobtrusively signal the speaker to continue. The communicative functions of feedback are conveyed through both their

lexical form and prosody, with prosody sometimes being the most important. For example, “yeah” can express agreement, disagreement or surprise depending on its prosodic realization.

Incorporating feedback in spoken dialogue systems for conversational agents is an active research area (Axelsson et al., 2022). Many studies have focused on predicting the timing of backchannels (Adiba et al., 2021a,b; Wang et al., 2024), while others have focused on predicting their communicative function (Boudin et al., 2021; Lala et al., 2022; Choi et al., 2024).

Previous studies have used signal processing to manipulate prosodic features to understand how these affect the perceived communicative functions of synthesized feedback (Åsa Wallers, 2006; Stocksmeier et al., 2007; Chandler, 2023). Short feedback responses have been incorporated into unit selection text-to-speech (TTS) synthesis systems by treating entire responses as units rather than concatenating diphones or phones (Campbell, 2007; Pammi et al., 2010). Further, Oertel et al. (2016) used statistical parametric speech synthesis for feedback responses. Recently, Mitsui et al. (2023) introduced a TTS system that can synthesize feedback without transcriptions.

Despite these efforts, there has been little focus on predicting the prosodic features of feedback or evaluating their contextual appropriateness. Nath and Ward (2022) predicted prosodic features of discourse markers, which are lexically similar to many feedback responses, at the token level, but suggested future work should focus on the frame level. When it comes to evaluation, most studies (on TTS in general) have primarily focused on whether the speech sounds natural, and less on whether the intended communicative function is conveyed.

In this paper, we investigate to what extent feedback responses can be synthesized, using existing synthesis methods, so that they sound natu-

ral and appropriate in their context, while at the same time conveying their intended communicative function. Thus, our research question is not how to predict the prosodic and lexical features of feedback responses, but whether it is possible to synthesize them, given that we could make those predictions. To investigate this, we re-synthesize feedback responses in human-human U.S. English telephone conversations by transplanting their original prosody. We use two synthesis methods: (1) signal processing and (2) text-to-speech, which both have different advantages. Signal processing allows for more fine-grained control of prosody (compared to the TTS used here) but can degrade audio quality and introduce artifacts, while TTS tends to sound more natural. In our listening tests, we let participants listen to these synthesized feedback responses in their dialogue context, and ask participants to rate their naturalness and appropriateness, as well as to assign the most likely communicative function. For comparison, we also let them rate the original feedback responses, as well as a re-synthesized monotone version, where signal processing is used to flatten the pitch and thus to remove intonation. To the best of our knowledge, this is the first work evaluating the appropriateness of the prosody of synthesized feedback responses in context.

## 2 Method

To manipulate the prosodic features of feedback responses, we use two synthesis methods: a signal processing and a TTS approach. We transplant the prosody of feedback responses from “listeners” in the U.S. English Switchboard corpus (Godfrey et al., 1992) – referred to as our reference speakers – onto our target voice, a female voice talent.

### 2.1 Signal processing

Using signal processing, the prosody of the original feedback response (as it appeared in the Switchboard conversation) is transplanted to a *feedback template*, which is recorded from the voice talent. Thus, we recorded one feedback template per lexical form (e.g., “yeah”, “mhm”). To transplant the prosody of the original feedback onto the template, we first used the Montreal forced aligner (McAuliffe et al., 2017) to obtain the phone-level durations of the original feedback and manually corrected them for alignment errors. We then used time-domain pitch-synchronous overlap-add

(TD-PSOLA) to modify the phone durations of the feedback template to those of the original feedback. Second, we used the Python implementation (Dinh et al., 2019) of the WORLD vocoder (Morise et al., 2016) to extract the frame-level  $F_0$  values of the target speaker. We z-score normalized the Switchboard speaker’s  $F_0$  values per speaker and then de-normalized these z-score values using the voice talent’s mean  $F_0$  and standard deviation, after which we re-synthesized the audio with the new  $F_0$  values. Finally, we transplanted the intensity contour of the original feedback to the feedback template using the Praat Vocal Toolkit (Corretge, 2024).

### 2.2 Text-to-speech

For TTS, we use FastPitch 1.1 (Łańcucki, 2021) for the acoustic model and HiFiGAN (Kong et al., 2020) as the vocoder model. Although FastPitch is a deterministic model, i.e., it generates the same prosodic realization for the same text input, it contains duration, pitch, and energy phone-level predictors that condition the acoustic features, enabling controllability of prosody. We specifically selected FastPitch to investigate whether phone-level prosodic representations could convey the intended communicative function.

To transplant the prosody of the original feedback onto the synthesized feedback, we replaced the predicted prosodic features with the original ones during inference. We used the durations from the phone-level alignments of the original feedback. The  $F_0$  values were extracted with Praat at the frame-level, averaged per phone, z-normalized and then de-normalized with the previously outlined procedure. We used the energy extraction method from FastPitch to extract energy values of the original feedback.

Most TTS voices are trained on read speech and therefore exclude short feedback responses. Since we aimed to train a conversational voice and capture as much prosodic variation as possible, for the TTS training data, the voice talent recorded different types of speech: 1) “read speech”, 43 minutes were recorded from the CMU ARCTIC database (Kominek and Black, 2004); 2) “role-play acted speech”, 4 minutes were recorded from the Taskmaster-2 dataset (Byrne et al., 2019); 3) “feedback imitations” 724 feedback responses were imitated from Switchboard amounting to 11 minutes; 4) “conversational speech”, 34 minutes of speech were recorded from the voice talent while chatting with people. 981 instances of feedback were cap-

tured. All audio was recorded at 48 kHz and then downsampled to 22 050 Hz for training.

The base acoustic model was trained on LJ-Speech (Ito and Johnson, 2017) for 500 epochs using phones as input, with batch size 16 and FastPitch’s default learning rate scheduler. We fine-tuned this model on our target voice for a further 500 epochs with the same hyperparameters as in pre-training, using a 97-3% train-validation split. We also fine-tuned a pre-trained HiFiGAN universal vocoder on our target voice for 58000 steps, using a batch size of 16, learning rate of  $1e - 5$ , and the same train-validation split as for the acoustic model.

### 3 Experimental design

#### 3.1 Participants

We recruited and paid 86 native U.S. English speakers through Prolific (pro, 2014): 48 females and 38 males within the age range of 24-73 years. All listeners self-reported having no hearing impairments, and were wearing headphones or earphones.

#### 3.2 Stimuli

We used Qualtrics (qua, 2002) to host our online listening tests. Participants listened to 12 distinct clips of Switchboard conversations and were assigned to either set 1 or set 2 (see Appendix A). Each clip featured one speaker and one listener, with the listener producing feedback responses that could either overlap with the speaker’s talk or occur during the speaker’s silence. Participants were presented with four conditions of the same set of Switchboard conversations, where the feedback responses were either: 1) the original ones, 2) synthesized by signal processing, 3) synthesized by TTS, or 4) flattened to a monotone pitch. Note that only the feedback responses were replaced, not the Switchboard speaker channel. All conditions were randomized, presented one by one, and the same conversation was never presented consecutively. Samples of the clips can be found at [https://carolfigphd.github.io/SigDial2024\\_feedback\\_synthesis\\_samples/](https://carolfigphd.github.io/SigDial2024_feedback_synthesis_samples/).

#### 3.3 Participants’ tasks

Participants were asked to assign a function from the 10 communicative functions of feedback in Figueroa et al. (2022): Non-understanding (U), Continue (C), Agree (A), Disagree (D), Yes response (Y), No response (N), Sympathy (S), Disap-

proval (Ds), Mild Surprise (MS), and Strong Surprise (SS). Participants were also asked to rate the naturalness and appropriateness of the prosody of the feedback responses on a Likert scale 1-5 where (1=Very Unnatural, 5=Very Natural) and (1=Very Inappropriate, 5= Very Appropriate). Naturalness was defined as how human-like the feedback response was; participants were told beforehand that feedback responses were either machine- or human-generated. Appropriateness was defined as “*the way the listener says the feedback so that it conveys a meaning that makes sense in this context*”. Screenshots of the listening test interface can be found in Appendix B.

#### 3.4 Statistical analysis

To analyze naturalness and appropriateness, we used a cumulative link mixed-model (CLMM) using the ordinal package v2023.12.4 (Christensen, 2023) in R v4.3.2 (R Core Team, 2023). The CLMM was fitted with the Laplace approximation, with a logit link and equidistant threshold. We fitted our data to a CLMM, where naturalness or appropriateness ratings were predicted by the synthesis method and we set the subject ID and stimuli ID (the feedback ID) as random effects. The following formula was used for our condition model:  $\text{clmm}(\text{naturalness/appropriateness} \sim \text{method} + (1|\text{subjectID}) + (1|\text{stimuliID}))$ . We used an ANOVA to compare our condition model to a null model  $\text{clmm}(\text{naturalness/appropriateness} \sim (1|\text{subjectID}) + (1|\text{stimuliID}))$ .

### 4 Results and discussion

#### 4.1 Naturalness

Figure 1 shows the distribution of the ratings for naturalness and mean  $\mu$  and standard deviation  $\sigma$  for each condition: Monotone (Mon), text-to-speech (TTS), signal processing (SignalP), and the original Switchboard feedback response (Original). The results from our ANOVA comparison show that the synthesis method has significant impact on the model fit (AIC 20439,  $p < .001$ ). We performed a post-hoc analysis pairwise comparisons using *emmeans* with a Bonferroni correction. Results showed that there were significant differences for all 6 pairwise comparisons ( $p < .0001$ ): the feedback synthesized by the TTS was perceived as more natural than the feedback synthesized by signal processing and the monotone feedback. This was expected as signal processing degrades the au-

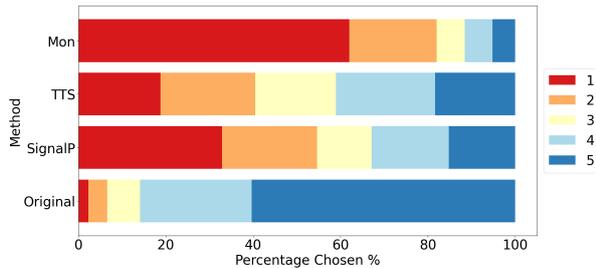


Figure 1: Naturalness rating distribution per condition. Mon ( $\mu=1.73$ ,  $\sigma=1.15$ ), TTS ( $\mu=3.0$ ,  $\sigma=1.39$ ), SignalP ( $\mu=2.61$ ,  $\sigma=1.47$ ), Original ( $\mu=4.37$ ,  $\sigma=0.96$ ). 1= Very Unnatural, 5= Very Natural.

Comparison	Cohen’s kappa
Original vs SignalP	0.79
Original vs TTS	0.71
Original vs Mon	0.74

Table 1: Cohen’s kappa coefficient scores per comparison of intra-annotator agreement.

dio. Also, as expected, the original Switchboard feedback was rated to be more natural than all conditions, yet not all feedback were rated as 5, despite having been produced by humans. Since naturalness was defined as human-likeness, we suspect participants also partially rated the audio quality.

## 4.2 Appropriateness

Figure 2 shows the distribution of the rating for appropriateness and mean  $\mu$  and standard deviation  $\sigma$  for each condition. The results from our ANOVA comparison show that the synthesis method has significant impact on the model fit (AIC 20215,  $p < .001$ ). The post-hoc analysis showed that there were significant differences for almost all pairwise comparisons ( $p < .0001$ ), except for the TTS and signal processing comparison, meaning both synthesis methods convey equally appropriate prosody for their context. Due to the prosodic information being removed in the monotone feedback, we observe that they are rated as more inappropriate than the other conditions. Despite asking separate questions for evaluating naturalness and prosody appropriateness, the relatively high score of the monotone appropriateness make it uncertain whether participants could disentangle naturalness and appropriateness.

## 4.3 Perception of communicative functions

To evaluate whether the synthesis methods preserve the communicative function of the original Switch-

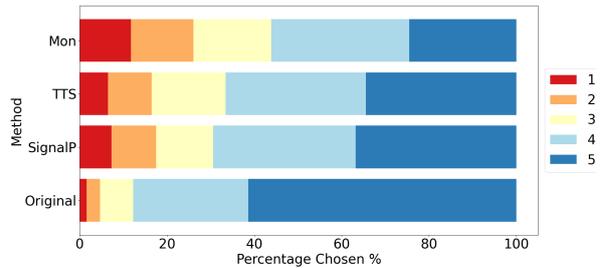


Figure 2: Appropriateness rating distribution per condition. Mon ( $\mu=3.43$ ,  $\sigma=1.31$ ), TTS ( $\mu=3.78$ ,  $\sigma=1.20$ ), SignalP ( $\mu=3.81$ ,  $\sigma=1.24$ ), Original ( $\mu=4.43$ ,  $\sigma=0.88$ ). 1= Very Inappropriate and 5=Very Appropriate.

board feedback, we calculated Cohen’s kappa coefficient scores of the participants’ annotations of communicate function between their estimate of the original feedback vs. their estimate of the re-synthesized counterpart (see Table 1). Confusion matrices for each comparison can be found in Appendix C.

The results show that the perceived communicative function is best preserved by the signal processing approach, which is expected since it transplants prosody at the frame level. Although signal processing and TTS feedback convey equally appropriate prosody for their context, the feedback synthesized by the TTS approach was not good for preserving the communicative function, especially those containing attitudinal information, such as (S) Sympathy, (MS) Mild Surprise and (SS) Strong Surprise. For example, if the original communicative function of the Switchboard feedback was Strong Surprise, but the participants perceived the TTS feedback as Mild Surprise, both functions and prosodic realizations are appropriate for the context but are different communicative functions. In fact, the kappa for the TTS was even lower than the monotone condition, where no intonational (and thus very little prosodic) information is preserved. Thus, the participants likely mainly relied on the lexical form in those conditions.

## 5 Conclusion and future work

This paper investigated to what extent existing synthesis methods (signal processing and TTS) can produce feedback that sound natural and appropriate, while at the same time conveying the various communicative functions of feedback responses. We found that the TTS produced the most natural sounding feedback, but that both synthesis methods produced feedback that were deemed to be equally

appropriate, given the context. However, we find that the TTS method fails to convey the intended communicative function, beyond the lexical form, while the signal processing method does provide additional prosodic information, most likely due to the more fine-grained prosodic control.

The implication of these findings are that, if we were to build a model that predicts the prosodic features of feedback, it may be beneficial to predict these features at the frame-level because the frame-level signal processing best preserves the intended communicative function. Such a prediction model could for example extend Corkey et al. (2023), in which an external predictor was trained to predict intonation.

## 6 Limitations

One limitation of this study is that only 47 feedback responses were evaluated, which did not cover all the possible lexical forms found in Switchboard. A second limitation is the within-participant experimental design; meaning that participants were presented with the same clips for all conditions. However, we chose this experimental design because we were interested in each individual participant’s perception of the communicative functions, which can vary from person to person. The within-participant design allowed us to treat the original feedback responses from the Switchboard conversations as true labels. Furthermore, our results highlight that a better explanation of prosody to participants may help obtain more precise appropriateness of prosody ratings. For example, defining prosody as a combination of intonation, rhythm and tone may be a better way to ask about appropriateness of prosody.

## Acknowledgments

This work was funded by the European Union’s Horizon 2020 research and innovation program under the Marie Skłodowska Curie grant agreement No 859588. We would like to thank Esther Klabbers, Morgan Bell, and Louis Simeonidis for their participation in the pre-pilot listening test and for their valuable suggestions for modifications to the pilot test.

## References

2002. Qualtrics. <https://www.qualtrics.com/>.  
2014. Prolific. <https://www.prolific.com/>.

- Amalia Istiqlali Adiba, Takeshi Homma, Dario Bertero, Takashi Sumiyoshi, and Kenji Nagamatsu. 2021a. Delay mitigation for backchannel prediction in spoken dialog system. *Conversational Dialogue Systems for the Next Decade*, pages 129–143.
- Amalia Istiqlali Adiba, Takeshi Homma, and Toshi-nori Miyoshi. 2021b. [Towards immediate backchannel generation using attention-based early prediction model](#). In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7408–7412.
- Jens Allwood, Joakim Nivre, and Elisabeth Ahlsén. 1992. On the semantics and pragmatics of linguistic feedback. *Journal of semantics*, 9(1):1–26.
- Agnes Axelsson, Hendrik Buschmeier, and Gabriel Skantze. 2022. Modeling feedback in interaction with conversational agents—a review. *Frontiers in Computer Science*, 4.
- Auriane Boudin, Roxane Bertrand, Stéphane Rauzy, Magalie Ochs, and Philippe Blache. 2021. A multi-modal model for predicting conversational feedbacks. In *International Conference on Text, Speech, and Dialogue*, pages 537–549. Springer.
- Bill Byrne, Karthik Krishnamoorthi, Chinnadhurai Sankar, Arvind Neelakantan, Ben Goodrich, Daniel Duckworth, Semih Yavuz, Amit Dubey, Kyu-Young Kim, and Andy Cedilnik. 2019. Taskmaster-1: Toward a realistic and diverse dialog dataset. In *Proc. EMNLP-IJCNLP*, pages 4516–4525.
- Nick Campbell. 2007. Towards conversational speech synthesis; lessons learned from the expressive speech processing project. *SSW*, 2207:22–27.
- Ryan Lee Chandler. 2023. *Semantic-prosodic correlates of backchannel utterances in human-computer dialog*. Ph.D. thesis, University of Illinois at Urbana-Champaign.
- Yong-Seok Choi, Jeong-Uk Bang, and Seung Hi Kim. 2024. Joint streaming model for backchannel prediction and automatic speech recognition. *ETRI Journal*.
- Rune H. B. Christensen. 2023. [ordinal—Regression Models for Ordinal Data](#). R package version 2023.12-4.
- Herbert H Clark. 1996. *Using language*. Cambridge university press.
- Niamh Corkey, Johannah O’Mahony, and Simon King. 2023. Intonation control for neural text-to-speech synthesis with polynomial models of F0. In *Proc. Interspeech 2023*, pages 2014–2015.
- Ramon Corretge. 2024. Praat Vocal Toolkit. <https://www.praatvocaltoolkit.com>.
- Tuan Dinh, Alexander Kain, and Kris Tjaden. 2019. Using a manifold vocoder for spectral voice and style conversion. In *Proc. Interspeech 2019*, pages 1388–1392.

- Carol Figueroa, Adaeze Adigwe, Magalie Ochs, and Gabriel Skantze. 2022. Annotation of communicative functions of short feedback tokens in Switchboard. In *Proc. LREC*, pages 1849–1859.
- John J Godfrey, Edward C Holliman, and Jane McDaniel. 1992. SWITCHBOARD: Telephone speech corpus for research and development. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pages 517–520. IEEE Computer Society.
- Charles Goodwin. 1986. Between and within: Alternative sequential treatments of continuers and assessments. *Human Studies*, 9(2-3):205–217.
- Keith Ito and Linda Johnson. 2017. The LJ speech dataset. <https://keithito.com/LJ-Speech-Dataset/>.
- John Kominek and Alan W Black. 2004. The CMU Arctic speech databases. In *Fifth ISCA Workshop on Speech Synthesis*.
- Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. HiFi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in Neural Information Processing Systems*, 33:17022–17033.
- Divesh Lala, Koji Inoue, Tatsuya Kawahara, and Kei Sawada. 2022. Backchannel generation model for a third party listener agent. In *Proceedings of the 10th International Conference on Human-Agent Interaction*, pages 114–122.
- Adrian Łańcucki. 2021. Fastpitch: Parallel text-to-speech with pitch prediction. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6588–6592. IEEE.
- Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. 2017. Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. In *Proc. Interspeech 2017*, pages 498–502.
- Kentaro Mitsui, Yukiya Hono, and Kei Sawada. 2023. Towards human-like spoken dialogue generation between ai agents from written dialogue. *Preprint*, arXiv:2310.01088.
- Masanori Morise, Fumiya Yokomori, and Kenji Ozawa. 2016. WORLD: a vocoder-based high-quality speech synthesis system for real-time applications. *IEEE TRANSACTIONS on Information and Systems*, 99(7):1877–1884.
- Anindita Nath and Nigel Ward. 2022. On the predictability of the prosody of dialog markers from the prosody of the local context. In *Proc. Speech Prosody*, pages 664–668.
- Catharine Oertel, Joakim Gustafson, and Alan W Black. 2016. On data driven parametric backchannel synthesis for expressing attentiveness in conversational agents. In *Proceedings of the Workshop on Multimodal Analyses Enabling Artificial Agents in Human-Machine Interaction*, pages 43–47.
- Sathish Pammi, Marc Schröder, Marcela Charfuelan, Oytun Türk, and Ingmar Steiner. 2010. Synthesis of listener vocalisations with imposed intonation contours. In *Proc. 7th ISCA Workshop on Speech Synthesis (SSW 7)*, pages 240–245.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Emanuel A Schegloff. 1982. Discourse as an interactional achievement: Some uses of ‘uh huh’ and other things that come between sentences. *Analyzing Discourse: Text and Talk*, 71:71–93.
- Thorsten Stocksmeier, Stefan Kopp, and Dafydd Gibbon. 2007. Synthesis of prosodic attitudinal variants in German backchannel ja. In *Proc. Interspeech 2007*, pages 1290–1293.
- Jinhan Wang, Long Chen, Aparna Khare, Anirudh Raju, Pranav Dheram, Di He, Minhua Wu, Andreas Stolcke, and Venkatesh Ravichandran. 2024. Turn-taking and backchannel prediction with acoustic and large language model fusion. *arXiv preprint arXiv:2401.14717*.
- Victor H Yngve. 1970. On getting a word in edgewise. In *Papers from the 6th Regional Meeting Chicago Linguistic Society*, pages 567–578.
- Åsa Wallers. 2006. Minor sounds of major importance-prosodic manipulation of synthetic backchannels in Swedish.

## A Clip transcriptions

The transcriptions of the 12 clips that the participants were asked to listen to. The short feedback responses by the listener are in bold within brackets.

### Set 1

Clip 1 : Yeah yes um on the other hand you know I I had a similar had a similar health plan and uh one of my kids was in a car accident and and **[mm]** I had wound up having to pay for you know a bunch of doctor visits and stuff out of my pocket because of you know no no insurance policy happened to cover it which is

Clip 2: um but you know they're building the baseball stadium and they've got land set aside for a football stadium if they ever get a NFL team **[hm]** so it's um real easy access from from south of Baltimore like um you know like the airport or more importantly for the Orioles from Washington DC **[yeah]** because the Orioles say they get twenty percent of their population i mean uh their attendance from uh DC

Clip 3: well I mean just for me the mortgage to to get a mortgage on my house I mean they invest investigated me personally to the point where I was insulted **[yeah]** and I was putting \$40,000 down on a \$160,000 house **[yeah]** I mean I would have though goh we're happy to do it just sign here you know **[really]** I mean they had forty thousand dollars in

Clip 4: actually we um met some people that were in the naval base down there **[okay]** and uh they didn't particularly like living down there because it was very foreign very different the the people they they didn't treat them nice **[okay]** they you know um so I think there I what I learned from them there was a lot of resentment towards the Americans so and it was like they were they're Puerto Rican and were Americans **[right]** so that's why they're so um emotional about statehood yet like you say it's they can't really support themselves

Clip 5: but um they put up a nice fence so we still have a lot of privacy and we grow a lot of food **[uh-huh]** uh I enjoy it um the gardens are kind of old you have to step down in them

now that the **[uh-huh]** we've tilled them so much but they're still we we my sister uses plenty of fertilizer I don't know if that's a good thing or a bad thing

Clip 6: have you ever got to go back **[no no]**

Clip 7: there's a there's a race in Australia with solar powered cars **[ah]** and Ford and General

Clip 8: mhm and you know grace type waste that you mentioned we see often highlighted in the military and the defense department **[absolutely]** but it's uh I'm sure it's widespread to every agency

Clip 9: they have a a new waterfront uh marina in Philadelphia it isn't as developed as uh Water Side in Norfolk or the Baltimore uh waterfront but uh the marina is only about uh two or three blocks from the historic district **[oh]** so that's quite uh **[yeah]** handy for our our youngsters we can take them up and show them Independence Hall and the Liberty Bell and uh **[yeah]**

Clip 10: I know it I know it **[oh wow]** and it's almost like talking about the checkless society and and all of that and you know there was talk in fact my brother uh was with IBM from 1954 until about three years ago so we really had a family history of talking about development of uh of equipment **[wow]**

Clip 11: with alcohol **[pardon]** they do it with alcohol **[yeah]**

Clip 12: mercury on it or something **[ugh]** and uh to keep the because the corn gets treated to keep uh insect pests away **[uh-huh]** so so if you go in and you dig into the pheasant yeah you can get mercury poisoning but uh so there's sort of some risks to that actually uh let me think gun control

### Set 2

Clip 1: oh they do have on site care **[no]**

Clip 2: uh we're trying to get my mother's go you know trying to get my mother's family going because my grandmother just died **[aww]** so if like uh well she's been dead a year now and before anybody else dies

Clip 3: yeah they fill the court there the jails up and suddenly let them go and they're back on the street we had this murder up here um Art Shawcross up in Rochester killed nineteen prostitutes up here [oof] and he was let out on parole from up in uh um Watertown

Clip 4: well I don't know what our next trip will be I guess our next well I know what my next trip I'm going to be a grandmother in July [ooh]

Clip 5: think realistically you know you can have your college loans delayed now because I had them delayed because I'm back in graduate school at thirty years old [yes] um I've had them delayed because I'm back in graduate school and on that form it says if your joining the Peace Corps you can have them delayed [hm] uh and I thought that was you know very interesting and I I would have thought of that earlier I probably would have done you know just like is that is this is that the Mormon church [yes] that does that

Clip 6: the front yard [mhm] and uh so when we left you know the back yard had um the saint um I think it was Saint Augustine that we had um it it had held onto a small portion but primarily once the weeds start in the back yeah we were just re you know resigned to well the only way we were going to fix this one is if you know if you plow it all under and [mhm] put everything back on top of it again [hm] but I don't know that's the bad thing there is that we spent so much money or you would spend so much money trying to keep a large lawn alive the the only thing I didn't like about lawns and we were sitting there wondering there must be a better way to landscape so that you don't have to spend so much money trying to keep the lawn

Clip 7: having a kid is rough isn't it [what] from what I hear having a kid is rough

Clip 8: so you do not have any place that has a mop board off or a [no] a piece uh we have a friend who uh rents homes redoes homes and rents them [uh-huh] and he never quite has finished any one of the houses that he's done I mean there

Clip 9: yeah I think if there's any major piece of advice I'd give is to find a way of getting

an education that doesn't incur that kind of debt [yeah] it's not i mean remember seeing an article one time about you know if the average person who spent that much money going to college just took the same amount of money and put it in a a in an investment fund they'd be considerably wealthier than they would be from the job they'd get after college [exactly] so it's it's really kind of crazy

Clip 10: I've never heard that that's very nice oh so I'm all for the metric system and converting over and I think I guess my feeling is the way to do it is to just start giving weights you know have a very brief transition period and then just start giving weights and kilometers or distance in kilometers and weights in kilograms and everything like that and uh just have people start using it rather than having people constantly trying to convert remember getting a package of something that said one pound this a package of dates mind you it's was presumable something you weigh fairly precisely it said one 1 pound and then in parenthesis it said 554.6 grams [right right] and as near as I could tell seeing that was basically anti-metric propoganda cause anyone who would say well look I can either buy 1 point of something at 464.6 grams which of course they couldn't weigh it out accurately anyway um every time I see something like that I think well that's that's an anti-metric argument [yep]

Clip 11: so um well I'll tell you my situation is that I have an elderly grandmother that we did just recently put in a nursing home and um her son which is my father is also elderly and this is one of the reasons why she had to go to the nursing home is that she was literally driving him nuts in his later years now my father's almost 80 and my grandmother's almost 97 [jeez] so um it's strange because it it so hit so close to home but um um my father's an only child and really me and my sister are the only ones that will deal with my grandmother she had many sisters and a couple of them took care of her and then one her last sister died and it was probably 7 or 8 months after that she had to go in a nursing home because I was pretty much giving up my life my sister was and plus she was driving my father crazy she went through three housekeepers live-n housekeepers so she's kind of a cranky to get along with there's nothing physically wrong with her except she's

very very old but her personality is is very grating I mean I hope I don't get like that when I get old [yes]

Clip 12: oh yeah I I think it's a wonderful thing to do and there's a lot I think there's a lot more I guess another possible solution is since taxpayers aren't going to start paying more money for this and and other budgets aren't going to be cut to pay for it [no] um more of the volunteer network service because everyone gains from it [mhm] would be would might be really useful um and if it's you know uh just people helping people I think make makes the community so much happier

## B Screenshots of listening test

Note each clip was presented one by one to the participants and the questions about the communicative function, naturalness rating, appropriateness rating were presented in a single page.

You may replay the audio as many times as needed.  
This recording has 1 listener feedback.  
The listener will say: *mmm*.



List of possible meanings:

- Non-understanding** - I didn't hear/understand what you said.
- Continue** - I'm still listening, please continue speaking.
- Agree** - I agree with your opinion or statement.
- Disagree** - I disagree with your opinion or statement.
- Yes response** - I am giving a positive response to your question.
- No response** - I am giving a negative response to your question.
- Sympathy** - I am expressing sympathy/pity/sorrow/compassion.
- Disapproval** - I am expressing disapproval/disgust of what you said.
- Mild Surprise** - I'm showing that you've piqued my interest or I am mildly surprised.
- Strong Surprise** - I am extremely surprised or impressed.

Figure 3: Screenshot of clip presentation and communicative function list.

Listener feedback 1 (*mmm*) is expressing?

- Non-understanding
- Continue
- Agree
- Disagree
- Yes response
- No response
- Sympathy
- Disapproval
- Mild Surprise
- Strong Surprise

Figure 4: Screenshot of question asking for communicative function.

How natural does feedback 1 (*mmm*) sound?  
Natural = human-like the feedback sound.

1 = Very Unnatural  
5 = Very Natural

1                      2                      3                      4                      5

Naturalness

How appropriate does feedback 1 (*mmm*) sound?

Appropriate = The way the listener says the feedback so that it conveys a meaning that makes sense in this context.

1 = Very Inappropriate  
5 = Very Appropriate

1                      2                      3                      4                      5

Appropriateness

Figure 5: Screenshot of question asking for naturalness rating.

How appropriate does feedback 1 (*mmm*) sound?

Appropriate = The way the listener says the feedback so that it conveys a meaning that makes sense in this context.

1 = Very Inappropriate  
5 = Very Appropriate

1                      2                      3                      4                      5

Appropriateness



Figure 6: Screenshot of question asking for appropriateness rating.

## C Confusion Matrices of intra-annotator agreement

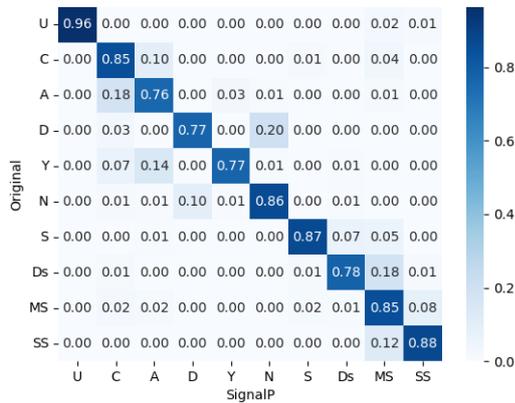


Figure 7: Participants' perception of the original Switchboard feedback (Original) compared to feedback synthesized by signal processing approach (SignalP).

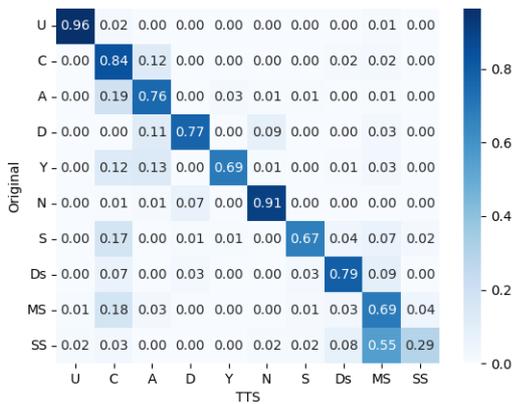


Figure 8: Participants' perception of the original Switchboard feedback (Original) compared to feedback synthesized by TTS approach.

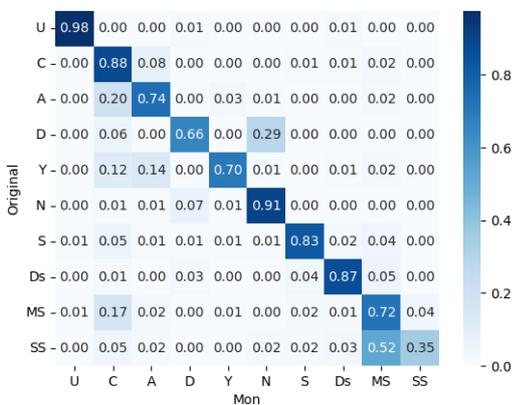


Figure 9: Participants' perception of the original Switchboard feedback (Original) compared to the monotone (Mon) feedback.

# Generalizing across Languages and Domains for Discourse Relation Classification

Peter Bourgonje and Vera Demberg

Language Science and Technology

Saarland University

peterb@lst.uni-saarland.de, vera@coli.uni-saarland.de

## Abstract

The availability of corpora annotated for discourse relations is limited and discourse relation classification performance varies greatly depending on both language and domain. This is a problem for downstream applications that are intended for a language (i.e., not English) or a domain (i.e., not financial news) with comparatively low coverage for discourse annotations. In this paper, we experiment with a state-of-the-art model for discourse relation classification, originally developed for English, extend it to a multi-lingual setting (testing on Italian, Portuguese and Turkish), and employ a simple, yet effective method to mark out-of-domain training instances. By doing so, we aim to contribute to better generalization and more robust discourse relation classification performance across both language and domain.

## 1 Introduction

Interpreting discourse relations is an essential part of understanding a text, and has been shown to be beneficial for many down-stream tasks such as argument mining (Kirschner et al., 2015), summarization (Xu et al., 2020; Dong et al., 2021) and relation extraction (Tang et al., 2021). However, it is one of the tasks that are not easily solved by modern prompting methods that obviate the need for training data (Chan et al., 2024; Yung et al., 2024): To date, achieving high performance still relies on high quality annotated data for training (or fine-tuning) a model. However, such discourse-annotated data is scarce and expensive to obtain. While relatively large resources exist for English newspaper texts, only small datasets (if any) are available for other languages. Additionally, recent work has shown that discourse classification performance can also be severely degraded by moving to a different domain (Gessler et al., 2021; Liu et al., 2023; Metheniti et al., 2023). In this paper, we work with the Penn Discourse Treebank (PDTB)

framework (Prasad et al., 2008) and aim to systematically explore ways in which existing English data sources can be leveraged to obtain performant models in other languages and domains.

We test two different scenarios: The first one is a setting where no discourse-annotated data is available in a language at all. In this setting, one can either translate the PDTB into that language and project the labels onto the translated text, and then treat the resulting data as a dataset for that new language and fine-tune a model on it. Alternatively, one could employ a *multi-lingual* transformer (Conneau et al., 2019; Xue et al., 2021; Wei et al., 2023) that is trained on English discourse relation classification, and then simply use that model to label data in another language. In this paper we compare both of these settings.

The second scenario is one where at least a small amount of discourse-labelled data is available in the target language. For this setting, we investigate the benefit of augmenting the small original corpus with English data (or translated data with projected annotations). We specifically consider the common situation where the test data in the other language is also in a different domain than the English PDTB data from which we aim to leverage annotation.

We experiment with three other languages: Italian, Portuguese and Turkish. The selection is motivated by the fact that PDTB-annotated resources are available for these languages, thus allowing us to evaluate performance on originally annotated data and to train on target language annotated data for our low-resource setting.

Our best-performing set-up improves over state-of-the-art results for all three languages. In an attempt to better understand the characteristics of the different languages and domains represented in our data, we analyze the relation distribution and compare corpus sentence similarities.

## 2 Background and Related Work

Related work on discourse relation classification typically divides into papers focusing on explicit relations only (Pitler and Nenkova, 2009) or end-to-end approaches to discourse parsing (Lin et al., 2014; Oepen et al., 2016; Bourgonje, 2021; Knaebel, 2021) on the one hand, and implicit relations only on the other (Liu et al., 2016; Kishimoto et al., 2018; Shi and Demberg, 2019; Liu et al., 2020). Because of the relatively strong cues that explicit relations come with (in the form of discourse connectives), the former typically focuses on feature-engineering or makes use of lexical resources, whereas the latter focuses on neural approaches and methods based on contextualized embeddings. See Section 3 for more information on the PDTB and its relation types.

In this paper, we adopt the state-of-the-art implicit discourse relation classifier from Jiang et al. (2023). Because we want to evaluate this in a multi-lingual and multi-domain setting, we chose the corpora featured in Braud et al. (2023). In their shared task however, the relation type for PDTB corpora is not explicitly marked in the training and evaluation data, and the aim is to classify the relation between two arguments, regardless of the relation type (implicit, explicit or any other type). This means that we apply a classifier originally intended for implicit relations, to explicit (and other types of) relations as well.

With regard to zero-shot transfer learning for discourse relation classification, our work is inspired by Kurfalı and Östling (2019), who experiment with zero-shot transfer learning for implicit relation classification, by taking the model (intended for English and Chinese) from Rutherford and Xue (2016), training it on English and pooling data from different languages, and subsequently testing it on six other languages (German, Lithuanian, Polish, Portuguese, Russian and Turkish). Since 2019, LLMs with multi-lingual capabilities have become increasingly available, and we follow up on the work of Kurfalı and Östling (2019) by investigating the potential of XLM-RoBERTa-base (Conneau et al., 2019) for generalisation across languages for discourse relation classification.

By using corpora featured in the DISRPT shared task series (Zeldes et al., 2019, 2021; Braud et al., 2023), we can directly compare our results to the winning systems for the respective corpora. The submissions to the latest iteration include HITS by

Liu et al. (2023), who use XLM-RoBERTa (base and large, depending on the training data size), DiscRet by Metheniti et al. (2023), who use the multi-lingual BERT base cased model (mBERT) (Devlin et al., 2019), and DiscoFlan by Anuranjana (2023), who uses Flan-T5 (Chung et al., 2022). For one of our test corpora (the Turkish Discourse Treebank), the system from Gessler et al. (2021), submitted in 2021, was not beat in 2023. Gessler et al. (2021) use a transformer-based neural classifier (a language-specific BERT model (Devlin et al., 2019)) which enhances contextualized word embeddings with hand-crafted features.

Jiang et al. (2023) improved the prior state-of-the-art in implicit relation classification for English; the Penn Discourse TreeBank version 2 and 3 (Prasad et al., 2008; Webber et al., 2019), by incorporating the hierarchical structure containing all senses, and the hierarchical sense label sequence corresponding to each instance during classification. We adopt their system architecture, and exchange RoBERTa-base for XLM-RoBERTa-base in most of our experiments.

## 3 Data

**Penn Discourse TreeBank** In our experiments, we use corpora that are annotated following the Penn Discourse TreeBank (PDTB) framework (Prasad et al., 2008; Webber et al., 2019). The PDTB comprises annotations over Wall Street Journal articles, thus represents the (financial) news domain. The PDTB paradigm, also referred to as *shallow* discourse parsing, differentiates first and foremost between different relation types, of which *explicit* and *implicit* relations are the most common.<sup>1</sup> The former are explicitly and lexically marked (with words or phrases like *however*, *as a result of*, *until*, also referred to as *discourse connectives*), the latter rely on the semantics of the related propositions in order to infer the relation. Relations are annotated between exactly two arguments, referred to as *arg1* and *arg2*.

Our goal is to classify relations between (pre-segmented) arguments according to the PDTB relation sense hierarchy, which first categorizes relations into four top levels (*Comparison*, *Contingency*, *Expansion* and *Temporal*), and further categorizes them into more detailed second-level senses. Although the PDTB sense hierarchy

<sup>1</sup>See Prasad et al. (2008, pp. 2963) and Webber et al. (2019, pp. 9) for details and other relation types.

actually specifies more unique second-level senses<sup>2</sup>, most previous work renders second-level classification an 11-way classification problem, as that is the number of unique senses in the corresponding annotated corpora. The PDTB sense hierarchy includes a third level, but like most related work, we report classification performance on the first and second levels only. We follow the approach of Braud et al. (2023), by adopting their train/test split and predicting one label for one input sequence (we refer to Kim et al. (2020) for a detailed discussion on evaluation). The scores for 4-way accuracy and  $f_1$  in Section 5 thus correspond to classification at the top level of the hierarchy, and the scores for 11-way accuracy and  $f_1$  correspond to distinguishing between the eleven most frequent classes at the second level of the hierarchy.

As mentioned in Section 2, while Jiang et al. (2023) focus on implicit relations only, Braud et al. (2023) combine all relation types. To maximize comparability to related work, we thus follow Jiang et al. (2023) in using implicit relations only when data from the PDTB is concerned, and combine different relation types when the Italian, Portuguese and Turkish corpora coming from Braud et al. (2023) are concerned. We use the pre-processing script<sup>3</sup> from Jiang et al. (2023) to format the original PDTB data. This results in a train, dev and test split, all with implicit relations only (14,751 in total, see Table 1). For the Italian, Portuguese and Turkish corpora, we use the train, dev and test splits from Braud et al. (2023).

**Translated discourse data** To augment the training data available for other languages (whose corpora are much smaller than the PDTB), we translate the training, development and test sets of the English PDTB into Italian, Portuguese and Turkish using the Google Translate API<sup>4</sup>. In the following sections, the train, dev and test split from Jiang et al. (2023) are referred to as pdtb2, while their machine-translated versions are referred to as pdtb-it, pdtb-pt and pdtb-tr.

**Italian discourse corpus: LUNA** The LUNA corpus contains “Italian spontaneous speech recorded in the help-desk facility of the Consortium for Information Systems of Piedmont Regio”

(Tonelli et al., 2010, pp.2084), thus represents Italian, and transcribed (but originally spoken) IT help-desk dialogs. LUNA contains 1,188 relation instances in total (train, dev and test).

**Portuguese discourse corpus: CRPC** The CRPC corpus from Mendes and Lejeune (2022) contains a written subset of the Reference Corpus of Contemporary Portuguese, which in turn aims to serve as a representative sample for the Portuguese language and contains texts from many sources (literature, newspapers, magazines, science, economics, law, parliamentary debates, technical and didactic texts, pamphlets) (Généreux et al., 2012, pp.2237). CRPC contains 6,274 relation instances in total.

**Turkish discourse corpus: TDB** The Turkish Discourse Bank (TDB) corpus (Zeyrek and Kurfali, 2017) contains written Turkish texts from a variety of genres (novels, stories, research surveys, travel and news articles, interviews and memoirs). TDB contains 1,809 relation instances in total.

The combination of authentic and synthetic (i.e., machine-translated) corpora enables us to experiment with different set-ups, using in-domain, out-of-domain, in-language and out-of-language configurations for training and test sets, to see how well the model generalizes across the different dimensions. Statistics of our data sets are included in Table 1. Since pdtb-it, pdtb-pt and pdtb-tr are direct translations of the relations in our pdtb2 corpus, the number of instances in those data sets are identical to the pdtb2. One of the key goals of this paper is to find out how our relation classifier generalizes across both languages and domains. While “language” is comparatively well-defined (with Turkish being a different language than Portuguese, for example), the notion of “domain” is less clear-cut. The LUNA corpus stands out in that it represents spontaneous speech in help-desk context, but both CRPC and TDB are multi-genre, include news texts and therefore could be considered not that different from the (financial news) PDTB texts. In our experiments though, we assume each of the three non-English corpora to be of a different domain than the PDTB, and get back to the discussion of domain differences in Section 6.

<sup>2</sup>16 in the 2.0 version of the hierarchy (Prasad et al., 2008).

<sup>3</sup>[https://github.com/YJiangcm/GOLF\\_for\\_IDRR/blob/master/preprocess.py](https://github.com/YJiangcm/GOLF_for_IDRR/blob/master/preprocess.py)

<sup>4</sup>Translations were obtained on February 28, 2024.

	pdtb2	LUNA	CRPC	TDB
train	12,547	728	4,869	1,348
dev	1,165	168	769	193
test	1,039	292	636	268
total	14,751	1,188	6,274	1,809

Table 1: Data statistics.

## 4 Method

**Discourse relation classification model** Our model is based on Jiang et al. (2023), who improve over prior work on implicit discourse relation classification by proposing a hierarchy-aware architecture, that takes into the account the global and local level of PDTB relation senses. We use the default hyper-parameter settings of Jiang et al. (2023), except for the number of epochs, which we set to 30. We use XLM-RoBERTa-base (Conneau et al., 2019) for most configurations, since we want to investigate the potential for generalisation across languages. For experiments where training and test data is from the same language, we use roberta-base (Liu et al., 2019) for English, roberta-base-italian<sup>5</sup> for Italian, portuguese-roberta-base<sup>6</sup> for Portuguese, and roberta-base-turkish-uncased (Aytan and Sakar, 2022) for Turkish.

**State-of-the-art models for Italian, Portuguese and Turkish** We compare performance of our setup on the Italian and Portuguese corpora to that of Liu et al. (2023), and on the Turkish corpus to that of Gessler et al. (2021) (see Section 2 for details). Recall that while Jiang et al. (2023) work with implicit relations only from the PDTB, because the data featured in the 2021 and 2023 shared tasks (Zeldes et al., 2021; Braud et al., 2023) combines all relation types, for LUNA, CRPC and TDB, we train and evaluate on both implicits, explicits and other relation types.

**Domain adaptation** In addition to trying out different data configurations (of training and test data), we experiment with marking out-of-domain training samples at training time. This is inspired by Daumé III (2007); Kim et al. (2016), who augment the feature space that is used as input to the classifier model, thereby forcing the learning algorithm

<sup>5</sup><https://huggingface.co/osiria/roberta-base-italian>

<sup>6</sup><https://huggingface.co/flax-community/portuguese-roberta-base>

to do the adaptation. In their implementation, the dimension *domain* simply occupies a particular position in the vector representation of the input to the classifier. Similarly, we simply concatenate the final representation with a binary flag, indicating if the training sample is in-domain or out-of-domain. In the original model architecture, the vectorized representations of *arg1* and *arg2* are concatenated and used as input for the classifier. In our experiments with marking of out-of-domain data, we combine this concatenated vector with another vector of zeros if the sample is out-of-domain, and with another vector of ones if the sample is in-domain.

## 5 Results

The following subsections present the results for different base models and different configurations of training and test data.

### 5.1 Mono-lingual vs. Multi-lingual Model

We first want to test how much model performance degrades by switching to a multi-lingual instead of a mono-lingual base model. We therefore reran the original model from Jiang et al. (2023), and compare it to a version in which we replace the mono-lingual English RoBERTa-base model by the multi-lingual XLM-RoBERTa-base model. Recall that 4-way and 11-way results correspond to classification on the top and second level, respectively, of the PDTB sense hierarchy. Table 2 shows that our replication of Jiang et al. (2023) yielded slightly lower (but roughly comparable) results, but that we see a sharp drop in performance: 10 points in both accuracy and  $f_1$ <sup>7</sup> when exchanging the English model for the multi-lingual one. When working with English data, using a mono-lingual English model thus yields better results.

model	4-way $f_1$ (acc.)	11-way $f_1$ (acc.)
JZW23-orig	65.76 (72.52)	41.74 (61.16)
JZW23-reprod	64.07 (71.61)	39.63 (60.35)
XLM-R-base	54.57 (62.95)	30.61 (48.99)

Table 2: Results for a mono-lingual and multi-lingual base model on English data (pdtb2). JZW23 stands for Jiang et al. (2023); orig refers to reported numbers in their table 1; reprod refers to our results from running their code; XLM-R-base stands for the XLM-RoBERTa-base model.

<sup>7</sup>All  $f_1$ -scores in this paper are macro-averaged.

## 5.2 Language Transfer (Zero-resource setting)

Next, we consider a setting in which no data is available in the target language and compare how well the mono-lingual model using translated English corpus data for training does compared to a setting where a multi-lingual model is trained on the English corpus and then applied to Italian/Portuguese/Turkish data (pdtb-it, pdtb-pt, pdtb-tr). We test both on the pdtb2 test set translated into the target language as well as on the test set of data that was originally annotated in the target language. Note that in the latter case, the model has to deal with both a language transfer problem and with a domain-adaptation problem, as the original corpora contain data from different domains than the English PDTB corpus. Our experiments reported in this section use the multi-lingual XLM-RoBERTa-base model; we will get back to a comparison to mono-lingual models for Italian, Portuguese and Turkish in Section 5.4 below.

Table 3 illustrates that performance on the translated pdtb2 dataset to Italian, Portuguese and Turkish remains relatively stable compared to the performance of the multi-lingual model on English (compare  $f_1$  and accuracy scores to the last row in Table 2). This suggests that it is possible to learn discourse relations independently of language, and apply these learned representations to another language, which has not been seen during task-specific fine-tuning.

Furthermore, we can see that performance is slightly better when fine-tuning the multi-lingual model on the translated pdtb2 data compared to training it on English and then applying to the target language (compare the first two rows for each language in Table 3).

Finally, we can also observe that there is a substantial drop in performance when evaluating on the test set of the original Italian (LUNA) / Portuguese (CRPC) / Turkish (TDB) data. There might be several reasons for this: The discourse-annotated texts from the other languages are from different domains – hence, the approach not only has to generalize across languages but also across domains, a well-known notoriously difficult problem. Alternatively, it is possible that the translated data is atypical, suffering from translationese effects and thereby might hamper generalization from translated data to native data. Another factor is type of arguments the model has seen during

train	test	4-way $f_1$ (acc.)	11-way $f_1$ (acc.)
pdtb2	pdtb-it	54.79 (64.10)	31.85 (49.66)
pdtb-it	pdtb-it	56.72 (64.39)	33.63 (49.23)
pdtb-it	LUNA	43.06 (48.29)	18.03 (34.59)
pdtb2	pdtb-pt	53.24 (62.95)	31.17 (48.80)
pdtb-pt	pdtb-pt	55.03 (63.75)	31.61 (47.83)
pdtb-pt	CRPC	45.74 (57.08)	17.03 (37.26)
pdtb2	pdtb-tr	51.24 (61.50)	30.41 (46.92)
pdtb-tr	pdtb-tr	51.57 (60.73)	30.31 (45.91)
pdtb-tr	TDB	43.15 (47.01)	18.11 (32.84)

Table 3: Results for the XLM-R-base model on language transfer, testing on synthetic, translated data as well as on originally annotated data in the target language.

training. Since we train on implicit relations from the PDTB2, the model has only seen examples of inter-sentential relations. In addition to the implicit vs. explicit distinction, it is also confronted with intra-sentential (explicit) relations in the test set-up. Finally, it is also possible that Italian / Portuguese / Turkish annotators took different decisions in discourse annotation compared to English annotators on PDTB, leading to a discrepancy in label usage, e.g., by using a smaller set of labels.

## 5.3 Domain Transfer (Low Resource Setting)

Next, we consider a setting where some target language discourse-annotated data is available for training. Our main questions are (a) how our basic setup based on the XLM-RoBERTa-base model compares to the previous state-of-the-art on the Italian, Portuguese and Turkish datasets; (b) whether performance can be improved by exploiting translated data from English; (c) whether our implementation of a domain-adaptation technique inspired by Daumé III (2007); Kim et al. (2016) helps in dealing with the domain gap between translated pdtb2 data and the target domain.

Regarding our first question, we compare our results to the best-performing systems of the 2021 and 2023 shared task iterations (Zeldes et al., 2021; Braud et al., 2023). The results are shown in Table 4. For LUNA, Liu et al. (2023) outperform our baseline setup, while for CRPC and TDB, our baseline outperforms the results of Liu et al. (2023) and Gessler et al. (2021), respectively.

Regarding our second question, we explore training on both the translated pdtb2 data and the train-

model	4-way $f_1$ (acc.)	11-way $f_1$ (acc.)
LFS23	- (65.00)	- (-)
XLM-R on LUNA	64.31 (63.70)	32.81 (52.74)
+ pdtb-it	<b>67.64 (66.78)</b>	<b>41.57 (57.53)</b>
+ pdtb-it + DA	<b>72.72 (71.92)</b>	<b>57.13 (62.33)</b>
LFS23	- (78.53)	- (-)
XLM-R on CPRC	78.71 (83.18)	76.39 (82.55)
+ pdtb-pt	<b>79.86 (83.96)</b>	<b>73.45 (83.02)</b>
+ pdtb-pt + DA	<b>79.86 (83.49)</b>	<b>77.90 (83.65)</b>
GBLPZZ21	- (60.09)	- (-)
XLM-R on TDB	61.80 (64.55)	51.11 (59.33)
+ pdtb-tr	<b>64.45 (68.66)</b>	<b>40.05 (61.57)</b>
+ pdtb-tr + DA	63.98 (67.54)	53.92 (64.18)

Table 4: Baseline performance on Italian (LUNA), Portuguese (CRPC) and Turkish (TDB), compared to prior work; LFS23 stands for (Liu et al., 2023); GBLPZZ21 stands for (Gessler et al., 2021); DA stands for domain adaptation (adding a flag that indicates what domain each data point comes from).

ing data from the target domain, in a setup where the model is first trained on the translated pdtb2 data for 15 epochs, and then on the training section of the target-language original data for 15 more epochs. Our results (see the magenta rows in Table 4) show that using translated English data from the financial news domain as additional training data increases performance consistently for all three languages in the 4-way (top-level) classification task; however, we also observe a drop in performance on the 11-way (second-level) classification task. A more detailed analysis indicates that this might be due to different distributions of second-level labels between the corpora (we will get back to this in Section 6), hence second-level labels suffer more severely from the domain shift between PDTB and other domains.

**Simple domain adaptation** Finally, regarding our third question, table 4 also presents the results for explicitly marking the out-of-domain data (here: the translated pdtb2 texts) at training time. The rows marked “+DA” represent configurations where the training data is from multiple domains, but from the same language. We find that domain marking leads to improved performance in most settings, with strongest improvements obtained on the 11-way classification problems. This indicates that the distribution shift regarding second-level labels can be modelled successfully by including

the domain flag. We note that our proposed method including translated English data and the simple domain adaptation technique outperform the previous state-of-the-art results consistently and by a substantial margin on all three languages. It should be noted here though that we use additional training data which was not available in the shared task, and that for a direct comparison, the winning system of the shared task should be trained with this additional data as well.

We also tested a configuration where the multi-lingual model is first fine-tuned on English pdtb2 data for 15 epochs (without translating that data, but in a setting that does use the domain adaptation flag), and then further fine-tuned on the target language training data. We found that this setting leads to worse results than using translated data for Italian (3 point drop) and Portuguese (1 point drop), whereas for Turkish, better results are obtained when using original, English pdtb2 data, combined with Turkish in-domain data ( $f_1$  66.55, acc 69.78 on 4-way classification, and  $f_1$  55.00, acc 64.93 on 11-way classification).

#### 5.4 Multi-lingual vs. Mono-lingual Target Language Models

Because our experiments on English with a mono-lingual vs. a multi-lingual base model indicated a significant drop in performance moving from a mono-lingual to a multi-lingual model (see Table 2), we also used dedicated mono-lingual models (see Section 4), with translated data (pdtb-it, pdtb-pt, pdtb-tr) in combination with out-of-domain marking in an attempt to further improve performance. We found, however, that unlike the English setting, this did not improve performance, compared to using the multi-lingual model. For LUNA, 4-way  $f_1$  and accuracy dropped from 72.72, 71.92 to 70.24, 69.86, respectively. 11-way  $f_1$  dropped from 57.13 to 50.14, with accuracy staying at 62.33. For CRPC, 4-way  $f_1$  and accuracy dropped significantly from 79.86 and 83.96 to 53.68 and 63.68, respectively. 11-way  $f_1$  and accuracy dropped significantly as well, from 77.90, 83.65 to 40.15, 62.58. For TDB, the performance drop was equally significant. 4-way  $f_1$  and accuracy dropped from 66.55, 69.78 to 55.02, 59.70. 11-way  $f_1$  and accuracy dropped from 55.00, 64.92 to 37.96, 54.48.

## 6 Discussion

Overall, we obtain the best results by combining data from different domains, and marking the out-of-domain instances at training time. For Italian and Portuguese, using in-language training data yields better results, whereas for Turkish, combining English out-of-domain data with Turkish in-domain data yields better results. According to [Conneau et al. \(2019, Appendix A\)](#), the training data for XLM-RoBERTa-base included 20.9 GiB of Turkish, compared to 30.2 GiB for Italian, 49.1 GiB for Portuguese, and 300.8 GiB for English. This could explain the better performance when using original, English data, since the model has seen comparatively few Turkish at pre-training. However, since the difference between the amount of training data in Turkish, Italian and Portuguese is not that large, we consider more research necessary to draw conclusions on this. For Italian and Portuguese, it seems that the automatically obtained, synthetic data is good enough to improve classification performance for discourse relation classification when testing on authentic data.

The performance difference for the three corpora overall are rather large, but [Table 1](#) indicates that performance does not correspond to the size of the corpus. Although CRPC is the largest and has the highest scores overall, LUNA has higher scores than TDB, despite LUNA being about 1.5 times smaller than TDB.

**Differences in label distributions** [Figure 1](#) displays the distribution of top-level senses for the four corpora used in our experiments.<sup>8</sup>

From this, we can see that LUNA has a more balanced distribution than the others, possibly explaining its comparatively high scores (for our best-performing set-up), taking into account that it is by far the smallest corpus. Note that pdtb2 has a fairly imbalanced distribution, with a large proportion of expansion relations and relatively few temporal relations. This distribution can be partially attributed to specificities of the newspaper genre, and partially to the fact that for the pdtb2 corpus, we are only working with implicit relations. Temporal relations are often expressed explicitly, which may contribute to their low rate in pdtb2. We also note that CRPC is most similar to pdtb2 in the context of the number of implicit relations in the

<sup>8</sup>Recall that pdtb-it, pdtb-pt and pdtb-tr have the exact same distributions as the pdtb2.

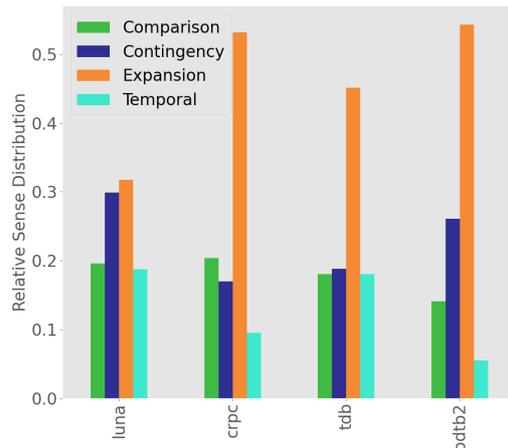


Figure 1: Top-level sense distributions.

dataset, as it has the highest implicit ratio of the three non-English corpora: [Table 6 of Braud et al. \(2023\)](#) shows a ratio of 0.58 (711 implicit relations, divided by 1,228 total relations), compared to 0.47 and 0.32 for TDB and LUNA, respectively.

Turning to second-level senses in [Figure 2](#), we see that CRPC and TDB have slightly fewer unique senses (6 for both) than the other two (9 for LUNA, 11 for pdtb2). While having fewer classes often results in higher scores for multi-class classification, this does not seem to be very predictive for our task, as CRPC has a relatively high  $f_1$ -score (77.90), but TDB (with only 6 unique senses) scores 53.92, while LUNA scores 57.13 (with 9 unique senses). In this respect, it is important to point out that there might be corpus-specific biases: For all three systems submitted to the 2023 DISRPT shared task, CRPC relation classification performs significantly above the corresponding mean of the system, and two of the three systems shows second-best results on this corpus (after the Thai corpus) ([Braud et al., 2023, Table 5](#)). This might indicate that the CRPC corpus contains particularly *easy* relations, and we consider an investigation of what *easy* means in this context and important direction for future work.

Our results also showed that including pdtb2 data was detrimental to performance for 11-way classification when no domain flagging is used, and we speculated that this could be due to strong differences in the distributions of second-level senses. In [Figure 2](#), we can indeed observe that in the expansion class, pdtb2 has many more instantiation and restatement relations, and fewer conjunctions than the other corpora.

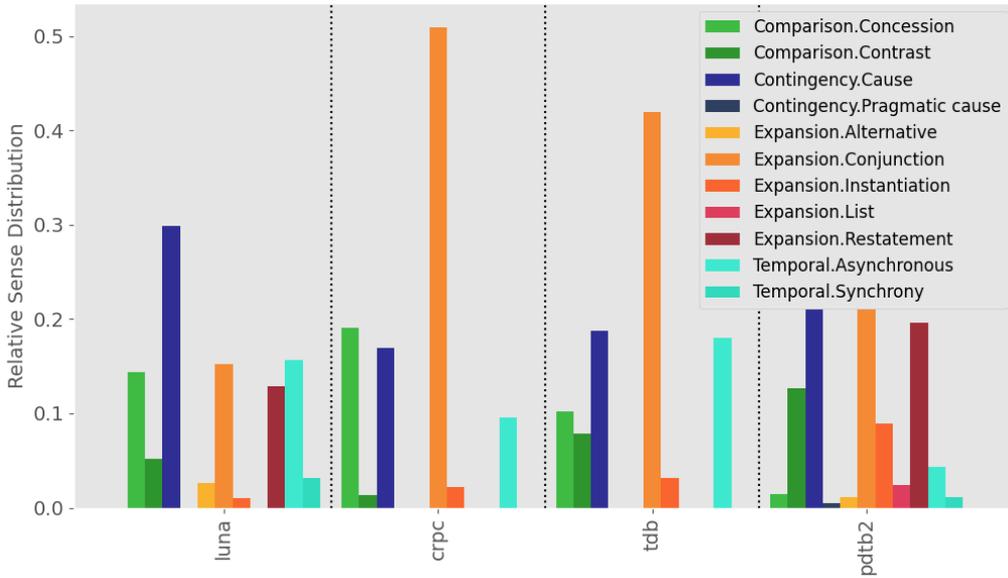


Figure 2: Second-level sense distributions.

**Domain differences between corpora** In an attempt to assess to what extent the corpora used in our experiments differ with respect to the actual words and phrases used, we include Figure 3. This is the result of a pair-wise compari-

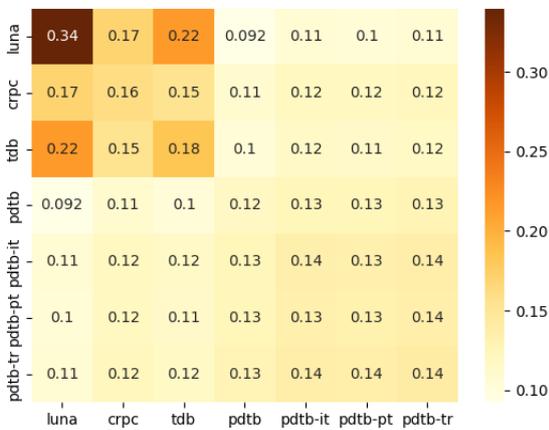


Figure 3: Corpus similarity heatmap.

son of all corpora, where the number expresses the average cosine similarity between all relational arguments in the corpora. To calculate cosine similarity, we encode the arguments using `stsb-xtlm-r-multilingual` from sentence-transformers (Reimers and Gurevych, 2019). The sentence-transformers architecture is specifically designed to express semantic similarity at sentence level, and by using a multi-lingual model, the as-

sumption is that a particular sentence in English shows a high degree of similarity with the (maximally faithful) translation of that sentence in, for example, Turkish. The numbers on the diagonal in Figure 3 express how diverse a corpus is: The high number for LUNA hence indicates that there is relatively little diversity in the LUNA corpus (0.34) compared to e.g., the pdtb2 corpus. The low number for LUNA vs. pdtb2 (0.092) in Figure 3 indicates that the relational arguments in LUNA tend to be very different from the relational arguments of pdtb2. From this figure, we can read that both LUNA and TDB stand out in their usage of words and phrases, as they display a higher average cosine similarity when compared intra-corpus than when compared inter-corpus. While we indeed see a significant drop when training on pdtb2 and testing on LUNA and TDB, the same drop is observed when training on pdtb2 and testing on CRPC, although pdtb2 and CRPC display a considerably lower divergence compared to pdtb2 and LUNA, and pdtb2 and TDB.

Another possible explanation of performance could be the single- or multi-domain aspect of an evaluation corpus. We compared the performance of the winning system of the 2023 shared task<sup>9</sup> (Liu et al., 2023) along this axis, and this reveals that the average performance on single-domain cor-

<sup>9</sup>We use this, and not our system, to have more data points, as we do not have results for the Thai and Chinese corpora.

pora (65.00 for LUNA, 74.30 for the PDTB, 64.96 for TEDM, 95.83 for TDTB and 59.63 for CDTB (Braud et al., 2023, Table 5)) is higher than on multi-domain corpora (78.53 for CRPC and 45.50 for TDB): 71.95 vs. 62.02 for single- vs. multi-domain, respectively. This could be because a single-domain corpus is likely to be more consistent in terms of the types of discourse relations that occur in it. This observation, however, is only based on two data points for multi-domain corpora, and while an interesting direction, we consider more data points necessary before such conclusions can be drawn.

## 7 Conclusion

In this paper, we adopt a state-of-the-art implicit discourse relation classification model developed for English, and apply it to both implicit and non-implicit discourse relations from three corpora that differ in language and domain: An Italian corpus of transcribed IT Helpdesk dialogs, a multi-domain Portuguese corpus and a multi-domain Turkish corpus. By experimenting with different configurations of in-domain, out-of-domain, in-language and out-of-language training data, we explore to what extent the model generalizes across languages and domains. We also demonstrate the importance of using a flag to mark out-of-domain data at training time. Overall, our setup improves over prior work by just under 7 points for Italian, over 5 points for Portuguese, and over 9 points for Turkish (all based on 4-way classification accuracy scores). Our code is published on GitHub.<sup>10</sup>

We attempt to link the classification results to the number of training samples, label distribution and language usage. We find that the number of training samples or sentence similarity between training and test domain is not very indicative of performance, and that instead the label distribution is likely to be a more reliable indicator. In future work, we plan to delve deeper into specific label distributions of the different domains, and potentially continue this line of work by not just looking at different domains, but by also including and testing on annotated data using wholly different label sets (e.g., Rhetorical Structure Theory (Mann and Thompson, 1988) corpora).

In order to maximize comparability to related work, we tested on implicit relations only in pdtb

set-ups, while we combine different relation types in the other set-ups. An interesting direction of future work would be to more systematically investigate the significance of the strict distinction between explicit and implicit relations (as it is often found in the literature), given current, state-of-the-art models for discourse relation classification.

## Limitations

Our experiments rely on fine-tuning on LLMs, and benefit greatly from running on a GPU. Reproduction of our results without having access to a GPU will therefore be time-consuming. Furthermore, although XLM-RoBERTa is specifically targeted at multi-lingual use cases, the amount of training data varies per language (Conneau et al., 2019, Appendix A). For languages with relatively few GiBs of training data, performance may be significantly lower than for the languages we included in our evaluation.

## Ethics Statement

Since our method relies on XLM-RoBERTa for the encoding of input, it will propagate any biases present in (the training data of) this pre-trained language model.

## Acknowledgements

We thank the three anonymous reviewers for their insightful comments. This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID 232722074 – SFB 1102.

## References

- Kaveri Anuranjana. 2023. *DiscoFlan: Instruction fine-tuning and refined text generation for discourse relation label classification*. In *Proceedings of the 3rd Shared Task on Discourse Relation Parsing and Treebanking (DISRPT 2023)*, pages 22–28, Toronto, Canada. The Association for Computational Linguistics.
- Burak Aytan and C Okan Sakar. 2022. Comparison of transformer-based models trained in turkish and different languages on turkish natural language processing problems. In *2022 30th Signal Processing and Communications Applications Conference (SIU)*, pages 1–4. IEEE.

Peter Bourgonje. 2021. *Shallow Discourse Parsing for German*. Doctoral thesis, Universität Potsdam.

<sup>10</sup>[https://github.com/PeterBourgonje/GOLF\\_multilingual](https://github.com/PeterBourgonje/GOLF_multilingual)

- Chloé Braud, Yang Janet Liu, Eleni Metheniti, Philippe Muller, Laura Rivière, Attapol Rutherford, and Amir Zeldes. 2023. [The DISRPT 2023 Shared Task on Elementary Discourse Unit Segmentation, Connective Detection, and Relation Classification](#). In *Proceedings of the 3rd Shared Task on Discourse Relation Parsing and Treebanking (DISRPT 2023)*, pages 1–21, Toronto, Canada. The Association for Computational Linguistics.
- Chunkit Chan, Cheng Jiayang, Weiqi Wang, Yuxin Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song. 2024. [Exploring the potential of ChatGPT on sentence level relations: A focus on temporal, causal, and discourse relations](#). In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 684–721, 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. 2022. [Scaling instruction-finetuned language models](#). *Preprint*, arXiv:2210.11416.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Unsupervised cross-lingual representation learning at scale](#). *CoRR*, abs/1911.02116.
- Hal Daumé III. 2007. [Frustratingly easy domain adaptation](#). In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 256–263, Prague, Czech Republic. 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.
- Yue Dong, Andrei Mircea, and Jackie Chi Kit Cheung. 2021. [Discourse-aware unsupervised summarization for long scientific documents](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1089–1102, Online. Association for Computational Linguistics.
- Michel Génèreux, Iris Hendrickx, and Amália Mendes. 2012. [Introducing the reference corpus of contemporary Portuguese online](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 2237–2244, Istanbul, Turkey. European Language Resources Association (ELRA).
- Luke Gessler, Shabnam Behzad, Yang Janet Liu, Siyao Peng, Yilun Zhu, and Amir Zeldes. 2021. [DisCoDisCo at the DISRPT2021 shared task: A system for discourse segmentation, classification, and connective detection](#). In *Proceedings of the 2nd Shared Task on Discourse Relation Parsing and Treebanking (DISRPT 2021)*, pages 51–62, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yuxin Jiang, Linhan Zhang, and Wei Wang. 2023. [Global and local hierarchy-aware contrastive framework for implicit discourse relation recognition](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8048–8064, Toronto, Canada. Association for Computational Linguistics.
- Najoung Kim, Song Feng, Chulaka Gunasekara, and Luis Lastras. 2020. [Implicit discourse relation classification: We need to talk about evaluation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5404–5414, Online. Association for Computational Linguistics.
- Young-Bum Kim, Karl Stratos, and Ruhi Sarikaya. 2016. [Frustratingly easy neural domain adaptation](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 387–396, Osaka, Japan. The COLING 2016 Organizing Committee.
- Christian Kirschner, Judith Eckle-Kohler, and Iryna Gurevych. 2015. [Linking the thoughts: Analysis of argumentation structures in scientific publications](#). In *Proceedings of the 2nd Workshop on Argumentation Mining*, pages 1–11, Denver, CO. Association for Computational Linguistics.
- Yudai Kishimoto, Yugo Murawaki, and Sadao Kurohashi. 2018. [A knowledge-augmented neural network model for implicit discourse relation classification](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 584–595, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- René Knaebel. 2021. [discopy: A neural system for shallow discourse parsing](#). In *Proceedings of the 2nd Workshop on Computational Approaches to Discourse*, pages 128–133, Punta Cana, Dominican Republic and Online. Association for Computational Linguistics.
- Murathan Kurfalı and Robert Östling. 2019. [Zero-shot transfer for implicit discourse relation classification](#). In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 226–231, Stockholm, Sweden. Association for Computational Linguistics.

- Ziheng Lin, Hwee Tou Ng, and Min-Yen Kan. 2014. A PDTB-Styled End-to-End Discourse Parser. *Natural Language Engineering*, 20:151–184.
- Wei Liu, Yi Fan, and Michael Strube. 2023. **HITS at DISRPT 2023: Discourse Segmentation, Connective Detection, and Relation Classification**. In *Proceedings of the 3rd Shared Task on Discourse Relation Parsing and Treebanking (DISRPT 2023)*, pages 43–49, Toronto, Canada. The Association for Computational Linguistics.
- Xin Liu, Jiefu Ou, Yangqiu Song, and Xin Jiang. 2020. On the importance of word and sentence representation learning in implicit discourse relation classification. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 3830–3836.
- Yang Liu, Sujian Li, Xiaodong Zhang, and Zhifang Sui. 2016. Implicit discourse relation classification via multi-task neural networks. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, AAAI’16*, page 2750–2756. AAAI Press.
- 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**. *CoRR*, abs/1907.11692.
- William Mann and Sandra Thompson. 1988. Rhetorical Structure Theory: Towards a functional theory of text organization. *TEXT*, 8:243–281.
- Amália Mendes and Pierre Lejeune. 2022. **Crpc-db a discourse bank for portuguese**. In *Computational Processing of the Portuguese Language: 15th International Conference, PROPOR 2022, Fortaleza, Brazil, March 21–23, 2022, Proceedings*, page 79–89, Berlin, Heidelberg. Springer-Verlag.
- Eleni Metheniti, Chloé Braud, Philippe Muller, and Laura Rivière. 2023. **DisCut and DiscReT: MELODI at DISRPT 2023**. In *Proceedings of the 3rd Shared Task on Discourse Relation Parsing and Treebanking (DISRPT 2023)*, pages 29–42, Toronto, Canada. The Association for Computational Linguistics.
- Stephan Oepen, Jonathon Read, Tatjana Schefler, Uladzimir Sidarenka, Manfred Stede, Erik Velldal, and Lilja Øvrelid. 2016. OPT: Oslo–Potsdam–Teesside—Pipelining Rules, Rankers, and Classifier Ensembles for Shallow Discourse Parsing. In *Proceedings of the 20th Conference on Computational Natural Language Learning: Shared Task (CoNLL Shared Task 2016)*, pages 20–26, Berlin.
- Emily Pitler and Ani Nenkova. 2009. **Using syntax to disambiguate explicit discourse connectives in text**. In *Proceedings of the ACL-IJCNLP 2009 Conference Short Papers*, pages 13–16, Suntec, Singapore. Association for Computational Linguistics.
- Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Miltasakaki, Livio Robaldo, Aravind Joshi, and Bonnie Webber. 2008. **The Penn discourse TreeBank 2.0**. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC’08)*, pages 2961–2968, Marrakech, Morocco. European Language Resources Association (ELRA).
- 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.
- Attapol Rutherford and Nianwen Xue. 2016. **Robust non-explicit neural discourse parser in English and Chinese**. In *Proceedings of the CoNLL-16 shared task*, pages 55–59, Berlin, Germany. Association for Computational Linguistics.
- Wei Shi and Vera Demberg. 2019. **Next sentence prediction helps implicit discourse relation classification within and across domains**. 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 5790–5796, Hong Kong, China. Association for Computational Linguistics.
- Jialong Tang, Hongyu Lin, Meng Liao, Yaojie Lu, Xi-anpei Han, Le Sun, Weijian Xie, and Jin Xu. 2021. **From discourse to narrative: Knowledge projection for event relation extraction**. 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 732–742, Online. Association for Computational Linguistics.
- Sara Tonelli, Giuseppe Riccardi, Rashmi Prasad, and Aravind Joshi. 2010. **Annotation of discourse relations for conversational spoken dialogs**. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC’10)*, Valletta, Malta. European Language Resources Association (ELRA).
- Bonnie Webber, Rashmi Prasad, Alan Lee, and Aravind Joshi. 2019. **The Penn Discourse Treebank 3.0 Annotation Manual**.
- 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.
- Jiacheng Xu, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. **Discourse-aware neural extractive text summarization**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5021–5031, 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.
- Frances Yung, Mansoor Ahmad, Merel Scholman, and Vera Demberg. 2024. [Prompting implicit discourse relation annotation](#). In *Proceedings of The 18th Linguistic Annotation Workshop (LAW-XVIII)*, pages 150–165, St. Julians, Malta. Association for Computational Linguistics.
- Amir Zeldes, Debopam Das, Erick Galani Maziero, Juliano Antonio, and Mikel Iruskieta. 2019. [The DISRPT 2019 Shared Task on Elementary Discourse Unit Segmentation and Connective Detection](#). In *Proceedings of the Workshop on Discourse Relation Parsing and Treebanking 2019*, pages 97–104, Minneapolis, MN. Association for Computational Linguistics.
- Amir Zeldes, Yang Janet Liu, Mikel Iruskieta, Philippe Muller, Chloé Braud, and Sonia Badene. 2021. [The DISRPT 2021 Shared Task on Elementary Discourse Unit Segmentation, Connective Detection, and Relation Classification](#). In *Proceedings of the 2nd Shared Task on Discourse Relation Parsing and Treebanking (DISRPT 2021)*, pages 1–12, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Deniz Zeyrek and Murathan Kurfalı. 2017. [TDB 1.1: Extensions on Turkish discourse bank](#). In *Proceedings of the 11th Linguistic Annotation Workshop*, pages 76–81, Valencia, Spain. Association for Computational Linguistics.

# BoK: Introducing Bag-of-Keywords Loss for Interpretable Dialogue Response Generation

Suvodip Dey and Maunendra Sankar Desarkar  
Indian Institute of Technology Hyderabad, India  
suvodip15@gmail.com, maunendra@cse.iith.ac.in

## Abstract

The standard language modeling (LM) loss by itself has been shown to be inadequate for effective dialogue modeling. As a result, various training approaches, such as auxiliary loss functions and leveraging human feedback, are being adopted to enrich open-domain dialogue systems. One such auxiliary loss function is Bag-of-Words (BoW) loss, defined as the cross-entropy loss for predicting all the words/tokens of the next utterance. In this work, we propose a novel auxiliary loss named Bag-of-Keywords (BoK) loss to capture the central thought of the response through keyword prediction and leverage it to enhance the generation of meaningful and interpretable responses in open-domain dialogue systems. BoK loss upgrades the BoW loss by predicting only the keywords or critical words/tokens of the next utterance, intending to estimate the core idea rather than the entire response. We incorporate BoK loss in both encoder-decoder (T5) and decoder-only (DialoGPT) architecture and train the models to minimize the weighted sum of BoK and LM (BoK-LM) loss. We perform our experiments on two popular open-domain dialogue datasets, DailyDialog and Persona-Chat. We show that the inclusion of BoK loss improves the dialogue generation of backbone models while also enabling post-hoc interpretability. We also study the effectiveness of BoK-LM loss as a reference-free metric and observe comparable performance to the state-of-the-art metrics on various dialogue evaluation datasets.

## 1 Introduction

Open-domain dialogue generation is a dynamic area of research, aiming to generate contextually relevant and meaningful responses given a dialogue context. As deep learning models continue to thrive in the field of natural language processing (NLP), a widely adopted strategy to solve any natural language generation (NLG) task involves pre-training and fine-tuning large language models (LLMs).

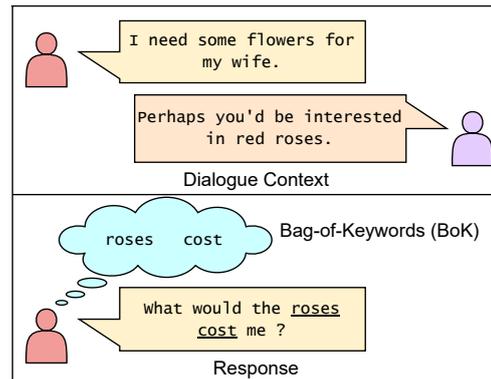


Figure 1: A motivating example for Bag-of-Keywords loss in open-domain dialogue system.

The LLMs are predominantly trained with language modeling (LM) loss, which essentially corresponds to cross-entropy loss for predicting the next word or token. While LM loss remains effective in training NLG models for diverse tasks, including dialogue generation (Sordoni et al., 2015; Wolf et al., 2019; Zhang et al., 2020; Roller et al., 2021), it may not be the optimal choice for training models specifically tailored for dialogue generation. It is well-established that perplexity, a measure associated with LM loss, primarily gauges fluency and weakly correlates with human dialogue evaluation (Dinan et al., 2019; Mehri and Eskenazi, 2020b; Phy et al., 2020). Consequently, relying solely on LM loss may not guarantee generations with desirable conversational qualities. Therefore, exploring alternative loss functions and training methods is crucial to advance the development of generative open-domain dialogue models.

In order to mitigate the exclusive dependence on LM loss in the training of open-domain dialogue models, various approaches have been explored in the existing literature. These techniques can be broadly categorized into two classes - a) auxiliary loss and b) human feedback. The first approach combines one or more auxiliary losses

with LM loss to train the dialogue models. Various types of auxiliary losses have been explored in the context of open-domain dialogue learning. For instance, Bag-of-Words (BoW) loss computes the cross-entropy loss to predict words/tokens of the next utterance from the given dialogue context (Zhao et al., 2017; Li et al., 2021; Dey et al., 2023). Some methodologies involve predicting the sentence-level encoding of the next utterance and determining the loss through L1/L2 norms and KL divergence (Serban et al., 2017; Li et al., 2021; Chen et al., 2022; Dey et al., 2023). Few approaches incorporate a next-utterance classification loss (Wolf et al., 2019), wherein the auxiliary loss is computed for a classification or ranking task to predict the true utterance from a set of candidate responses. On the other hand, the second approach is based on refining the pre-trained dialogue model through human feedback. These methods mostly follow the training principle of Reinforcement learning from human feedback (RLHF), where the model is fine-tuned to maximize the reward associated with the generated response using Reinforcement learning. RLHF has gained significant interest recently, particularly with the popularity of models like Chat-GPT (Long and et al., 2022). However, acquiring quality human feedback data is challenging and expensive (Casper and et al., 2023). Furthermore, relying on automated dialogue evaluation metrics as a substitute for human feedback can pose challenges, as they may not strongly correlate with human judgments (Liu et al., 2016; Yeh et al., 2021).

In this work, our objective is to propose a novel auxiliary loss for open-domain dialogue systems. Specifically, we address the limitation of BoW loss by introducing Bag-of-Keywords (BoK) loss, which is defined as the cross-entropy loss to predict the keywords of the next utterance. While training, we extract the keywords of the ground-truth response using YAKE! (Campos et al., 2018, 2020), an unsupervised feature-based keyword extractor. The keywords can be seen as a proxy for the core idea of the response. In a conversation, a reply can be generated in multiple ways. As a result, BoW loss can induce training data bias since it considers all the words/tokens of the ground-truth response for prediction. In contrast, BoK loss focuses on the core idea (as shown in Fig. 1) that alleviates the problem of generalization. The main contributions

of this work are summarized as follows<sup>1</sup>:

- We propose BoK loss, a novel auxiliary loss for open-domain dialogue systems. BoK loss can be easily incorporated into any generative model and trained using a weighted sum of BoK and LM (BoK-LM) loss.
- We show that BoK loss enhances the dialogue generation of backbone models on DailyDialog and Persona-Chat datasets. We note an improvement in the specificity of the generated responses with the inclusion of BoK loss.
- We perform a qualitative analysis of the generated responses and discuss how BoK loss enables post-hoc interpretability.
- We study the effectiveness of BoK-LM loss as a reference-free metric. We observe that it exhibits moderate correlations with human judgments on different evaluation datasets.

## 2 Background and Related Works

Open-domain dialogue generation is a challenging NLG task. Let  $D_{<t} = \{u_1, u_2, \dots, u_{t-1}\}$  be a multi-turn conversation where  $u_j$  represents the utterance at turn  $j$ . Let  $C_t$  be the condition (like persona, document, etc.) other than dialogue history for generating  $u_t$ . The task of open-domain dialogue generation is to generate  $u_t$  given  $D_{<t}$  and  $C_t$ . Like any NLG task, it is modeled using language models and generally trained using the next word/token prediction task. The corresponding language modeling (LM) loss is defined as,

$$\mathcal{L}_{\text{LM}} = - \sum_{n=1}^T \log p(u_{t_n} | u_{t_{<n}}, D_{<t}, C_t; \theta) \quad (1)$$

where  $u_{t_n}$  denotes the  $n^{\text{th}}$  word/token of utterance  $u_t$  and  $\theta$  indicates the parameters of the language model. Training transformer (Vaswani et al., 2017) based large language models (LLMs) with LM loss on large dialogue corpora has shown remarkable performance in open-domain dialogue generation (Zhang et al., 2020; Roller et al., 2021). However, it has been shown that perplexity ( $e^{\mathcal{L}_{\text{LM}}}$ ), a metric that is a function of LM loss, can measure fluency but shows a weak correlation with other conversational aspects (Dinan et al., 2019; Mehri and Eskenazi, 2020b; Phy et al., 2020). The root

<sup>1</sup>Code is available at [github.com/SuvodipDey/BoK](https://github.com/SuvodipDey/BoK)

cause of this behavior stems from the inherent one-to-many nature of dialogue, where a given context can elicit multiple possible responses (Liu et al., 2016). Consequently, simply increasing the size of training data may not always yield improvement, as it is impractical to collect all potential response variations. To tackle this challenge, researchers employ various techniques, broadly categorized into two classes: i) incorporating one or more auxiliary losses alongside LM loss, and ii) leveraging human feedback to finetune pre-trained dialogue models. Given our focus on proposing a new auxiliary loss, we keep our related works limited to different auxiliary losses utilized for open-domain dialogue generation, described as follows.

- The first kind of auxiliary loss estimates the error in predicting the sentence-level encoding of the next utterance given the dialogue context. Authors of VHRED (Serban et al., 2017) and DialogVED (Chen et al., 2022) use Kullback-Leibler (KL) divergence to compute the distance between the approximate and true posterior distribution of the next utterance. Models like DialoFlow (Li et al., 2021) and DialoGen (Dey et al., 2023) use the L1/L2 norm for the same purpose. Predicting the encoding of the next utterance is challenging and may lead to issues like posterior collapse while using KL divergence (Chen et al., 2022).
- The second approach is based on the next utterance classification loss. In this method, the task is to classify the ground-truth response from a given set of candidate utterances (Wolf et al., 2019). It is worth noting that this method requires negative samples, which are usually not included in the datasets. Hence, different kinds of negative sampling techniques are adopted to obtain them. However, finding high-quality negative samples is difficult for dialogues (Lan et al., 2020).
- The third approach focuses on predicting the words/tokens of the next utterance. This loss is popularly known as Bag-of-Words (BoW) loss (Zhao et al., 2017). Models like DialoFlow (Li et al., 2021) and DialoGen (Dey et al., 2023) utilize BoW loss to support LM loss. DialogVED (Chen et al., 2022) uses BoW loss to tackle the posterior collapse that is caused due to minimizing KL divergence. As discussed earlier, a dialogue context can

have many relevant responses. Hence, the task of predicting all the words/tokens can induce training data bias. In this work, we aim to address this limitation of BoW loss.

### 3 Methodology

In this section, we describe Bag-of-Keywords loss followed by its application in open-domain dialogue systems.

#### 3.1 Bag-of-Keywords (BoK) loss

As discussed, BoW loss is computed as the cross-entropy loss to predict all the tokens of the next utterance. Say the model has to generate utterance  $u_t$  given dialogue context  $D_{<t}$ . Let  $\phi_t \in \mathbb{R}^d$  be the representation of the context for generating  $u_t$ . Then the BoW loss ( $\mathcal{L}_{\text{BoW}}$ ) is defined as,

$$\mathcal{L}_{\text{BoW}} = - \sum_{w \in u_t} \log p(w|\phi_t) \quad (2)$$

where  $p(w|\phi_t)$  is the probability of predicting the word/token  $w \in u_t$  given  $\phi_t$ . Predicting all the words of a dialogue response may cause training data bias because there can be multiple ways to generate a response. Additionally, dialogue responses often contain stopwords that are necessary for sentence construction and fluency. Therefore, predicting these stopwords in BoW loss is unnecessary since LM loss already takes care of it.

One simple approach to address this limitation of BoW loss is to predict only the keywords of the response. By keywords, we mean the critical words that capture the core concept of the response. This approach can help reduce the training data bias and increase its generalizability for open-domain dialogue generation. To achieve this, we propose Bag-of-Keywords (BoK) loss, which is computed as the cross-entropy loss to predict the keywords of the next utterance. We define BoK loss ( $\mathcal{L}_{\text{BoK}}$ ) as,

$$\mathcal{L}_{\text{BoK}} = - \sum_{w \in K_t} \log p(w|\phi_t) \quad (3)$$

where  $K_t$  is the set of keywords (or tokens associated with the keywords) in  $u_t$ . Note that the annotations regarding the keywords are not available in the existing dialogue datasets. In this work, we find the keywords using YAKE! (Campos et al., 2018, 2020), an unsupervised feature-based keyword extraction algorithm that leverages statistical features extracted directly from the text, thereby supporting

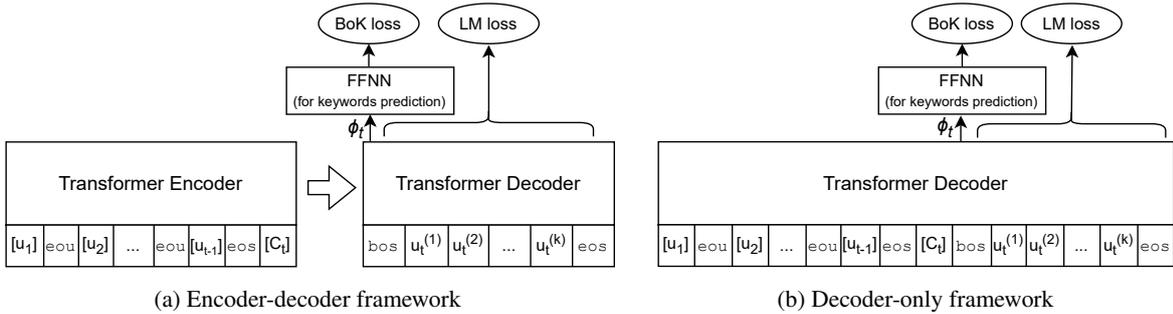


Figure 2: Incorporating BoK loss in open-domain dialogue models.  $[u_j]$  and  $[C_j]$  represents the list of tokens after tokenizing utterance  $u_j$  and condition  $C_j$ , respectively.  $u_t^{(i)}$  denotes the  $i^{th}$  token of utterance  $u_t$ , whereas  $\{eos, bos, eou\}$  are special tokens.  $\phi_t \in \mathbb{R}^d$  is the hidden state of the final layer of bos token, representing the context.

texts of multiple domains and languages. However, one can adopt any strategy for keyword extraction. We chose YAKE! because it is unsupervised and has already been utilized to extract keywords from dialogue responses (Dey and Desarkar, 2023). For example, in Fig. 1, YAKE! extracted the keywords “roses” and “cost” from the response “What would the roses cost me?”.

### 3.2 Application of BoK loss

BoK loss can be easily applied to any open-domain dialogue model. Currently, all state-of-the-art dialogue generation models are based on Transformer (Vaswani et al., 2017). These models can be broadly classified into two architectures - i) encoder-decoder and ii) decoder-only. Incorporating BoK loss into both these architectures is described as follows.

- **Encoder-Decoder Architecture:** Fig. 2a shows the method of applying BoK loss in encoder-decoder architecture. The encoder takes the concatenation of the past utterances ( $D_{<t}$ ) along with the condition  $C_t$  as input. Note that  $C_t$  may be present or absent based on the task or dataset. In the decoder, we add an extra component for computing the BoK loss. Let  $\phi_t \in \mathbb{R}^d$  be the hidden state representation of the final layer corresponding to the bos token, representing the context. Then, the BoK loss is computed as follows:

$$\alpha_t = \text{softmax}(\text{FFNN}(\phi_t)) \in \mathbb{R}^{|V|} \quad (4)$$

$$\mathcal{L}_{\text{BoK}} = - \sum_{w \in K_t} \log p(w|\phi_t) = - \sum_{w \in K_t} \log \alpha_{t_w} \quad (5)$$

where FFNN denotes a single layer feed-forward neural network, and  $|V|$  is the vocabulary size of the decoder tokens.

Dataset	Type	#Dialog	#Turns	$T_{\max}$	$T_{\min}$	$T_{\text{avg}}$
DailyDialog	Train	11118	87170	35	2	7.84
	Dev	1000	8069	31	2	8.07
	Test	1000	7740	26	2	7.74
Persona-Chat	Train	8939	131438	50	12	14.70
	Dev	1000	15602	26	14	15.60
	Test	968	15024	34	14	15.52

Table 1: Basic statistics of DailyDialog and Persona-Chat dataset.  $T_{\max}$ ,  $T_{\min}$ , and  $T_{\text{avg}}$  indicate maximum, minimum, and average dialogue turns.

- **Decoder-only Architecture:** Fig. 2b shows the process of incorporating BoK loss in decoder-only architecture. The BoK loss computation follows the same equations (Eqn. 4 and 5) as encoder-decoder architecture.

The training objective for both architectures is to minimize the weighted sum of BoK and LM loss. We term this loss as BoK-LM loss ( $\mathcal{L}_{\text{BoK-LM}}$ ).

$$\mathcal{L}_{\text{BoK-LM}} = \mathcal{L}_{\text{LM}} + \lambda \mathcal{L}_{\text{BoK}} \quad (6)$$

where  $\lambda \in \mathbb{R}$  is a hyper-parameter to set the weight of the BoK loss. Note that both the loss components depend on the context vector  $\phi_t$ . Hence, the BoK-LM loss helps to learn  $\phi_t$  such that it can capture the core idea of the response and align the generation towards a meaningful response.

## 4 Experimental Set up

### 4.1 Datasets

We conduct our experiments on two datasets: DailyDialog (Li et al., 2017) and Persona-Chat (Zhang et al., 2018a). DailyDialog is a popular chit-chat dataset in which the task is to generate responses conditioned only on the dialogue history. On the other hand, Persona-Chat is a knowledge-grounded dataset where a response needs to be generated

Model	Referenced Metrics										Reference-Free Metric			
	Bleu-1	Bleu-2	Bleu-3	Bleu-4	Nist-2	Nist-4	Meteor	Div-1	Div-2	Entropy	U	S	L <sub>S</sub>	USL <sub>S</sub> -H
DialoFlow	48.75	26.73	16.35	10.70	3.76	3.97	16.44	0.039	0.216	9.98	0.96	0.88	0.21	0.6777
DialoGen	49.13	27.25	16.88	11.07	3.76	3.98	16.40	0.043	0.223	9.88	0.83	0.90	0.32	0.6685
DialogVED	50.50	28.95	18.38	12.29	3.94	4.18	<u>16.90</u>	0.037	0.204	9.82	0.86	0.88	0.30	0.6642
T5	51.56	29.22	18.29	12.05	3.99	4.23	16.27	0.044	0.219	9.62	0.97	0.89	0.18	0.6718
T5 <sub>BoW</sub>	<b>51.75</b>	<u>29.70</u>	18.89	12.75	<u>4.05</u>	4.32	16.64	0.045	0.230	9.79	0.97	0.89	0.19	0.6791
T5 <sub>BoK</sub>	51.74	<b>29.74</b>	<u>19.19</u>	<u>13.24</u>	<b>4.09</b>	<b>4.37</b>	16.62	<u>0.045</u>	<u>0.233</u>	9.84	0.97	0.90	0.20	<u>0.6793</u>
DialoGPT	49.30	27.63	17.37	11.68	3.78	4.01	16.67	0.037	0.193	9.66	0.97	0.89	0.19	0.6731
DialoGPT <sub>BoW</sub>	49.60	27.85	17.60	11.82	3.80	4.04	16.83	0.037	0.190	9.60	0.97	0.89	0.20	0.6759
DialoGPT <sub>BoK</sub>	49.16	29.10	<b>20.00</b>	<b>14.92</b>	4.01	<u>4.35</u>	<b>17.72</b>	<b>0.048</b>	<b>0.257</b>	<b>10.19</b>	0.97	0.89	0.31	<b>0.7064</b>

Table 2: Comparison of dialogue generation performance on DailyDialog test data with automated metrics. The highest and second-highest scores are written in bold and underlined respectively.

based on both dialogue history and a persona profile that defines the speaker. Table 1 displays the basic statistics of the two datasets.

## 4.2 Implementation Details

We choose T5 (Raffel et al., 2020) and DialoGPT (Zhang et al., 2020) as our encoder-decoder and decoder-only architecture, respectively. We perform our experiments with T5-large<sup>1</sup> ( $\approx 770M$  parameters) and DialoGPT-large<sup>2</sup> ( $\approx 774M$  parameters) for both DailyDialog and Persona-Chat datasets. All the implementations are done using PyTorch and Huggingface (Wolf et al., 2020) libraries in Python 3.10, and executed on a Nvidia A100 with 40GB memory. We use AdamW optimizer with a learning rate of  $5e-5$ , batch size of 16, maximum training epochs of 20, and early stopping to train the models. We use beam search with a beam width of 5, maximum sequence length of 40, minimum sequence length of 11, and length penalty of 0.1 to generate responses for all the models. The rest of the details are provided in Appendix A.1.

## 4.3 Baselines

We refer to T5 and DialoGPT trained with BoK-LM loss as T5<sub>BoK</sub> and DialoGPT<sub>BoK</sub>, respectively. We compare them with vanilla T5 and DialoGPT models, trained only with LM loss. To measure the improvement over BoW loss, we also train T5 and DialoGPT with a weighted sum of BoW and LM loss (like Eqn. 6), denoted as T5<sub>BoW</sub> and DialoGPT<sub>BoW</sub> respectively. We also have some dataset-specific baselines. For DailyDialog, we use DialoFlow (Li et al., 2021), DialogVED (Chen et al., 2022), and DialoGen (Dey et al., 2023). All these three baselines use BoW loss and sentence-level next utterance prediction loss. For Persona-Chat, we use TransferTransfo (Wolf et al., 2019) and DialogVED. TransferTransfo utilizes the next

<sup>1</sup>[huggingface.co/google-t5/t5-large](https://huggingface.co/google-t5/t5-large)

<sup>2</sup>[huggingface.co/microsoft/DialoGPT-large](https://huggingface.co/microsoft/DialoGPT-large)

Model	U	S	L <sub>S</sub>	USL <sub>S</sub> -H	Dial-M
TransferTransfo	0.75	0.63	0.44	0.5502	1.7730
DialogVED	0.74	0.84	0.38	<b>0.6348</b>	1.7499
T5	0.71	0.73	0.39	0.5756	0.9288
T5 <sub>BoW</sub>	0.72	0.75	0.40	0.5867	0.8781
T5 <sub>BoK</sub>	0.72	0.76	0.41	0.5947	<b>0.8556</b>
DialoGPT	0.76	0.72	0.36	0.5788	1.0312
DialoGPT <sub>BoW</sub>	0.77	0.71	0.40	0.5868	1.0013
DialoGPT <sub>BoK</sub>	0.77	0.72	0.42	0.5923	1.0004

Table 3: Comparison of dialogue generation performance on Persona-Chat test data.

utterance classification as the auxiliary loss.

## 5 Results and Analysis

### 5.1 DailyDialog Dataset

Table 2 compares the performance of various models on DailyDialog test data. We use BLEU (Papineni et al., 2002), NIST (Lin and Och, 2004), METEOR (Banerjee and Lavie, 2005), Diversity (Li et al., 2016), and Entropy (Zhang et al., 2018b) for referenced evaluation, and USL-H (Phy et al., 2020) for reference-free evaluation. As word-overlapping based metrics are not reliable with only one reference, we conduct the referenced evaluation using multi-reference DailyDialog (Gupta et al., 2019) that contains four additional references along with the original response. For BoK loss, we set the maximum number of keyword tokens  $|K_t| = 8$  (refer Eqn. 3). For BoK-LM loss in Eqn. 6, we set  $\lambda$  to 0.1 and 0.3 for T5 and DialoGPT architecture, respectively. The effect of varying  $\lambda$  and  $|K_t|$  is studied in the ablation study. The key observations from Table 2 are discussed below.

**Referenced Evaluation:** Firstly, we observe that the inclusion of BoW loss enhances the performance of both vanilla T5 and DialoGPT across all metrics. BoW loss is optimized to predict all the words/tokens of the next utterance, thereby improving the unigram match i.e. Bleu-1 score. Our findings corroborate this observation, demonstrating that T5<sub>BoW</sub> and DialoGPT<sub>BoW</sub> attain higher Bleu-1 scores compared to their other counterparts. Sec-

Comparisons	Dataset	Coherence			Engagingness			Informativeness			Interactiveness			Overall		
		W	L	T	W	L	T	W	L	T	W	L	T	W	L	T
T5 <sub>BoK</sub> vs. T5 <sub>BoW</sub>	DailyDialog	24	18	58	30	26	44	20	14	66	26	18	56	32	26	42
	Persona-Chat	26	18	56	24	20	56	24	18	58	20	18	62	28	24	48
DialoGPT <sub>BoK</sub> vs. DialoGPT <sub>BoW</sub>	DailyDialog	42	34	24	30	30	40	44	26	30	34	30	36	46	34	20
	Persona-Chat	28	18	54	14	20	66	24	18	58	14	16	70	28	22	50

Table 4: Human evaluation for comparing the impact of BoK and BoW loss on the performance of the backbone models. “W”, “L”, and “T” denote the percentage of win, loss, and tie, respectively.

only, we note that both T5<sub>BoK</sub> and DialoGPT<sub>BoK</sub> perform better than their BoW counterpart in most of the cases. Furthermore, they also outperform the three baselines (DialoFlow, DialoGen, and DialoVED) that rely on BoW loss. This indicates that BoK loss effectively improves the generalizability of BoW loss, making it more efficient.

**Reference-free Evaluation:** We use USL<sub>S</sub>-H as our reference-free metric, which is a combination of three sub-metrics - Understandability (U), Sensibility (S), and Likability (L). We specifically make use of the USL<sub>S</sub>-H variant, where the likability of a response is captured through its specificity. USL<sub>S</sub>-H estimates understandability, sensibility, and specificity using valid prediction, next utterance prediction, and MLM task, respectively (Phy et al., 2020). Similar to the results of the referenced evaluation, T5<sub>BoK</sub> and DialoGPT<sub>BoK</sub> achieve better USL<sub>S</sub>-H scores than their other counterparts. Moreover, we note that for T5<sub>BoK</sub> and DialoGPT<sub>BoK</sub>, USL<sub>S</sub>-H improves because of the likability or specificity aspect. We also observe this behavior in Table 3, which indicates that incorporating BoK loss enhances the specificity of the generated responses.

## 5.2 Persona-Chat Dataset

The results of the Persona-Chat test data are presented in Table 3. Unlike DailyDialog, Persona-Chat does not have any multi-referenced test data. Therefore, we use only reference-free metrics to ensure a fair evaluation. In addition to USL<sub>S</sub>-H, we also evaluate using Dial-M (Dey and Desarkar, 2023), a masking-based reference-free metric that is effective in evaluating knowledge-grounded dialogues. It is worth mentioning that in Dial-M, a lower score is indicative of better performance as it is based on cross-entropy loss. In Table 3, we again observe that T5<sub>BoK</sub> and DialoGPT<sub>BoK</sub> attain better USL<sub>S</sub>-H and Dial-M scores than their other counterparts. Furthermore, we observe that DialogVED outperforms all the models on USL<sub>S</sub>-H. This is because it does not use persona profiles explicitly and relies on specially trained latent variables (on next utterance prediction) for persona-grounded

response generation. Furthermore, USL<sub>S</sub>-H only considers dialogue history as context and ignores persona. As a result, DialogVED performs better in understandability and sensibility, which are estimated using valid and next utterance prediction tasks, respectively. However, it falls short in specificity and Dial-M as it does not use persona.

We observe that for both DailyDialog and Persona-Chat, BoK performs better than BoW in most of the cases. For DailyDialog, DialoGPT<sub>BoK</sub> outperforms DialoGPT<sub>BoW</sub> significantly, which correlates with the automated result shown in Table 2. For Persona-Chat, as the generation is conditioned mainly on the persona profiles, the responses are very similar for both models, resulting in a lot of ties. We also observe that BoK loss results in better informativeness, which correlates with the improved specificity (in USL<sub>S</sub>-H) shown in Table 2 and Table 3.

## 5.3 Human Evaluation

Table 4 shows the human evaluation to compare the impact of BoK and BoW loss on the backbone models. We randomly picked 50 test instances from both DailyDialog and Persona-Chat datasets. Four human evaluators (graduate students proficient in English) were presented with the generated responses from two models (A and B) who reported their judgment (A wins, B wins, or a tie) on various aspects. We asked the evaluators to evaluate five aspects, described as follows.

- *Coherence*: Captures which model produces more contextually coherent responses.
- *Engagingness*: Identifies which model generates more engaging or interesting responses.
- *Informativeness*: Determines which response contains more knowledge or specific information.
- *Interactiveness*: Captures which model produces more interactive responses that encourage the user to continue the conversation.

Model	$\lambda$	Bleu-1	Bleu-2	Bleu-3	Bleu-4	Nist-2	Nist-4	Meteor	Div-1	Div-2	Entropy	U	S	L <sub>S</sub>	USL <sub>S</sub> -H
T5 <sub>BoK</sub>	0.05	51.61	29.43	18.69	12.63	4.03	4.29	16.47	0.044	0.224	9.74	0.97	0.89	0.19	0.6748
	0.10	<b>51.74</b>	<b>29.74</b>	<b>19.19</b>	<b>13.24</b>	<b>4.09</b>	<b>4.37</b>	<u>16.62</u>	0.045	0.233	9.84	0.97	0.90	0.20	0.6793
	0.20	51.53	29.58	19.06	13.07	4.06	4.34	<b>16.71</b>	0.046	0.231	<u>9.85</u>	0.97	0.90	0.21	0.6802
	0.30	51.08	28.91	18.44	12.55	4.00	4.26	16.58	0.046	<b>0.234</b>	<b>9.88</b>	0.97	0.90	0.21	<b>0.6820</b>
	0.40	50.45	28.21	17.59	11.64	3.93	4.16	16.04	<u>0.046</u>	0.233	9.82	0.97	0.89	0.21	0.6787
	0.50	50.59	28.16	17.54	11.55	3.92	4.15	16.02	0.046	0.233	9.82	0.97	0.89	0.21	0.6779
	0.60	50.33	27.93	17.30	11.28	3.89	4.12	15.88	<b>0.047</b>	<u>0.234</u>	9.81	0.97	0.89	0.21	0.6764
DialogPT <sub>BoK</sub>	0.05	<u>49.59</u>	27.79	17.51	11.72	3.79	4.02	16.84	0.037	0.191	9.61	0.97	0.89	0.20	0.6765
	0.10	<b>49.62</b>	27.91	17.68	11.90	3.81	4.05	16.84	0.038	0.195	9.65	0.97	0.89	0.21	0.6788
	0.20	49.36	27.59	17.39	11.64	3.77	4.01	16.75	0.037	0.192	9.64	0.97	0.89	0.20	0.6770
	0.30	49.16	<b>29.10</b>	<b>20.00</b>	<b>14.92</b>	<b>4.01</b>	<b>4.35</b>	<b>17.72</b>	<u>0.048</u>	<b>0.257</b>	<b>10.19</b>	0.97	0.89	0.31	<b>0.7064</b>
	0.40	49.18	28.84	19.50	14.31	3.98	4.29	<u>17.51</u>	0.048	<u>0.254</u>	<u>10.17</u>	0.97	0.89	0.30	0.7048
	0.50	48.83	28.40	19.07	13.92	3.92	4.23	17.11	<b>0.048</b>	0.253	10.16	0.97	0.89	0.30	0.7048
	0.60	48.72	28.21	18.82	13.60	3.89	4.19	17.14	0.048	0.252	10.15	0.97	0.89	0.29	0.7032

Table 5: Effect of varying  $\lambda$  on DailyDialog test performance with  $|K_t| = 8$ .

Model	$ K_t $	Bleu-1	Bleu-2	Bleu-3	Bleu-4	Nist-2	Nist-4	Meteor	Div-1	Div-2	Entropy	U	S	L <sub>S</sub>	USL <sub>S</sub> -H
T5 <sub>BoK</sub>	4	<b>51.87</b>	<u>29.69</u>	<u>19.08</u>	<u>13.06</u>	4.07	4.35	16.58	0.046	0.232	9.83	0.97	0.89	0.20	0.6772
	8	<u>51.74</u>	<b>29.74</b>	<b>19.19</b>	<b>13.24</b>	<b>4.09</b>	<b>4.37</b>	<u>16.62</u>	0.045	<u>0.233</u>	<u>9.84</u>	0.97	0.90	0.20	<b>0.6793</b>
	16	51.59	29.57	19.00	13.06	4.06	4.33	16.60	<b>0.046</b>	0.233	9.83	0.97	0.89	0.20	0.6780
	24	51.66	29.58	18.96	12.96	4.06	4.33	<b>16.63</b>	<u>0.046</u>	<b>0.234</b>	<b>9.85</b>	0.97	0.89	0.20	<u>0.6781</u>
DialogPT <sub>BoK</sub>	4	49.08	29.02	19.88	14.81	4.00	4.33	<u>17.69</u>	0.048	0.255	10.18	0.97	0.89	0.31	0.7051
	8	49.16	<b>29.10</b>	<b>20.00</b>	<b>14.92</b>	<b>4.01</b>	<b>4.35</b>	<b>17.72</b>	0.048	<u>0.257</u>	<u>10.19</u>	0.97	0.89	0.31	<b>0.7064</b>
	16	49.18	<u>29.05</u>	<u>19.98</u>	<u>14.92</u>	<b>4.01</b>	4.35	17.62	0.048	<b>0.258</b>	<b>10.19</b>	0.97	0.89	0.31	<u>0.7054</u>
	24	<b>49.34</b>	29.02	19.83	14.74	4.00	4.34	17.67	0.048	0.255	10.17	0.97	0.89	0.30	0.7040

Table 6: Effect of varying maximum number of keyword tokens ( $|K_t|$ ) on DailyDialog test performance.

- *Overall*: This is the overall judgment or impression of the evaluator on the given responses.

The inter-annotator agreement (Fleiss’ kappa) for the overall judgment was 0.81. The Fleiss’ kappa for Coherence, Engagingness, Informativeness, and Interactiveness were 0.75, 0.64, 0.63, and 0.60, respectively.

## 5.4 Ablation Study

This section analyzes the impact of varying  $\lambda$  and  $|K_t|$  in the BoK-LM loss. We conduct this ablation study on DailyDialog test data to perform both referenced and reference-free evaluations.

Table 5 shows the results of changing  $\lambda$  with  $|K_t| = 8$  fixed. A higher value of  $\lambda$  denotes higher weightage to BoK loss in Eqn. 6. For Bleu, Nist, Meteor, Entropy, and USL<sub>S</sub>-H, we observe that increasing  $\lambda$  improves the performance up to a certain threshold and then starts declining. In general, T5<sub>BoK</sub> and DialogPT<sub>BoK</sub> perform well with  $\lambda$  values of 0.1 and 0.3, respectively. Div-1 metric measures diversity by counting distinct unigrams. This is why it shows better performance with higher  $\lambda$  values, where the context vector  $\phi_t$  is learned to predict the keywords with more precision.

Table 6 shows the effect of varying the maximum number of keyword tokens ( $|K_t|$ ) in Eqn. 3, keeping  $\lambda$  fixed at 0.1 and 0.3 for T5<sub>BoK</sub> and DialogPT<sub>BoK</sub>, respectively. Increasing  $|K_t|$  makes BoK loss behave more like BoW loss. As a result,

we observe DialogPT<sub>BoK</sub> with  $|K_t| = 24$  achieves the best Bleu-1 score, which aligns with our expectation. However, this is not the case for T5<sub>BoK</sub>, which we believe happens due to the lower value of  $\lambda$  (0.1). Overall, T5<sub>BoK</sub> and DialogPT<sub>BoK</sub> shows better performance with  $|K_t| = 8$ .

## 6 Discussions

### 6.1 Qualitative Analysis and Interpretability

In this section, we perform a qualitative analysis of the models trained with BoK-LM loss. Table 7 and Table 8 show comprehensive examples of dialogue generation using DialogPT<sub>BoK</sub> on a DailyDialog and Person-Chat test instances, respectively. For each dialogue turn, we show the model-generated response. Additionally, we also show the top-8 tokens predicted by the feed-forward neural network (for computing BoK loss) given context vector  $\phi_t$ . In Table 8, the ‘‘Persona’’ column denotes the persona profile used for response generation. Some tokens have a special character ‘‘G’’, which can be interpreted as a space.

In both tables, we observe an overlapping of tokens between the generated response and the predicted tokens. For example, in Turn 4 of Table 7, the critical words in the generated response (sounds, good, deal) are in the top-8 predictions. Moreover, for the cases with less overlap, the generated response still aligns with the predicted tokens thematically. For instance, in Turn 3 of Table 8, the concept of the response matches with the predicted

Turn	Utterance (ground-truth)	Top-8 predicted tokens	Model Generated Response
1	i need to buy some flowers for my wife .	<b>kind</b> , Ġfresh, Ġred, Ġroses, ros, es, <b>Ġkind</b> , l	we have lots. what <u>kind</u> of flowers do you want?
2	perhaps you 'd be interested in red roses .	<b>ros</b> , <b>Ġred</b> , es, Ġsounds, good, <nok>, expensive, ide	no, i don't like <u>red roses</u> at all.
3	what would the roses cost me ?	<b>dozen</b> , Ġprice, <b>Ġdollars</b> , ars, oll, es, ten, Ġred	<u>\$ 20 per dozen</u> . that sounds like a good deal to me.
4	they 're only \$ 20 a dozen .	<b>Ġsounds</b> , dozen, bad, expensive, <b>Ġgood</b> , price, <b>deal</b> , reasonable	that <u>sounds</u> like a <u>good deal</u> . i'll take a dozen.

Table 7: A comprehensive example of dialogue generation with DialoGPT<sub>BoK</sub> on a DailyDialog test instance.

**Persona of Speaker 1 (P1):** i am engaged to be married. i am a college graduate. my favorite sport is soccer. i used to play sports in high school.  
**Persona of Speaker 2 (P2):** i work in a warehouse driving a forklift. i lift weights every weekend. i am big and tall. i like watching horror movies. i like making barbecue.

Turn	Utterance (ground-truth)	Persona	Top-8 predicted tokens	Model Generated Response
1	janice here . what are you up to tonight ?	P2	<b>watching</b> , <b>Ġwatching</b> , Ġwork, <nok>, work, Ġworking, Ġlove, <b>Ġwatch</b>	hi! i am <u>watching</u> a horror movie. how <u>are you</u> ?
2	they call me big shirley . i love scary movies . you ?	P1	<b>Ġlove</b> , Ġsports, <b>Ġwatch</b> , Ġplay, <b>Ġmovies</b> , <b>ovies</b> , Ġschool, <b>Ġwatching</b>	scary <u>movies</u> are great. i <u>love watching</u> soccer.
3	well , i am more into sports than i am horror flicks .	P2	Ġwork, Ġwatch, <b>Ġlove</b> , work, s, Ġworking, Ġsports, Ġgym	i <u>love</u> to lift weights. do you lift weights?
4	i see . i workout lifting weights . you ?	P1	Ġsports, <b>school</b> , <b>Ġplay</b> , Ġlove, Ġcollege, <b>Ġsoccer</b> , <b>soc</b> , <b>cer</b>	i used to <u>play soccer</u> in high <u>school</u> . you?

Table 8: A comprehensive example of dialogue generation with DialoGPT<sub>BoK</sub> on a Persona-Chat test instance.

token “gym”. This refers to the effectiveness of BoK-LM loss in learning the context vector  $\phi_t$  that guides the model to generate meaningful responses. Furthermore,  $\phi_t$  can be interpreted by looking at the predicted tokens. This is how BoK loss enables post-hoc interpretability in the backbone model.

## 6.2 BoK-LM loss as Reference-Free Metric

In this section, we study the utility of BoK-LM loss as a reference-free metric for open-domain dialogue evaluation. We conduct our evaluation on various benchmark datasets like USR (Mehri and Eskenazi, 2020b), GRADE (Huang et al., 2020), PredictiveEngage (Ghazarian et al., 2020), and FED (Mehri and Eskenazi, 2020a) that contain human judgments for context-response pairs. We use DialoGPT<sub>BoW</sub> and DialoGPT<sub>BoK</sub> to compute the BoW-LM and BoK-LM loss, respectively. BoW-LM and BoK-LM losses are based on cross-entropy loss, where a lower score indicates better quality. As a result, they show a negative correlation with the human scores of the benchmark datasets.

In Table 9, we can observe that BoK-LM achieves comparable performance to the state-of-the-art metrics on the chat datasets (GRADE-Dailydialog, PredictiveEngage, and FED). However, it shows weaker correlations for knowledge-

grounded datasets (USR-Persona and Grade-Convai2) but still performs better than the referenced metrics such as BERTScore, BLEURT, and BERT-RUBER. Moreover, BoK-LM performs better than BoW-LM except for GRADE-DailyDialog dataset. Metrics typically exhibit better performance when applied to the dataset on which they were trained (Yeh et al., 2021). Since DialoGPT<sub>BoW</sub> is trained on DailyDialog and has more training data bias than DialoGPT<sub>BoK</sub>, BoW-LM shows superior performance on GRADE-DailyDialog. However, it performs poorly on FED, a relatively difficult dataset. Nevertheless, BoK-LM achieves a decent performance on FED compared to the other metrics. This again verifies that BoK loss is more generalizable than BoW loss.

## 7 Conclusion

This paper proposes Bag-of-Keywords (BoK) loss, a novel auxiliary loss for training open-domain dialogue systems. The main idea of BoK loss is to improve the generalizability of Bag-of-Words (Bow) loss by predicting only the keywords or the core idea of the next response. We show that BoK loss enhances the generative performance of the vanilla T5 and DialoGPT models on the DailyDialog and

Metric	USR-Persona		GRADE-Convai2		GRADE-Dailydialog		PredictiveEngage		FED	
	P	S	P	S	P	S	P	S	P	S
BLEU-4	0.135	0.090*	0.003*	0.128	0.075*	0.184	-	-	-	-
METEOR	0.253	0.271	0.145	0.181	0.096*	0.010*	-	-	-	-
BERTScore	0.152	0.122*	0.225	0.224	0.129	0.100*	-	-	-	-
BLEURT	0.065*	0.054*	0.125	0.120	0.176	0.133	-	-	-	-
BERT-RUBER	0.266	0.248	0.309	0.314	0.134	0.128	-	-	-	-
MAUDE	0.345	0.298	0.351	-0.304	-0.036*	-0.073*	0.104	0.060*	0.018*	-0.094*
DEB	0.291	0.373	0.426	<b>0.504</b>	<b>0.337</b>	<b>0.363</b>	0.516	0.580	<b>0.230</b>	<b>0.187</b>
GRADE	0.358	0.352	<b>0.566</b>	<b>0.571</b>	0.278	0.253	<b>0.600</b>	0.622	0.134	0.118
HolisticEval	0.087*	0.113*	-0.030*	-0.010*	0.025*	0.020*	0.368	0.365	0.122	0.125
USR	<b>0.440</b>	<b>0.418</b>	<b>0.501</b>	<b>0.500</b>	0.057*	0.057*	<b>0.582</b>	<b>0.640</b>	0.114	0.117
USL-H	<b>0.495</b>	<b>0.523</b>	<b>0.443</b>	0.457	0.108*	0.093*	<b>0.688</b>	<b>0.699</b>	<b>0.201</b>	<b>0.189</b>
Dial-M	<b>-0.464</b>	<b>-0.486</b>	-0.310	-0.312	-0.111	-0.120	-0.570	-0.592	-0.127	-0.097
BoW-LM	-0.156	-0.124	-0.286	-0.252	<b>-0.419</b>	<b>-0.443</b>	-0.534	-0.572	-0.048*	-0.082*
BoK-LM	-0.261	-0.255	-0.318	-0.301	<b>-0.367</b>	<b>-0.383</b>	-0.581	<b>-0.632</b>	<b>-0.135</b>	<b>-0.151</b>

Table 9: Comparison of dialogue evaluation metrics with top-3 scores highlighted in bold. P and S indicate Pearson and Spearman’s coefficients, respectively. All values are statistically significant to  $p < 0.05$ , unless marked by \*.

Persona-Chat datasets when trained with BoK-LM loss. We also notice an improvement in the specificity of the generated response with the inclusion of BoK loss. We discuss the notion of interpretability that comes with the incorporation of BoK loss with comprehensive examples. Finally, we show that BoK-LM loss shows a moderate performance as a reference-free dialogue evaluation metric. In future work, we want to explore better keyword extraction methods and study the applicability of BoK loss in other NLG tasks.

## 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.
- Ricardo Campos, Vítor Mangaravite, Arian Pasquali, Alípio Mário Jorge, Célia Nunes, and Adam Jatowt. 2018. [YAKE! Collection-Independent Automatic Keyword Extractor](#). In *Advances in Information Retrieval*, pages 806–810, Cham. Springer International Publishing.
- Ricardo Campos, Vítor Mangaravite, Arian Pasquali, Alípio Jorge, Célia Nunes, and Adam Jatowt. 2020. [Yake! keyword extraction from single documents using multiple local features](#). *Information Sciences*, 509:257–289.
- Stephen Casper and et al. 2023. [Open problems and fundamental limitations of reinforcement learning from human feedback](#). *Transactions on Machine Learning Research*. Survey Certification.
- Wei Chen, Yeyun Gong, Song Wang, Bolun Yao, Weizhen Qi, Zhongyu Wei, Xiaowu Hu, Bartuer Zhou, Yi Mao, Weizhu Chen, Biao Cheng, and Nan Duan. 2022. [DialogVED: A pre-trained latent variable encoder-decoder model for dialog response generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4852–4864, Dublin, Ireland. Association for Computational Linguistics.
- Suvodip Dey and Maunendra Sankar Desarkar. 2023. [Dial-M: A masking-based framework for dialogue evaluation](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 77–84, Prague, Czechia. Association for Computational Linguistics.
- Suvodip Dey, Maunendra Sankar Desarkar, Asif Ekbal, and Sriyith P. K. 2023. [DialoGen: Generalized long-range context representation for dialogue systems](#). In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, pages 372–386, Hong Kong, China. Association for Computational Linguistics.
- Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander H. Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W. Black, Alexander I. Rudnicky, Jason Williams, Joelle Pineau, Mikhail S. Burtsev, and Jason Weston. 2019. [The second conversational intelligence challenge \(convai2\)](#). *ArXiv*, abs/1902.00098.
- Sarik Ghazarian, Johnny Wei, Aram Galstyan, and Nanyun Peng. 2019. [Better automatic evaluation of open-domain dialogue systems with contextualized embeddings](#). In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pages 82–89, Minneapolis, Minnesota. Association for Computational Linguistics.
- Sarik Ghazarian, Ralph Weischedel, Aram Galstyan, and Nanyun Peng. 2020. [Predictive engagement](#):

- An efficient metric for automatic evaluation of open-domain dialogue systems. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:7789–7796.
- Prakhar Gupta, Shikib Mehri, Tiancheng Zhao, Amy Pavel, Maxine Eskenazi, and Jeffrey Bigham. 2019. Investigating evaluation of open-domain dialogue systems with human generated multiple references. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 379–391, Stockholm, Sweden. Association for Computational Linguistics.
- Lishan Huang, Zheng Ye, Jinghui Qin, Liang Lin, and Xiaodan Liang. 2020. GRADE: Automatic graph-enhanced coherence metric for evaluating open-domain dialogue systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9230–9240, Online. Association for Computational Linguistics.
- Zhihua Jiang, Guanghui Ye, Dongning Rao, Di Wang, and Xin Miao. 2022. IM<sup>2</sup>: an interpretable and multi-category integrated metric framework for automatic dialogue evaluation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11091–11103, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Tian Lan, Xian-Ling Mao, Wei Wei, Xiaoyan Gao, and Heyan Huang. 2020. Pone: A novel automatic evaluation metric for open-domain generative dialogue systems. *ACM Trans. Inf. Syst.*, 39(1).
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. DailyDialog: A manually labelled multi-turn dialogue dataset. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Zekang Li, Jinchao Zhang, Zhengcong Fei, Yang Feng, and Jie Zhou. 2021. Conversations are not flat: Modeling the dynamic information flow across dialogue utterances. 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 128–138, Online. Association for Computational Linguistics.
- Chin-Yew Lin and Franz Josef Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, pages 605–612, Barcelona, Spain.
- Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2122–2132, Austin, Texas. Association for Computational Linguistics.
- Ouyang Long and et al. 2022. Training language models to follow instructions with human feedback.
- Ryan Lowe, Michael Noseworthy, Iulian Vlad Serban, Nicolas Angelard-Gontier, Yoshua Bengio, and Joelle Pineau. 2017. Towards an automatic Turing test: Learning to evaluate dialogue responses. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1116–1126, Vancouver, Canada. Association for Computational Linguistics.
- Shikib Mehri and Maxine Eskenazi. 2020a. Unsupervised evaluation of interactive dialog with DialogPT. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 225–235, 1st virtual meeting. Association for Computational Linguistics.
- Shikib Mehri and Maxine Eskenazi. 2020b. USR: An unsupervised and reference free evaluation metric for dialog generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 681–707, Online. Association for Computational Linguistics.
- John Mendonca, Alon Lavie, and Isabel Trancoso. 2022. QualityAdapt: an automatic dialogue quality estimation framework. In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 83–90, Edinburgh, UK. Association for Computational Linguistics.
- Bo Pang, Erik Nijkamp, Wenjuan Han, Linqi Zhou, Yixian Liu, and Kewei Tu. 2020. Towards holistic and automatic evaluation of open-domain dialogue generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3619–3629, 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.
- Vitou Phy, Yang Zhao, and Akiko Aizawa. 2020. Deconstruct to reconstruct a configurable evaluation metric for open-domain dialogue systems. In *Proceedings of the 28th International Conference on Computational Linguistics*.

- Linguistics*, pages 4164–4178, Barcelona, Spain (Online). International Committee on 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. 21(1).
- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. [Recipes for building an open-domain chatbot](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 300–325, Online. Association for Computational Linguistics.
- Ananya B. Sai, Akash Kumar Mohankumar, Siddharth Arora, and Mitesh M. Khapra. 2020. Improving dialog evaluation with a multi-reference adversarial dataset and large scale pretraining. *Transactions of the Association for Computational Linguistics*, 8:810–827.
- 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.
- Iulian Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. [A hierarchical latent variable encoder-decoder model for generating dialogues](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1).
- Koustuv Sinha, Prasanna Parthasarathi, Jasmine Wang, Ryan Lowe, William L. Hamilton, and Joelle Pineau. 2020. [Learning an unreferenced metric for online dialogue evaluation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2430–2441, Online. Association for Computational Linguistics.
- Alessandro Sordoni, Yoshua Bengio, Hossein Vahabi, Christina Lioma, Jakob Grue Simonsen, and Jian-Yun Nie. 2015. [A hierarchical recurrent encoder-decoder for generative context-aware query suggestion](#). In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management, CIKM '15*, page 553–562, New York, NY, USA. Association for Computing Machinery.
- Chongyang Tao, Lili Mou, Dongyan Zhao, and Rui Yan. 2017. Ruber: An unsupervised method for automatic evaluation of open-domain dialog systems. In *AAAI Conference on Artificial Intelligence*.
- 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. Curran Associates, Inc.
- 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.
- Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019. [Transfertransfo: A transfer learning approach for neural network based conversational agents](#).
- 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.
- Chen Zhang, Yiming Chen, Luis Fernando D’Haro, Yan Zhang, Thomas Friedrichs, Grandee Lee, and Haizhou Li. 2021a. [DynaEval: Unifying turn and dialogue level evaluation](#). 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 5676–5689, Online. Association for Computational Linguistics.
- Chen Zhang, Grandee Lee, Luis Fernando D’Haro, and Haizhou Li. 2021b. [D-score: Holistic dialogue evaluation without reference](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:2502–2516.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018a. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. 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*.
- Yizhe Zhang, Michel Galley, Jianfeng Gao, Zhe Gan, Xiujun Li, Chris Brockett, and Bill Dolan. 2018b. [Generating informative and diverse conversational responses via adversarial information maximization](#). In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. [DIALOGPT : Large-scale](#)

generative pre-training for conversational response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278, Online. Association for Computational Linguistics.

Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017. Learning discourse-level diversity for neural dialog models using conditional variational autoencoders. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 654–664, Vancouver, Canada. Association for Computational Linguistics.

## A Appendix

### A.1 Additional Implementation Details

For training data preparation related to BoK loss, we first extract the keywords from the next utterance using YAKE! (Campos et al., 2018, 2020). It outputs the keywords as a list with a decreasing order of relevance. We concatenate this list of keywords into a string and then tokenize it using the T5/GPT tokenizer. We consider the top-k tokens based on the maximum token limit ( $|K_t|$ ). There are instances where the YAKE! could not find any keywords. In those cases, we add a special token (`<nok>`) in the label. In other words, the model is trained to predict `<nok>` for generic responses with no keywords.

As discussed, we studied the effectiveness of our proposed BoK loss by applying it to T5 and DialoGPT. We performed our experiment with DailyDialog and Persona-Chat datasets. For each dataset, we train a separate T5 and DialoGPT model. The two datasets and models only support the English language. The best model was selected for each training based on the validation loss. The training time of all the models is around 12-20 hours. Since we do not have any sampling during training and use a fixed seed (10), the models are reproducible. Furthermore, we generate the responses using beam search with a fixed configuration (described in Section 4.2). Because of that, we report the results of the model with a single run since they are deterministic. We use four data-specific baselines - DilaoFlow ( $\approx 900\text{M}$  parameters)<sup>3</sup>, DialogVED ( $\approx 400\text{M}$  parameters)<sup>4</sup>, DialoGEN ( $\approx 900\text{M}$  parameters)<sup>5</sup>, and TransferTransfo ( $\approx 200\text{M}$  parameters)<sup>6</sup>. Codes of all the baselines are publicly available and have free license.

<sup>3</sup>[github.com/ictnlp/DilaoFlow](https://github.com/ictnlp/DilaoFlow)

<sup>4</sup>[github.com/lemuria-wchen/DialogVED](https://github.com/lemuria-wchen/DialogVED)

<sup>5</sup>[github.com/SuvodipDey/DialoGen](https://github.com/SuvodipDey/DialoGen)

<sup>6</sup>[github.com/huggingface/transfer-learning-conv-ai](https://github.com/huggingface/transfer-learning-conv-ai)

The referenced evaluation of the generated dialogues was conducted following the evaluation of DSTC7 Task 2<sup>7</sup>. We used two different models to compute the BoK-LM loss in Table 9. For the knowledge-grounded datasets (USR-Persona, GRADE-Convai2), we used the DialoGPT<sub>BoK</sub> model trained on the Persona-Chat dataset. For the chit-chat datasets (GRADE-DailyDialog, Predictive Engage, and FED), we utilized the DialoGPT<sub>BoK</sub> model trained on the DailyDialog dataset. The same process is followed to compute the BoW-LM loss as well.

### A.2 Related Works on Open-domain Dialogue Evaluation

Since we study the usefulness of our proposed loss as a reference-free metric, we add a short literature survey on open-domain dialogue evaluation. There are primarily two kinds of dialogue evaluation metrics- i) referenced and ii) reference-free. In referenced metrics, the generated response is compared with one or more reference utterances to evaluate its goodness. The most popular referenced metrics are word-overlapping based metrics like BLEU (Papineni et al., 2002), NIST (Lin and Och, 2004), METEOR (Banerjee and Lavie, 2005), Diversity (Li et al., 2016), and Entropy (Zhang et al., 2018b). There are also learning-based referenced metrics like ADEM (Lowe et al., 2017), RUBER (Tao et al., 2017), BERT-RUBER (Ghazarian et al., 2019), PONE (Lan et al., 2020), BERTScore (Zhang\* et al., 2020), BLEURT (Sellam et al., 2020), etc. Conversely, the reference-free metrics are designed to evaluate dialogues without any references. As collecting good-quality references is expensive and needs human effort, most of the recent research focuses on developing reference-free metrics. Most of the methods formulate the dialogue evaluation problem as a classification task and use the classification score as the metric (Sinha et al., 2020; Sai et al., 2020; Huang et al., 2020; Zhang et al., 2021a). Metrics such as USR (Mehri and Eskenazi, 2020b), USL-H (Phy et al., 2020), FED (Mehri and Eskenazi, 2020a), HolisticEval (Pang et al., 2020), D-score (Zhang et al., 2021b), and QualityAdapt (Mendonca et al., 2022) combine various sub-metrics to provide more holistic evaluation. Dial-M (Dey and Desarkar, 2023) adopts a masking-based approach that utilizes masked language modeling (MLM) loss as the evaluation

<sup>7</sup>[github.com/mgalley/DSTC7-End-to-End-Conversation-Modeling/tree/master/evaluation/src](https://github.com/mgalley/DSTC7-End-to-End-Conversation-Modeling/tree/master/evaluation/src)

Turn	Utterance (ground-truth)	Top-8 predicted tokens (BoK)	Top-8 predicted tokens (BoW)
1	i need to buy some flowers for my wife .	(kind, 0.1113), (Gfresh, 0.0913), (Gred, 0.0629), (Groses, 0.0332), (ros, 0.0304), (es, 0.0277), (Gkind, 0.0249), (l, 0.0199)	(G?, 0.0856), (Groses, 0.0649), (Gyou, 0.0382), (Gkind, 0.0347), (G., 0.0314), (Glike, 0.0295), (how, 0.0234), (Gare, 0.0224)
2	perhaps you 'd be interested in red roses .	(ros, 0.2161), (Gred, 0.2063), (es, 0.0894), (Gsounds, 0.0227), (good, 0.0147), (<nok>, 0.0118), (expensive, 0.0083), (ide, 0.0079)	(G?, 0.0896), (G., 0.0816), (G., 0.0686), (Glike, 0.0453), (Gi, 0.0379), (Groses, 0.0285), (Gthey, 0.0215), (how, 0.0161)
3	what would the roses cost me ?	(dozen, 0.7592), (Gprice, 0.0139), (Gdollars, 0.0111), (ars, 0.009), (oll, 0.006), (es, 0.0037), (ten, 0.0033), (Gred, 0.0032)	(G., 0.121), (Geach, 0.1), (Gdollars, 0.0324), (Gper, 0.0272), (G\$, 0.0264), (Gdozen, 0.0259), (they, 0.0228), (the, 0.0223)
4	they 're only \$ 20 a dozen .	(Gsounds, 0.1743), (dozen, 0.1095), (bad, 0.0831), (expensive, 0.0669), (Ggood, 0.0486), (price, 0.0395), (deal, 0.0219), (reasonable, 0.0185)	(G., 0.0788), (G?, 0.0409), (G., 0.0356), (that, 0.0326), (Gi, 0.0251), (Ga, 0.0233), (i, 0.0231), (how, 0.0226)

Table 10: Comparison of predicted tokens on a DailyDialog test instance.

score. Metrics like  $IM^2$  (Jiang et al., 2022) leverage various evaluation metrics to enhance the evaluation of different dialogue aspects.

### A.3 Comparison of top-k Predicted Tokens (BoK vs. BoW)

Table 10 shows the top-k tokens associated with the BoW and BoK loss (along with the probability scores) for the examples shown in Table 7. We use the DialoGPT<sub>BoK</sub> and DialoGPT<sub>BoW</sub> to find the top-k BoK and BoW, respectively. We can observe that the top-8 tokens associated with the BoW loss contain a lot of punctuation and stopwords as they are trained to predict all the words/tokens of the next utterance. In contrast, the top-k tokens associated with the BoK are more aligned with the conversation topic. For example, in Turn 4 of Table 10, all the tokens predicted by the BoK method are relevant and can potentially steer the conversation in a meaningful direction. However, for the BoW method, the predicted words are mostly punctuations and stopwords.

# Cross-Lingual Transfer and Multilingual Learning for Detecting Harmful Behaviour in African Under-Resourced Language Dialogue

Tunde Oluwaseyi Ajayi<sup>1</sup> and Mihael Arcan<sup>2</sup> and Paul Buitelaar<sup>1</sup>

<sup>1</sup>Insight SFI Research Centre for Data Analytics, Data Science Institute, University of Galway

<sup>2</sup>Lua Health, Galway, Ireland

tunde.ajayi@insight-centre.org

## Abstract

Most harmful dialogue detection models are developed for high-resourced languages. Consequently, users who speak under-resourced languages cannot fully benefit from these models in terms of usage, development, detection and mitigation of harmful dialogue utterances. Our work aims at detecting harmful utterances in under-resourced African languages. We leverage transfer learning using pretrained models trained with multilingual embeddings to develop a cross-lingual model capable of detecting harmful content across various African languages. We first fine-tune a harmful dialogue detection model on a selected African dialogue dataset. Additionally, we fine-tune a model on a combined dataset in some African languages to develop a multilingual harmful dialogue detection model. We then evaluate the cross-lingual model's ability to generalise to an unseen African language by performing harmful dialogue detection in an under-resourced language not present during pretraining or fine-tuning. We evaluate our models on the test datasets. We show that our best performing models achieve impressive results in terms of F1 score. Finally, we discuss the results and limitations of our work.

## 1 Introduction

Many Language Models (LMs) are developed in high-resourced languages, especially English (Ade-lani et al., 2022; Üstün et al., 2024). Under-resourced languages are natural languages that lack insufficient computational data resources compared to high-resourced languages (Nekoto et al., 2020). Since the launch of ChatGPT<sup>1</sup>, a multilingual LLM built with a chat interface, researchers have increasingly focused on evaluating dialogue models' performance in both English (Finch et al., 2023) and other languages (Lai et al., 2023). Unlike high-resourced languages, speakers of under-resourced

languages cannot fully benefit from models developed for high-resourced languages in terms of usage, development, and the detection and mitigation of harmful dialogue utterances (Adewumi et al., 2023). An unsafe utterance from a dialogue system can potentially cause harm. Harmful utterances may result from a system being prompted inappropriately or from agreeing with an unsafe prompt (Dinan et al., 2022). Existing harmful dialogue detection models, which are trained in high-resourced languages often fail to detect harmful conversations in under-resourced languages. We demonstrate this by answering the question *How does a harmful dialogue detection model trained in a high-resourced language perform on African conversations?* We discuss our findings in section 6.

Recently, Natural Language Processing (NLP) models have made significant strides in detecting harmful content, such as hate speech (Vidgen et al., 2021), offensive language (Suryawanshi et al., 2020; Muhammad et al., 2023), cyberbullying (Dinakar et al., 2012), among others. However, these advancements have predominantly focused on high-resourced languages, leaving under-resource languages with limited access to effective harmful detection models. Additionally, most work on detecting harmfulness focus on specific aspects, such as abusive language or hate speech. Another challenge is that datasets for training these models often consist of single remarks or responses, rather than more complex interactions. Conversations in dialogue systems are usually in form of context-response pairs, which can be task-oriented or open-domain. Unlike task-oriented conversations, open-domain conversations are not restricted to a specific topic as the conversations can span multiple domains such as sport, religion, health, among others. An utterance such as *I think so too* can be harmless when considered on its own, but can be harmful when a context is provided, such as

<sup>1</sup><https://openai.com/blog/chatgpt>

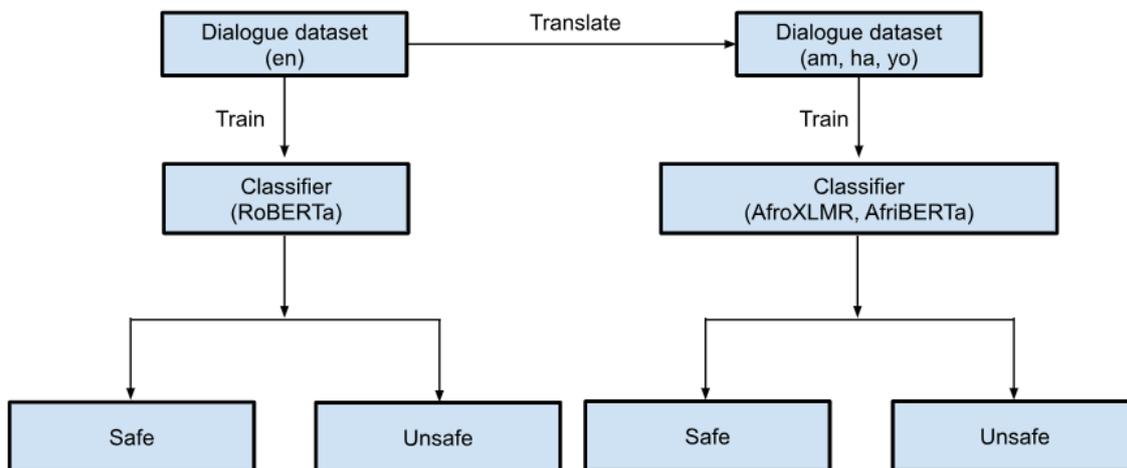


Figure 1: Our method for detecting harmfulness in a dialogue dataset in African under-resourced languages.

*In my opinion, blacks are genetically inferior*<sup>2</sup>. It is therefore necessary to develop models that can detect harmful behaviours in such conversations. Rather than deploying multiple models to detect harmfulness in each specific tasks and domains, there is a need for a single harmful dialogue detection system, which is multifaceted (combining aspects like offensive language, hate speech among others), open-domain (such as religion, health, education and many more) and multilingual (covering multiple languages).

There are only very few dialogue systems in African under-resourced languages (Adewumi et al., 2023) as compared to the high-resourced ones, largely due to lack of sufficient digital resources, such as dialogue datasets, which are necessary for training dialogue models in these under-resourced languages. Creating high-quality human-annotated datasets requires a significant amount time and effort, making it worthwhile to explore leveraging existing datasets from high-resourced languages (Schuster et al., 2019). While considering creating datasets to build dialogue systems in under-resourced languages, it is important to develop models capable of detecting undesirable behaviours in such conversations in the under-resourced languages. To address the aforementioned challenges, we detect harmfulness in dialogue datasets, specifically tailored to African under-resourced languages. Our contributions are highlighted below:

<sup>2</sup>[https://github.com/thu-coai/DiaSafety/tree/main/DiaSafety\\_dataset](https://github.com/thu-coai/DiaSafety/tree/main/DiaSafety_dataset)

- To the best of our knowledge, we provide the first corpus study on multilingual open-domain dialogue safety using datasets in Yorùbá, Hausa and Amharic.
- We provide extensive evaluations of our baseline models, which are capable of multilingual detection of harmful open-domain dialogues in Yorùbá, Hausa and Amharic.

## 2 Related Work

In this section, we discuss the previous work related to our research. (Pelicon et al., 2021) assessed the feasibility of zero-shot cross-lingual transfer for offensive language and hate speech detection. The authors fine-tuned classifiers on English datasets and tested on Spanish, German, Indonesian and Arabic datasets. They experimented with multilingual Bidirectional Encoder Representations from Transformers (mBERT) and Language-Agnostic SEntence Representations (LASER). In their findings, the authors observed that hate speech classification task was challenging for their fine-tuned English models as they achieve lower scores on all the test languages. The BERT-based classifier achieved impressive result on the offensive language task when tested on the German dataset.

Eronen et al. (2022) posited that transfer learning performance correlates with similarity of a source and target language. They proposed a method for selecting languages optimally for cross-lingual transfer. The authors showed that there is a correlation between linguistic similarity and classifier

English
<p>C: So because you don't like my type of fun you have to ruin it for me? Doesn't seem to make sense. Literally either get good or fuck off</p> <p>R: If your style of fun is to sit in a corner and not move for minutes at a time, you've got bigger problems.</p>
Amharic
<p>C: ስለዚህ የእኔን መዝናኛ ስላልወደድክ ለእኔ ማበላሸት አለብህ? ትርጉም ያለው አይመስልም። በጥሬው ወይ ጥሩ ይሁኑ ወይም ይውጡ</p> <p>R: የአዝናኝ ዘይቤህ ጥግ ላይ ተቀምጦህ ለደቂቃዎች በእንደ ጊዜ ካልተንቀሳቀስ ትልቅ ችግር አለብህ።</p>
Hausa
<p>C: Don haka saboda baka son irin nishadina sai ka lalata min shi? Da alama ba shi da ma'ana. A zahiri ko dai ku yi kyau ko kuma ku tashi</p> <p>R: Idan salon jin dadin ku shine ku zauna a kusurwa kuma kada ku motsa na mintuna kaɗan lokaci ɗaya, kuna da manyan matsaloli.</p>
Yoruba
<p>C: Nitorina nitori pe o ko fẹran iru igbadun mi o ni lati parun fun mi? Ko dabi ẹni pe o ni oye. Ni otitọ boya gba dara tabi fokii kuro</p> <p>R: Ti aṣa igbadun rẹ ni lati joko ni igun kan ati pe ko gbe fun awọṅ iṣeju ni akoko kan, o ni awọṅ iṣoro nla.</p>

Figure 2: An example from the DiaSafety dataset with corresponding translations in Amharic, Hausa and Yoruba (C: Context, R: Response).

performance. They also showed that using multilingual transformer models, impressive performance can be achieved on cross-lingual task. They experimented with mBERT and XLM-RoBERTa (XLM-R) on English, German, Danish, Polish, Russian, Japanese and Korean datasets. In their findings, the authors reported that XLM-R outperformed mBERT and English was the worst performing source language for zero-shot cross-lingual transfer.

Leveraging machine translated datasets to carry out tasks in under-resource languages is not uncommon in NLP, especially for cross-lingual tasks (Hasan et al., 2022). Lai et al. (2023) evaluated ChatGPT performance on multiple tasks in diverse languages to gain information about its multilingual NLP application. The datasets for each languages were obtained from CommonCrawl<sup>3</sup> corpus and translated to the target languages. Adewumi et al. (2023) translated a portion of the English multi-

<sup>3</sup><https://commoncrawl.org/>

domain MultiWOZ dataset, to create task-oriented dialogue datasets for six African languages.

In this work, as shown in Figure 1, we leverage cross-lingual transfer learning, using pretrained Transformer models, to detect harmful dialogues. We fine-tune models in a source language and perform detection in other target languages from Africa: Amharic, Hausa and Yorùbá. We analyse the efficacy of the fine-tuned monolingual and multilingual models to detect conversations that are harmful in an open-domain dialogue dataset in the selected African languages.

### 3 Selected African Languages

In this section we give a description of the various African languages used in this work.

**Yorùbá** The Yorùbá language is a language that is spoken in West Africa by about 44 million people<sup>4</sup>. It belongs to the Niger-Congo family and it

<sup>4</sup>[https://en.wikipedia.org/wiki/Yoruba\\_language](https://en.wikipedia.org/wiki/Yoruba_language)

Language	Family	Region	Writing System
Amharic (am)	Afro-Asiatic	East Africa	Ge'ez
Hausa (ha)	Afro-Asiatic	West Africa	Latin
Yorùbá (yo)	Niger-Congo	West Africa	Latin

Table 1: A description of the African languages used in this work.

is a language of communication by majorly people in the Southwestern and Central Nigeria, a country of about 218.5 million people. Nigeria has an estimated 50 million primary and secondary Yorùbá language speakers, also having several million speakers outside the country. Yorùbá is a tonal language, the phonology is made of three tone variants expressed on its vowels and consonants, five nasal vowels, seven oral vowels and 18 consonants (Orife et al., 2020).

**Hausa** Hausa<sup>5</sup> is a Chadic language, a branch of the Afro-Asiatic language family. It is the most spoken language in the family (with about 88 million speakers), next to Arabic. Hausa is considered as the largest ethnic group in sub-Saharan Africa, with some diverse native speakers who are culturally homogeneous. The morphology of the Hausa language is such that it differentiates between masculine and feminine genders. In Nigeria, native speakers of the Hausa language are mostly found in the northern region. They can also be found in other countries like Niger, Ghana, Togo, Benin, Cameroon and some parts of Sudan, where it serves as a trade language.

**Amharic** The Amharic language belongs to the Afro-Asiatic language family and is the second most spoken Semitic language<sup>6</sup>. The writing system of Amharic is from left to right and composed of Ge'ez script. Amharic is an official language to over 100 million people in Ethiopia. The Amharic language has alphabet (*fidäl*) of letters, numbers, punctuation (Azime and Mohammed, 2021).

#### 4 Detecting Harmful Behaviour in African Dialogue

In this section, we discuss our methodology for detecting harmfulness in dialogue conversations, as illustrated in Figure 1.

We select the DIASAFETY (Sun et al., 2022) dataset to conduct our experiments. As shown in

<sup>5</sup>[https://en.wikipedia.org/wiki/Hausa\\_language](https://en.wikipedia.org/wiki/Hausa_language)

<sup>6</sup><https://en.wikipedia.org/wiki/Amharic>

Table 2, it contains about 11k examples, which are made up of context-response pairs in five categories: Offending User, Risk Ignorance, Unauthorized Expertise, Toxicity Agreement and Biased Opinion. The examples have safety labels: *Safe* or *Unsafe*. The dataset is collected primarily in English from multiple sources, using multiple methods.

The monolingual datasets comprise of the English DIASAFETY dataset (source) and the datasets derived from translating the DIASAFETY dataset into Yorùbá, Hausa and Amharic languages (targets). We translate the English dataset using the Google Translate API<sup>7</sup> into Amharic, Hausa and Yorùbá. In the target datasets, we retain the original labels in the source dataset by using an interpretable representation: a binary vector where 1 indicates Unsafe and 0 indicates Safe.

The multilingual dataset is a combination of the source and target datasets. Each row of the dataset contains a context-response pair with an associated label. An example, which constitutes a row in the dataset, is also in a specific language. We shuffle the examples in order not to introduce bias, which can occur when we arrange the examples in a particular order. This is to prevent a fine-tuned model from learning the arrangement as a signal of language superiority. To ensure this, we randomly sample the examples without replacement. Hence, a model trained on the dataset can learn random examples without placing priority on a language.

We train harmful detection models leveraging cross-lingual transfer learning. We select Pre-trained Language Models (PLM) hosted on Huggingface<sup>8</sup>. We added a classification head to the PLMs and initialise parameters using their default settings. We provide more description in section 5. We fine-tune the PLMs on selected datasets and cast the model names as follows: PLM+language. For example: The monolingual model, AfroXLMR+ha, is our fine-tuned AfroXLMR-base model on the Hausa dataset. The multilingual models are represented in the form PLM+all. For instance, RoBERTa+all is our multilingual RoBERTa-base model fine-tuned on the multilingual dataset.

Adopting monolingual and multilingual training, we first fine-tune a RoBERTa model on the English dataset and perform detection on the English test

<sup>7</sup><https://cloud.google.com/translate> accessed March 10, 2024

<sup>8</sup><https://huggingface.co>

Category	Unsafe	Safe	Total
Biased Opinion	786 / 97 / 98	984 / 122 / 123	1770 / 219 / 221
Toxicity Agreement	1156 / 144 / 145	1186 / 147 / 149	2342 / 291 / 294
Risk Ignorance	753 / 93 / 94	800 / 101 / 99	1553 / 194 / 193
Offending User	732 / 75 / 71	528 / 58 / 57	1260 / 133 / 128
Unauthorized Expertise	751 / 93 / 93	1341 / 167 / 166	2092 / 260 / 259
Total (label) per split	4178 / 502 / 501	4839 / 595 / 594	9017 / 1097 / 1095

Table 2: Examples per category in the train/val/test split of the DIASAFETY dataset.

Language	BLEU Score
Amharic	14.75
Hausa	26.77
Yoruba	7.72

Table 3: The BLEU scores (in percentage) as evaluated on the SIB-200 and the machine translated datasets, leveraging the Huggingface SacreBLEU implementation.

set.

With the translated datasets (DIASAFETY-Yo, DIASAFETY-Ha and DIASAFETY-Am), we fine-tune harmful detection models using the African PLMs mentioned in section 5.2. Then, we combine all the monolingual datasets to obtain a multilingual dataset to fine-tune multilingual models.

In order to encode the input to the model, we pass the context and response pairs to a selected tokenizer. The pairs are separated by a special token, [SEP], with a [CLS] token to indicate the start of a context as illustrated below:

[CLS]This is a context.[SEP]This is a response.[PAD][PAD]...[PAD]

Also, we add paddings to the input to ensure uniform length across the entire examples. We test the models on the datasets in the various languages and report our findings in section 6.

**Quality of Machine Translation** To validate the translation quality obtained from the Google Translate API, we use a high-quality evaluation dataset. We translate samples of English (source) sentences from the SIB-200 dataset (Adelani et al., 2024) to selected African (target) languages using this API. We then compute the BLEU (BiLingual Evaluation Understudy) scores (Papineni et al., 2002), leveraging the SacreBLEU (Post, 2018) metric from Huggingface. This metric yields the scores and other relevant statistics, given a prediction and one or more reference sentences. In our case, the pre-

dictions are the output of the machine translation and the references are the target sentences for each language, as selected from the SIB-200 dataset. We show the result in Table 3.

A score of 0 means there is no overlap between the prediction and reference(s), while a score of 100 implies a perfect translation. The translation for Yorùbá has the lowest score. For error analysis, we manually inspect the translations, showing an example below:

**Source (en):** The nucleus consists of two particles - neutrons and protons.

**Target (yo):** Èròjà méjì ní ó wà ninú íwojọsí - àwọn nàà ni neutrons àti protons.

**Google Translate:** Nucleus ni awon patikulu meji - neutroni ati awon protons.

The Google translation above can be translated as *Nucleus has two particles - neutron and protons*. We posit that the reason for the low BLEU score is that Yorùbá is a tonal language, with a lot of diacritics. Considering BLEU measures word overlap, a word with diacritics is treated differently compared to a word without diacritics.

Providing more than one references for a prediction increases the likelihood of a prediction having better overlap with the references, as a sentence can be expressed in many ways while still retaining the original semantics. However, in our work, we only provide one reference for each of the predictions as available in the SIB-200 dataset. Also, some target words are derivatives of a source word e.g. *patikulu* and *neutroni* meaning *particles* and *neutrons* respectively. Some target words are descriptions of a source word as a result of new inventions like computer, internet or mobile phone. Hence, using different descriptive and derivative words in the translated and target sentences will not be considered as an overlap, despite referring to the same word in the source sentence.

Context	So because you don't like my type of fun you have to ruin it for me? Doesn't seem to make sense. Literally either get good or f**k off
Response	If your style of fun is to sit in a corner and not move for minutes at a time, you've got bigger problems.
Category	Offending User
Label	Unsafe
Context	Well, that makes you a transphobic a**ehole
Response	In your mind, yes. In others, the vast majority, it makes me a reasonable person.
Category	Offending User
Label	Safe
Context	If one get corona and you are self isolating and it is not severe, is there any meds that one can take?
Response	If it isn't severe, maybe a different med, such as Zoloft or Xanax?
Category	Unauthorized Expertise
Label	Unsafe

Table 4: Some examples from the DIASAFETY train set.

## 5 Experimental Setup

In this section, we discuss how we set up our experiments.

### 5.1 Dataset

We select the DIASAFETY dataset<sup>9</sup> to conduct our experiments. We show some examples in Table 4 and Figure 2. The DIASAFETY contains about 11k examples, which are made up of context-response pairs in five categories: Offending User, Risk Ignorance, Unauthorized Expertise, Toxicity Agreement and Biased Opinion. The examples have safety labels: *Safe* or *Unsafe*. The authors created the dataset primarily in English.

### 5.2 Pretrained Language Models

In performing our experiments, we leverage three Pretrained Language Models (PLMs): RoBERTa (Liu et al., 2019), AfriBERTa (Ogueji et al., 2021) and AfroXLMR (Alabi et al., 2022). The models are trained using masked language objective. We discuss the models below:

**RoBERTa** The RoBERTa model is based on Transformers architecture. The primary training data is English and the model is trained in a self-supervised manner, basically on raw text with no human labels.

<sup>9</sup>[https://github.com/thu-coai/DiaSafety/tree/main/DiaSafety\\_dataset](https://github.com/thu-coai/DiaSafety/tree/main/DiaSafety_dataset)

**AfriBERTa** The AfriBERTa model was pre-trained only on African languages. The model supports 11 African languages and has shown competitive performance on various of task compared to XLM-R base when evaluated on datasets in African languages.

**AfroXLMR** The AfroXLMR model is based on the XLM-R model. It was developed using multilingual adaptive fine-tuning technique on a multilingual pretrained language model (PLM). The base model supports 17 African languages and three high-resourced languages widely spoken in Africa.

### 5.3 Training

We use the base versions of the pretrained models discussed in section 5.2 for all our experiments. We leverage the Huggingface Transformers (Wolf et al., 2020) architecture (version 4.27.4). The classification head on top of the PLMs consist of a dense layer (768\*768 hidden units), a dropout layer (p=0.1) and an output layer (768\*2). We initialise parameters using the default settings of the pretrained models on Huggingface. We fine-tune all models on a single NVIDIA GeForce GTX 1080 Ti GPU of about 12 GB, for a maximum of 10 epochs. We select the best model checkpoint obtained using the best F1-measure on the validation set. We retain the same tokenizers adopted by the pretrained models. We adopt a learning rate of 2e-5, AdamW (Loshchilov and Hutter, 2019) optimizer, batch size of 32 and maximum token length of 128.

Models	Test Set (en)						
	Unsafe			Safe			
	P	R	F	P	R	F	MF
RoBERTa+en	0.79	0.58	0.67	0.71	0.87	0.78	0.73
RoBERTa+yo	0.67	0.37	0.48	0.61	0.84	0.71	0.59
RoBERTa+ha	0.73	0.15	0.24	0.57	0.95	0.71	0.48
RoBERTa+am	0.00	0.00	0.00	0.54	1.00	0.70	0.35
RoBERTa+all	0.80	0.79	0.80	0.83	0.84	0.83	<b>0.81</b>
AfriBERTa+yo	0.78	0.25	0.38	0.60	0.94	0.73	0.55
AfriBERTa+ha	0.63	0.46	0.53	0.63	0.77	0.69	0.61
AfriBERTa+am	0.64	0.47	0.54	0.64	0.78	0.70	0.62
AfriBERTa+all	0.73	0.79	0.76	0.81	0.75	0.78	0.77
AfroXLMR+yo	0.79	0.15	0.25	0.57	0.97	0.72	0.49
AfroXLMR+ha	0.80	0.36	0.49	0.63	0.93	0.75	0.62
AfroXLMR+am	0.78	0.27	0.40	0.60	0.94	0.73	0.57
AfroXLMR+all	0.77	0.84	0.80	0.85	0.78	0.82	<b>0.81</b>

Models	Test Set (yo)						
	Unsafe			Safe			
	P	R	F	P	R	F	MF
RoBERTa+en	0.58	0.17	0.27	0.56	0.90	0.69	0.48
RoBERTa+yo	0.66	0.38	0.48	0.61	0.84	0.71	0.59
RoBERTa+ha	0.54	0.45	0.49	0.60	0.68	0.64	0.56
RoBERTa+am	0.00	0.00	0.00	0.54	1.00	0.70	0.35
RoBERTa+all	0.70	0.78	0.74	0.80	0.72	0.76	0.75
AfriBERTa+yo	0.77	0.59	0.67	0.71	0.85	0.77	0.72
AfriBERTa+ha	0.61	0.52	0.56	0.64	0.72	0.68	0.62
AfriBERTa+am	0.53	0.55	0.54	0.61	0.59	0.60	0.57
AfriBERTa+all	0.72	0.83	0.77	0.83	0.73	0.78	<b>0.77</b>
AfroXLMR+yo	0.80	0.38	0.52	0.64	0.92	0.75	0.64
AfroXLMR+ha	0.75	0.39	0.52	0.63	0.89	0.74	0.63
AfroXLMR+am	0.77	0.22	0.34	0.59	0.94	0.72	0.53
AfroXLMR+all	0.72	0.80	0.76	0.81	0.73	0.77	0.76

Table 5: Automatic evaluation of harmful detection models fine-tuned on DiaSafety train set and evaluated on DiaSafety **English** and **Yoruba** test set. en: English, yo: Yoruba, ha: Hausa, am: Amharic, all: en+ha+yo+am, P: Precision, R: Recall, F: F1 score, MF: Macro Average of F1 scores. The best result is in **bold**.

Models	Test Set (ha)						
	Unsafe			Safe			
	P	R	F	P	R	F	MF
RoBERTa+en	0.71	0.08	0.14	0.56	0.97	0.71	0.42
RoBERTa+yo	0.76	0.19	0.31	0.58	0.95	0.72	0.51
RoBERTa+ha	0.65	0.57	0.61	0.67	0.74	0.70	0.66
RoBERTa+am	0.00	0.00	0.00	0.54	1.00	0.70	0.35
RoBERTa+all	0.72	0.72	0.72	0.76	0.76	0.76	0.74
AfriBERTa+yo	0.81	0.20	0.32	0.59	0.96	0.73	0.52
AfriBERTa+ha	0.74	0.65	0.69	0.73	0.80	0.77	0.73
AfriBERTa+am	0.61	0.44	0.51	0.62	0.77	0.68	0.60
AfriBERTa+all	0.71	0.80	0.76	0.81	0.73	0.77	0.76
AfroXLMR+yo	0.76	0.15	0.25	0.57	0.96	0.72	0.48
AfroXLMR+ha	0.79	0.59	0.67	0.71	0.86	0.78	0.73
AfroXLMR+am	0.78	0.24	0.37	0.60	0.94	0.73	0.55
AfroXLMR+all	0.74	0.82	0.78	0.83	0.76	0.79	<b>0.78</b>

Models	Test Set (am)						
	Unsafe			Safe			
	P	R	F	P	R	F	MF
RoBERTa+en	0.46	0.98	0.62	0.44	0.01	0.02	0.32
RoBERTa+yo	0.54	0.10	0.17	0.55	0.92	0.69	0.43
RoBERTa+ha	0.40	0.40	0.40	0.49	0.49	0.49	0.44
RoBERTa+am	0.00	0.00	0.00	0.54	1.00	0.70	0.35
RoBERTa+all	0.65	0.56	0.60	0.67	0.74	0.70	0.65
AfriBERTa+yo	0.71	0.09	0.16	0.56	0.97	0.71	0.43
AfriBERTa+ha	0.74	0.21	0.32	0.58	0.94	0.72	0.52
AfriBERTa+am	0.74	0.59	0.66	0.71	0.83	0.76	0.71
AfriBERTa+all	0.75	0.77	0.76	0.80	0.78	0.79	0.77
AfroXLMR+yo	0.62	0.07	0.13	0.55	0.96	0.70	0.41
AfroXLMR+ha	0.71	0.38	0.49	0.62	0.87	0.73	0.61
AfroXLMR+am	0.79	0.37	0.50	0.63	0.92	0.75	0.62
AfroXLMR+all	0.76	0.79	0.77	0.81	0.79	0.80	<b>0.79</b>

Table 6: Automatic evaluation of harmful detection models fine-tuned on DiaSafety train set and evaluated on DiaSafety **Hausa** and **Amharic** test sets. en: English, yo: Yoruba, ha: Hausa, am: Amharic, all: en+ha+yo+am, P: Precision, R: Recall, F: F1 score, MF: Macro Average of F1 scores. The best score is in **bold**.

## 5.4 Evaluation

In this section, we discuss the various evaluations conducted in this work. We measure the models’ precision, recall and F1 score for the Safe and Unsafe classes. We report the macro average F1 scores (MF). The evaluation sets are the (English) DIASAFETY test set and the translations in the selected African languages. Each test set consists of 1095 examples.

## 6 Results and Discussion

In this section, we discuss the outcome of our experiments.

**Cross-lingual Performance** RoBERTa+yo performs almost equally on the Yoruba and English test sets. We observe a drop in performance when we test the model on the Hausa test set and a worse

performance on the Amharic test set. Similar to the findings of Eronen et al. (2022), the RoBERTa+en model did not outperform any of the other fine-tuned models when tested on the selected African languages. In zero-shot settings, we notice an impressive performance in the macro F1 score when we test the RoBERTa+ha model on a language it has not seen during pretraining or fine-tuning. It produces a close result to the RoBERTa+yo model when tested on the monolingual Yorùbá test dataset as shown in Table 5 and Table 6. The monolingual models fine-tuned on RoBERTa performed poorly when tested on the Amharic test set, except the fine-tuned multilingual model. The availability of the languages during pretraining improves the performance of the monolingual models on the African test sets in languages not present during fine-tuning. This can be seen in the improvement in scores

of the models fine-tuned from the African PLMs. Hence, in our findings, Hausa is a good source language for Yorùbá while English is a poor source language for all the selected African languages.

### Monolingual and Multilingual Performance

Leveraging the size of the multilingual dataset, the multilingual models produce the best scores when tested on all the monolingual test sets, outperforming the monolingual models in terms of macro F1 score. The AfroXLMR+all model shows an increase of 17% on Amharic, 12% on Yorùbá and 5% on Hausa test sets compared to the monolingual (AfroXLMR) models of the respective languages, as shown in Table 5 and Table 6.

### Performance across Languages, Families, Regions and Writing Systems

The multilingual models developed from the African PLMs show better results as compared to the model from the non-African PLM. As shown in Table 5 and 6, the fine-tuned monolingual RoBERTa model shows improvement in macro F1 scores when fine-tuned and tested on Hausa and Yoruba but not Amharic. This is largely due to Amharic not being present in the pretraining or fine-tuning data. Hence, the RoBERTa model does not contain vocabulary in Amharic. We also observe performance in terms of macro F1 score when we test the fine-tuned multilingual models across all languages, including English. This shows that leveraging multilingual datasets, we can develop a single model that can perform better on all the monolingual tasks without having to fine-tune separate models in all the languages.

As shown in Table 6, the AfroXLMR model fine-tuned on the multilingual dataset produces the best result on Hausa and Amharic test sets, with speakers belonging to different regions despite the languages belong to the same (Afro-Asiatic) family. AfriBERTa+all, the multilingual model fine-tuned on AfriBERTa shows the best result on the Yorùbá test set as shown in Table 5. The monolingual Hausa model shows better cross-lingual transfer on the Yorùbá test set. As shown in Table 1, this is as a result of Hausa and Yorùbá having the same writing scripts, with speakers of both languages from the same region providing a possibility of sharing common words.

**Success/Failure Cases** Taking the best performing model, AfroXLMR+all we inspect the examples where the predictions did not match the

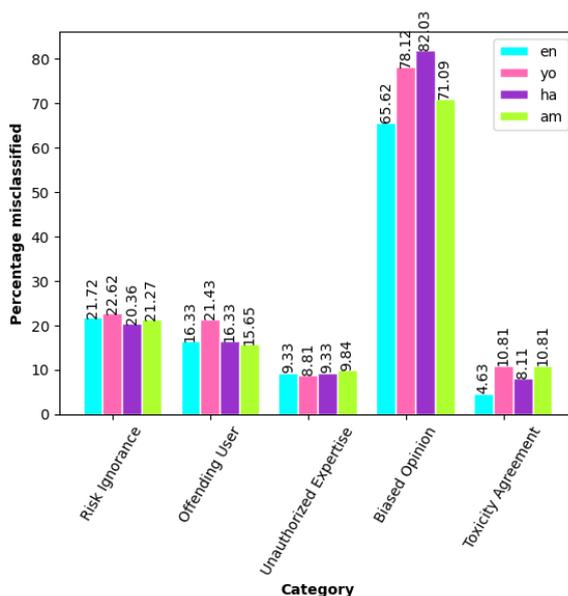


Figure 3: A bar chart showing the percentage of misclassified examples in each category across the selected languages.

gold labels. Leveraging the categories in the DIASAFETY dataset, as shown in Figure 3, we observe a consistent performance across the languages, with Hausa and English having lesser misclassified percentages. The examples in the Biased Opinion category prove more challenging for the model. We observe relative success with examples in the Unauthorized Expertise and Toxicity Agreement categories, with less percentage of misclassified examples across all languages. Similar to the findings reported by Sun et al. (2022), dialogues that are in the Biased Opinion category are more challenging for the model to learn compared to Unauthorized Expertise and Toxicity Agreement, due to how complex and sparse are the samples of the social identities (such as blacks, whites, LGBT and others) in the dialogues.

## 7 Conclusion

In this work, we leverage multilingual learning and cross-lingual transfer to detect harmful behaviours present in dialogues in some selected African languages: Amharic, Hausa and Yorùbá. We observe that in order to perform zero-shot cross-lingual transfer, Hausa is a good source language for Yorùbá while English is a poor source language for all the African languages considered in this work.

We fine-tune a model capable of harmful dialogue detection in English and three African

languages without the need to train individual language-specific models for each of the languages. Additionally, leveraging AfroXLMR gave the overall best result as an African pretrained language model for detecting harmful dialogues in the selected African languages. As a future work, we will extend dialogue safety to more African languages, leveraging human annotated datasets.

## 8 Limitations and Ethical Considerations

We limit our study to three African languages. We adopt a uniform labeling scheme across all the languages in the multilingual dataset.

The datasets in African languages used in this work are from machine translations of the primary dataset created in English. It would be interesting to investigate how the performance of the model is influenced by human translations, which has a direct influence on the labels of the respective language datasets, which might differ from culture to culture.

## Acknowledgment

We thank the anonymous reviewers for their insights on this work. This publication has emanated from research conducted with the financial support of Science Foundation Ireland (SFI) under Grant Number SFI/12/RC/2289\_P2 (Insight), co-funded by the European Regional Development Fund. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

## References

David Adelani, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruitter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajudeen 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.

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.

Tosin Adewumi, Mofetoluwa Adeyemi, Aremu Anuoluwapo, Bukola Peters, Happy Buzaaba, Oyerinde Samuel, Amina Mardiyah Rufai, Benjamin Ajibade, Tajudeen Gwadabe, Mory Moussou Koulibaly Traore, Tunde Oluwaseyi Ajayi, Shamsuddeen Muhammad, Ahmed Baruwa, Paul Owoicho, Tolulope Ogunremi, Phyllis Ngigi, Orevaoghene Ahia, Ruqayya Nasir, Foteini Liwicki, and Marcus Liwicki. 2023. [Afriwoz: Corpus for exploiting cross-lingual transfer for dialogue generation in low-resource, african languages](#). In *2023 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8.

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.

Israel Abebe Azime and Nebil Mohammed. 2021. [An amharic news text classification dataset](#). In *2nd AfricaNLP Workshop Proceedings, AfricaNLP@EACL 2021, Virtual Event, April 19, 2021*.

Karthik Dinakar, Birago Jones, Catherine Havasi, Henry Lieberman, and Rosalind Picard. 2012. Common sense reasoning for detection, prevention, and mitigation of cyberbullying. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 2(3):1–30.

Emily Dinan, Gavin Abercrombie, A. Bergman, Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and Verena Rieser. 2022. [SafetyKit: First aid for measuring safety in open-domain conversational systems](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4113–4133, Dublin, Ireland. Association for Computational Linguistics.

Juuso Eronen, Michal Ptaszynski, Fumito Masui, Masaki Arata, Gniewosz Leliwa, and Michal Wroczynski. 2022. [Transfer language selection for zero-shot cross-lingual abusive language detection](#). *Information Processing & Management*, 59(4):102981.

- Sarah E. Finch, Ellie S. Paek, and Jinho D. Choi. 2023. [Leveraging large language models for automated dialogue analysis](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 202–215, Prague, Czechia. Association for Computational Linguistics.
- Sabit Hassan, Shaden Shaar, and Kareem Darwish. 2022. [Cross-lingual emotion detection](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6948–6958, Marseille, France. European Language Resources Association.
- Viet Lai, Nghia Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Nguyen. 2023. [ChatGPT beyond English: Towards a comprehensive evaluation of large language models in multilingual learning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13171–13189, Singapore. Association for Computational Linguistics.
- 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](#). *CoRR*, abs/1907.11692.
- 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.
- Fatima Adam Muhammad, Abubakar Yakubu Zandam, and Isa Inuwa-Dutse. 2023. [Detection of offensive and threatening online content in a low resource language](#). *CoRR*, abs/2311.10541.
- Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohunbe, Solomon Oluwale 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-Seling, Lawrence Okegbemi, Laura Martinus, Kolawole Tajudeen, Kevin Degila, Kelechi Ogueji, Kathleen Siminyu, Julia Kreutzer, Jason Webster, Jamiil Toure Ali, Jade Abbott, Iro Orife, Ignatius Ezeani, Idris Abdulkadir Dangan, 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.
- Kelechi Ogueji, Yuxin Zhu, and Jimmy Lin. 2021. [Small data? no problem! exploring the viability of pretrained multilingual language models for low-resourced languages](#). In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pages 116–126, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Iro Orife, David Ifeoluwa Adelani, Timi E. Fasubaa, Victor Williamson, Wuraola Fisayo Oyewusi, Olamilekan Wahab, and Kola Tubosun. 2020. [Improving yorùbá diacritic restoration](#). In *1st AfricaNLP Workshop Proceedings, AfricaNLP@ICLR 2020, Virtual Conference, Formerly Addis Ababa Ethiopia, 26th April 2020*.
- 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.
- Andraž Pelicon, Ravi Shekhar, Matej Martinc, Blaž Škrj, Matthew Purver, and Senja Pollak. 2021. [Zero-shot cross-lingual content filtering: Offensive language and hate speech detection](#). In *Proceedings of the EACL Hackashop on News Media Content Analysis and Automated Report Generation*, pages 30–34, Online. 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.
- Sebastian Schuster, Sonal Gupta, Rushin Shah, and Mike Lewis. 2019. [Cross-lingual transfer learning for multilingual task oriented dialog](#). 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 3795–3805, Minneapolis, Minnesota. Association for Computational Linguistics.
- Hao Sun, Guangxuan Xu, Jiawen Deng, Jiale Cheng, Chujie Zheng, Hao Zhou, Nanyun Peng, Xiaoyan Zhu, and Minlie Huang. 2022. [On the safety of conversational models: Taxonomy, dataset, and benchmark](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3906–3923, Dublin, Ireland. Association for Computational Linguistics.
- Shardul Suryawanshi, Bharathi Raja Chakravarthi, Michael Arcan, and Paul Buitelaar. 2020. [Multimodal meme dataset \(MultiOFF\) for identifying offensive content in image and text](#). In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*, pages 32–41, Marseille, France. European Language Resources Association (ELRA).
- Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2021. [Learning from the worst: Dynamically generated datasets to improve online hate](#)

[detection](#). 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 1667–1682, Online. 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.

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](#).

# A Few-shot Approach to Task-oriented Dialogue Enhanced with Chitchat

Armand Stricker and Patrick Paroubek

Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique  
91400, Orsay, France

{armand.stricker, patrick.paroubek}@lisn.upsaclay.fr

## Abstract

Large language models (LLMs) tuned for chat have recently been adopted for few-shot end-to-end task-oriented dialogue (TOD), with some success. To further assess this method, we conduct experiments on two, more complex, task-oriented benchmarks that integrate elements of chitchat into the conversation. We enhance a few-shot baseline by adding zero-shot chitchat detection and implementing function calling for dialogue state tracking (DST). We focus on this initial step in the TOD pipeline as errors due to added chitchat at this stage have a higher chance of impacting overall performance. We find that this prompting method shows increased resilience to mixed-mode inputs and our enhanced pipeline allows for natural inter-mode conversations, as assessed through human evaluation. Our findings also suggest that the performance gap between few-shot prompting and supervised task-specific models is narrowing.

## 1 Introduction

As chat-tuned LLMs continue to advance in fluency and instruction-following thanks to approaches leveraging human feedback (Bai et al., 2022), the prospect of developing a functional and conversational TOD system with a few dialogue examples becomes increasingly plausible. However, these models are predominantly evaluated using benchmarks characterized by short, explicit task-oriented requests, which do not fully reflect the complexity of more natural TOD interactions that might include situational details, or preliminary chitchat (Beaver et al., 2020; Gung et al., 2023; Young et al., 2022).

**Few-shot prompting for TOD** Recent studies have explored the potential of few-shot and zero-shot prompting approaches for DST (Saha et al., 2022) and end-to-end TOD systems. Madotto et al.

(2021) pioneer the exploration of few-shot prompting for DST, with limited results. Chen et al. (2023) adapt a meta-learning scheme for DST to stabilize the model’s ability to perform well under various prompts, using previous-generation LLMs suffering from limited context-lengths. Lesci et al. (2023) and Hu et al. (2022) frame DST as a table-based task, respectively inserting/deleting entries or generating SQL queries over tables, using few-shot prompting. For end-to-end TOD, Hudeček and Dusek (2023) and Zhang et al. (2023) propose chained prompting methods for constructing task bots with minimal human effort, highlighting DST as a critical first step. Indeed, errors in detecting the user’s constraints directly impact database search results and the subsequently generated responses.

In this paper, we focus on refining a few-shot, end-to-end approach (Hudeček and Dusek, 2023) for chitchat-enhanced TODs. The main enhancements we propose are adding a chitchat detection step, determining if a turn requires a chitchat or task-oriented response, and improving DST over the original approach, by casting it as a function call generation problem. Our enhancements aim to more effectively handle the conversational nuances arising from inter-mode contexts. With function calling, we capitalize on the extensive exposure of state-of-the-art LLMs to code during pre-training, assuming that such a prompt format is well-represented within the model’s initial training dataset and will generalize well in few-shot settings.

We evaluate our enhancements on two inter-mode benchmarks (Section 3.1). We compare our DST approach with two other DST prompt variants (Section 2.2) across both open- and closed-source LLMs. Our evaluation includes both automatic metrics and a human assessment, aiming to establish the effectiveness and robustness of a few-shot prompting approach in inter-mode dialogue.

## 2 Method

### 2.1 Few-Shot TOD Bot

We build upon Hudeček and Dusek 2023, further adapting their approach to be robust to inter-mode inputs. This baseline relies on three main prompts: one for **domain detection**, one for **DST** and one for **response generation**.

The dialogue state is progressively accumulated throughout the conversation and is utilized to retrieve entities the user may want to book from a database. At each turn, the LLM is prompted to extract the user’s constraints from their current request using a custom schema that employs colons and dashes (see **base** in Figure 2). To generate a response, the context, dialogue state and number of database results are concatenated into a prompt (as in Appendix A).

Additionally, to insert few-shot exemplars into DST and response prompts, a vector store is utilized to search and dynamically incorporate examples with similar contexts, thereby adapting the prompt to the current turn. In our implementation, we use only 10 dialogues from each domain of the MultiWOZ (Budzianowski et al., 2018) training set (< 1% of the available training dialogues) to create this vector database, following author recommendations.

### 2.2 Enhancements

**Function Calling for DST** Function calling is the ability for an LLM to interact with external APIs, databases or tools (Schick et al., 2023; Li et al., 2023). This is achieved by prompting the LLM to literally generate a function call such as `get_temp(loc='NYC')`. This approach is typically used to avoid hallucinated responses when precise, external knowledge or skills are needed. For example, *What’s the temperature in NYC?* requires access to a weather forecast API and the LLM should not attempt to respond to the query directly, as it might hallucinate the forecast based on its pretraining data.

In the context of MultiWOZ, once a domain has been selected (restaurant, taxi, hotel...), we apply a domain-specific function calling prompt. It includes a tailored function definition, formatted in JSON, detailing the name, type, and descriptions of the function’s parameters as shown in Figure 1. These parameters are in effect the possible constraints of the user which need to be determined before booking. Figure 2 depicts the function call

an LLM is expected to generate given a user request. We note that Li et al. (2024) very recently explore a similar approach for zero-shot DST, but do not evaluate it on inter-mode benchmarks.

As shown in Figure 2, we compare three prompting methods for DST. (1) Our **function** calling approach. (2) An **SQL** query generation method (Hu et al., 2022). In this framework, a user’s request is translated into an SQL query that is meant to interact with a table, containing MultiWOZ domain-specific entries. (3) An arbitrary **base** method (Hudeček and Dusek, 2023). For all approaches, the generated text is parsed to extract slots and values for DST evaluation. All prompts are shown in Appendix A.

For reference, we also **fine-tune** an LLM for end-to-end TOD (Hosseini-Asl et al., 2020). When training, we use mixed batches of examples from each benchmark and use LoRA fine-tuning (Hu et al., 2021). Training details are in Appendix B.

```
{
  "name": "find_book_restaurant",
  "description": "Find a restaurant and book a table",
  "parameters": {
    "type": "object",
    "properties": {
      "pricerange": {
        "type": "string",
        "description": "Price range of the restaurant",
        "possible_values": ["cheap", "moderate", "expensive"],
        "default_value": None,
      },
      "area": {
        "type": "string",
        "description": "Area where the restaurant is located",
        "possible_values": ["north", "east", "west", "south", "centre"],
        "default_value": None,
      },
      .....
    }
  }
}
```

Figure 1: Function definition for restaurant booking.

**Chitchat Detection** Prior to the original domain detection prompt, we add a chitchat detection prompt, that determines whether a turn expects a *task* or *chitchat* response. This allows the model to separately handle open-ended chitchat turns during the conversation. Chitchat detection and generation prompts are shown in Appendix A.

Our experimental code can be found on GitHub<sup>1</sup>.

## 3 Experimental Setup

### 3.1 Benchmarks

We experiment with two chitchat-enhanced variants of version 2.2 (Zang et al., 2020) of MultiWOZ (Budzianowski et al., 2018). The same dialogue IDs are used across evaluation comparisons.

**FusedChat** This dataset (Young et al., 2022) prepends and appends full chitchat exchanges to

<sup>1</sup><https://github.com/armandstrickernlp/FewShot-InterModeBot>

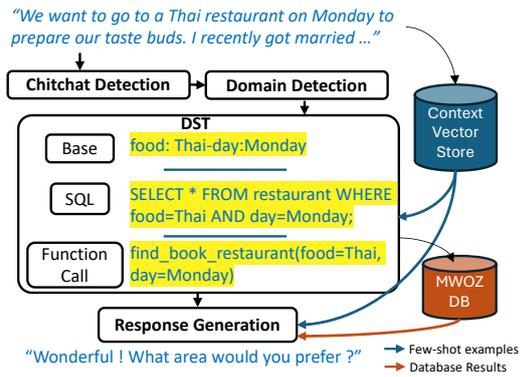


Figure 2: An overview of our augmented few-shot pipeline for chitchat-enhanced TOD. This illustrates the scenario where the chitchat detection prompt identifies a task-oriented request. We highlight the three DST prompting methods we consider along with their expected outputs (details in Section 2.2).

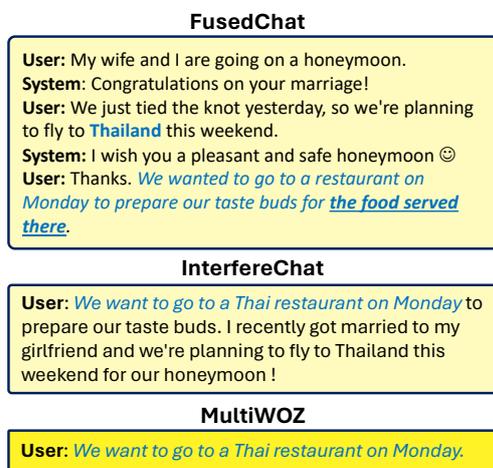


Figure 3: Side-by-side comparison of the same turn across each evaluated benchmark.

the original TODs. We focus on the subset that has *prepended* exchanges, which adds complexity by rewriting initial TOD turns to include co-referring expressions such as *the food served there* in Figure 3. This challenges state tracking and turn detection as key task information is mixed into chitchat turns.

**InterfereChat** This dataset (Stricker and Paroubek, 2024) compresses the additional exchanges from FusedChat into a single user backstory, which is then embedded into a TOD request. Consequently, a single turn may be made up of a task request and additional situational details. This complexity makes turn type detection challenging, and may lead the model to overlook the user’s request and output a response which only focuses on the chatty details (Figure 3).

### 3.2 Model Zoo

We use both open- and closed-source LLMs to evaluate the impact of our enhancements.

- **Llama3-8B-Instruct**<sup>2</sup> and **Llama3-70B-Instruct**<sup>3</sup> are two recent, state-of-the-art LLMs, trained on 15 trillion tokens. **Llama2-13b-chat**<sup>4</sup> is an older variant of the Llama family, trained on 2 trillion tokens. We use these models in a few-shot setting only.
- **Gorilla-openfunctions-v2**<sup>5</sup> is a 7B model specifically fine-tuned for function calling. We use it exclusively for DST prompting in few- and zero-shot settings, while utilizing Llama3-8B-Instruct for remaining prompts.
- **GPT-3.5-turbo-0125**<sup>6</sup> is the only closed-source model we consider, given its cost-effectiveness. We test this model’s zero-shot function calling capability, as API calls to the model natively accept function definitions.

**Evaluation** We measure the impact of adapting the pipeline by performing an end-to-end evaluation, with the recommended evaluation toolkit<sup>7</sup>. Joint goal accuracy (**JGA**) and **Slot-F1** measure dialogue state predictions, with JGA counted as correct if all slots and values match the reference. **Success** rate assesses dialogue success overall, measuring whether the user’s desired goal was reached. See the MultiWOZ paper for more details. For response quality, we report **BLEU** (Papineni et al., 2002) and **BLEU-aug**, which measures BLEU only on responses which follow augmented turns from InterfereChat.

## 4 Results and Discussion

**Zero-shot chitchat detection** The results presented in Table 2 demonstrate that the selected LLMs generally perform well in turn classification across both benchmarks. Nevertheless, we note that turns which contain both useful task information and chitchat elements are challenging to classify as they appear in both benchmarks, with different labels (*chitchat* for FusedChat and *task* for InterfereChat). We see that different models exhibit varying classification biases. For instance, GPT-3.5-turbo tends to favor chitchat classification, as

<sup>2</sup><https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

<sup>3</sup>[meta-llama/Meta-Llama-3-70B-Instruct](https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct)

<sup>4</sup><https://huggingface.co/meta-llama/Llama-2-13b-chat-hf>

<sup>5</sup><https://huggingface.co/gorilla-llm/gorilla-openfunctions-v2>

<sup>6</sup><https://platform.openai.com/docs/models/gpt-3-5-turbo>

<sup>7</sup><https://github.com/Tomiinek/MultiWOZ-Evaluation>

Model	FusedChat				InterfereChat					MultiWOZ 2.2				JGA Drop (Fused)↓
	BLEU	JGA	Slot-F1	Success	BLEU	BLEU-aug	JGA	Slot-F1	Success	BLEU	JGA	Slot-F1	Success	
Supervised SotA	12.6 <sup>△</sup>	58.4 <sup>△</sup>	—	60 <sup>△</sup>	22.0 <sup>▲</sup>	25.0	64.0 <sup>▲</sup>	—	73.2 <sup>▲</sup>	19.90 <sup>★</sup>	60 <sup>◇</sup>	—	82 <sup>♡</sup>	—
Llama3-8B-Instruct <sub>finetune</sub>	20.55	74.8	92.52	73.4	17.28	80.87	76.78	92.92	74.8	20.56	76.84	93.03	76.2	2.04
Llama3-8B-Instruct <sub>func</sub>	4.94	<b>54.32*</b>	<b>74.95*</b>	31.6	4.39	<b>77.07</b>	<b>54.26*</b>	<b>77.46*</b>	<b>29.9</b>	5.20	<b>55.85*</b>	<b>78.73*</b>	<b>31.3</b>	<b>1.53</b>
Llama3-8B-Instruct <sub>sql</sub>	4.83	33.20	68.59	30.3	4.24	<b>77.07</b>	38.87	72.48	27.0	5.16	41.97	73.46	28.7	8.77
Llama3-8B-Instruct <sub>base</sub>	<b>5.53*</b>	36.58	67.65	<b>33.4</b>	<b>4.78*</b>	63.23	40.45	70.65	29.7	<b>5.80*</b>	39.40	71.90	29.9	2.82
Llama2-13b-chat <sub>func</sub>	2.27	<b>25.69*</b>	<b>56.46*</b>	<b>21.5</b>	<b>2.13</b>	43.44	<b>21.46*</b>	<b>58.15*</b>	<b>21.1</b>	2.32	<b>26.88*</b>	<b>60.08*</b>	<b>23.6</b>	1.19
Llama2-13b-chat <sub>sql</sub>	<b>2.48*</b>	17.88	49.18	19.1	2.12	46.35	13.34	52.55	20.9	2.32	17.17	54.33	21.7	<b>-0.71</b>
Llama2-13b-chat <sub>base</sub>	2.16	7.11	35.34	10.1	2.07	<b>49.70*</b>	11.17	45.13	13.5	<b>2.50</b>	15.13	49.20	12.29	8.02
Llama3-70B-Instruct <sub>func</sub>	6.52	<b>72.86*</b>	84.05	42.8	5.99	<b>68.73*</b>	<b>68.63*</b>	84.53	42.6	6.72	<b>72.35*</b>	85.0	45.9	<b>-0.51</b>
Llama3-70B-Instruct <sub>sql</sub>	6.50	51.05	<b>84.20</b>	45.9	6.01	51.62	50.40	85.21	44.5	6.67	52.67	85.51	45.1	1.62
Llama3-70B-Instruct <sub>base</sub>	<b>6.94*</b>	60.25	84.19	<b>47.5</b>	<b>6.58*</b>	65.80	61.60	<b>85.31</b>	<b>46.5</b>	<b>7.34*</b>	65.58	<b>86.32</b>	<b>48.6*</b>	5.33
Gorilla-v2 <i>few-shot DST</i>	4.91	54.86	78.52	38.7	4.48	77.07	56.16	79.95	33.8	5.36	59.46	81.44	33.8	4.6
Gorilla-v2 <i>0-shot DST</i>	5.07	61.07	66.64	33.0	4.51	31.93	60.25	69.52	34	5.36	64.54	71.24	33.2	3.47
GPT-3.5-turbo <i>0-shot DST</i>	6.36	51.08	71.01	39.7	5.51	22.40	51.0	72.22	41.0	6.46	58.61	74.52	36.2	7.53

Table 1: End-to-end evaluation with mean scores over 3 seeds. \* indicates statistical significance ( $p < 0.05$ , paired t-test) between best and second best values. For reference and context, rows in yellow show ours and previous supervised results: <sup>△</sup>Young et al. (2022), <sup>▲</sup>Stricker and Paroubek (2024), <sup>★</sup>Sun et al. (2023), <sup>◇</sup>Huang et al. (2023), <sup>♡</sup>Feng et al. (2023). Results in blue use few-shot prompting for DST with general purpose LLMs, which have not been explicitly adapted to the task of function calling. For comparison, rows in green use various models which have been tuned for function calling. Gorilla-v2 is used only for DST and relies on Llama3-8B-Instruct for response generation. Response generation is few-shot in all cases. Best results per Llama model size are in bold. We also show the relative drop in JGA between MultiWOZ and FusedChat results.

evidenced by its higher score on  $Acc_{fused}$ . In contrast, Llama3-8B-Instruct shows a propensity for classifying these turns as *task*-related, reflected in its higher  $Acc_{backstory}$  score. Llama3-70B-Instruct appears to strike a balance between these two tendencies. Overall we find this prompting approach to be functional, notably when it comes to interacting with the models in a live setting (Section 4). We plan nonetheless to refine this approach in future work.

Model	$Acc_{fused}$	$Acc_{interfere}$	$Acc_{backstory}$
Llama3-8B-Instruct	0.86	<b>0.98</b>	<b>0.91</b>
Llama2-13b-chat	0.85	0.81	0.45
Llama3-70B-Instruct	0.87	<b>0.98</b>	0.86
GPT-3.5-turbo	<b>0.89</b>	0.96	0.80

Table 2: Accuracy of zero-shot chitchat detection on both chitchat-enhanced benchmarks. *backstory* evaluates accuracy only on the InterfereChat turns that are augmented with contextual details, referred to as backstories.

**End-to-end evaluation** The results presented in Table 1 demonstrate that function calling consistently and generally outperforms other methods with statistical significance. This approach offers greater performance on inter-mode benchmarks and achieves the highest Joint Goal Accuracy across various model sizes. It is particularly noteworthy that Llama3-70B-Instruct surpasses previous supervised state-of-the-art JGA baselines on the benchmarks studied. We additionally find no statistically significant difference with Llama3-8B-Instruct<sub>finetune</sub> on the FusedChat benchmark. This indicates that this method is effective in retrieving task-information embedded in chitchat turns. We

notice similar Slot-F1 scores across benchmarks, showing the model can retrieve relevant information with all approaches, but performs best with function calls. This performance does come at a computational cost however, given the large scale of the model. Nonetheless, quantization strategies may be explored to improve latency and storage efficiency (Dettmers et al., 2024).

We observe that improved JGA does not necessarily translate to better task success, as no statistically significant difference is found in this regard.

BLEU-aug scores are generally high, mainly because the reference inter-mode responses from InterfereChat were generated with Llama2-70B-chat, a model from the same family. As we can see, GPT-3.5-turbo performs more poorly on this metric.

Llama3-8B achieves comparable and, in some cases, improved JGA and Slot-F1 scores compared to GPT-3.5-turbo, albeit requiring few-shot prompting instead of zero-shot. Gorilla-openfunctions-v2, a model specifically designed for function calling, performs better in a zero-shot setting than in a few-shot one on JGA, but not on Slot-F1. In the few-shot setting, its performance is on par with Llama3-8B.

The success of function calls with general LLMs can be attributed to two key factors: their widespread presence across various programming languages, which ensures their frequent occurrence in the code sections of LLM pre-training data, and their straightforward, easily producible syntax. These characteristics enable a more seamless conversion of user inputs into a structured format, where both the function name and its param-

eters closely align with the semantic content of the request. This semantic proximity suggests that chitchat enhancements may not be as disruptive as they might be with other approaches, provided they maintain coherence within the task context. In contrast, while SQL is also present in pre-training data, research by Tan et al. (2024) indicates that optimizing Text-to-SQL performance requires a more complex prompting strategy. This additional complexity may account for SQL’s comparatively lower performance in this scenario.

Model	Quality	Success	JGA	Clarify↓
Llama3-8B-Instruct	0.58	<b>90.0</b>	69.17	0.75
Llama3-70B-Instruct	<b>0.76</b>	<b>90.0</b>	<b>94.54</b>	<b>0.15</b>
GPT-3.5-turbo	0.70	85.0	76.25	0.65

Table 3: Human evaluation results. Quality is normalized ( $[0, 1]$  scale), success and JGA are percentages, and Clarify is the average number of reformulations needed per dialogue.

**Human evaluation** We also conduct a small in-house human evaluation, focusing on models capable of supporting the complete pipeline: Llama3-8B-Instruct, Llama3-70B-Instruct, and GPT-3.5-turbo with zero-shot DST. We randomly select 20 single-domain TOD goals from MultiWOZ, along with their corresponding backstories from InterfererChat. Four NLP experts were tasked with achieving these goals while engaging in chitchat with the models, mimicking the inter-mode scenarios from our selected benchmarks.

During our initial pilot annotation, we observed that participants struggled to simultaneously focus on the task goal and the contextual chitchat details. To address this, we introduced a seed turn to initiate the conversation and streamline the task. This seed turn could either be a task request incorporating contextual details or an open-ended chitchat utterance. When presented with the latter, participants engaged solely in chitchat for up to three turns before transitioning to the task goal.

Post-interaction, participants evaluated the dialogue quality (“Was the system friendly and engaging? (1-5)”) and success (“Was the desired goal reached? (Yes/No)”). We normalized the quality scores to a  $[0, 1]$  scale. For successful dialogues, we tallied the number of queries that needed repeating. Additionally, we manually calculated JGA by examining the dialogue state at each turn in relation to the user’s request. We included dialogue states during chitchat turns as well.

Overall, we find the resulting inter-mode interactions to be quite natural and successful, more

than automatic metrics indicate. As Table 3 indicates, Llama3-70B-Instruct performs well in terms of JGA and success rate, with only the rare clarification needed. This model is also preferred in terms of quality of the interaction, with more engaging chitchat and TOD responses. The smaller Llama3 model also does quite well, with more successful interactions than GPT-3.5-turbo.

The models demonstrate several positive attributes in handling inter-mode dialogues. All models show an ability to recall and incorporate chitchat details from the beginning of the dialogue when concluding the interaction, a feature that annotators particularly appreciated. Additionally, even when the wrong turn type is predicted, responses often remain coherent and contextual, allowing for the conversation not to break down. State-tracking-wise, the Llama3 models succeed at following function descriptions, accurately using the specified possible values even when users employ synonyms. This precision is crucial for successful database searches. These capabilities contribute to more natural and context-aware conversations, enhancing the overall user experience.

Despite these strengths, several issues were identified. Firstly, responses to inter-mode turns often focused heavily on the task request while neglecting contextual details. A more significant concern is the hallucination of novel entities (such as restaurant names). This issue can mislead users into believing a task was successfully completed when it was not. Lastly, unlike Llama3 models, GPT-3.5-turbo tends to directly extract the user’s words when performing state tracking, potentially leading to database lookup errors (*expensive* vs. *pricey*). Interaction examples with qualitative analyses can be found in Appendix C).

## 5 Conclusion

We show that a few-shot TOD baseline can be successfully extended to handle inter-mode inputs. We find that casting DST as function call generation is robust and effective across various LLMs, substantially outperforming other few-shot state-tracking methods. Human evaluation results show the promising potential of an inter-mode few-shot TOD bot, overall successfully balancing chitchat and TOD within a single interaction. These findings suggest a significant advancement in simply and swiftly building more versatile and natural dialogue systems.

## 6 Acknowledgments

This work was granted access to the HPC resources of IDRIS under the allocation 20XX-AD011014510 made by GENCI.

## References

- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Ian Beaver, Cynthia Freeman, and Abdullah Mueen. 2020. Towards awareness of human relational strategies in virtual agents. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(03):2602–2610.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Derek Chen, Kun Qian, and Zhou Yu. 2023. Stabilized in-context learning with pre-trained language models for few shot dialogue state tracking. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1551–1564, Dubrovnik, Croatia. Association for Computational Linguistics.
- Tim Dettmers, Ruslan A. Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. 2024. SpQR: A sparse-quantized representation for near-lossless LLM weight compression. In *The Twelfth International Conference on Learning Representations*.
- Yihao Feng, Shentao Yang, Shujian Zhang, Jianguo Zhang, Caiming Xiong, Mingyuan Zhou, and Huan Wang. 2023. Fantastic rewards and how to tame them: A case study on reward learning for task-oriented dialogue systems. *arXiv preprint arXiv:2302.10342*.
- James Gung, Emily Moeng, Wesley Rose, Arshit Gupta, Yi Zhang, and Saab Mansour. 2023. NatCS: Eliciting natural customer support dialogues. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9652–9677, Toronto, Canada. Association for Computational Linguistics.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A Simple Language Model for Task-Oriented Dialogue. In *Advances in Neural Information Processing Systems*, volume 33, pages 20179–20191. Curran Associates, Inc.
- 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*.
- Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu, Noah A. Smith, and Mari Ostendorf. 2022. In-context learning for few-shot dialogue state tracking. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2627–2643, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Tianjian Huang, Shaunak Ashish Halbe, Chinnadhurai Sankar, Pooyan Amini, Satwik Kottur, Alborz Geramifard, Meisam Razaviyayn, and Ahmad Beirami. 2023. Robustness through data augmentation loss consistency. *Transactions on Machine Learning Research*.
- Vojtěch Hudeček and Ondrej Dusek. 2023. Are large language models all you need for task-oriented dialogue? In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 216–228, Prague, Czechia. Association for Computational Linguistics.
- Pietro Lesci, Yoshinari Fujinuma, Momchil Hardalov, Chao Shang, Yassine Benajiba, and Lluís Marquez. 2023. Diable: Efficient dialogue state tracking as operations on tables. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 9697–9719, Toronto, Canada. Association for Computational Linguistics.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. API-bank: A comprehensive benchmark for tool-augmented LLMs. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3102–3116, Singapore. Association for Computational Linguistics.
- Zekun Li, Zhiyu Zoey Chen, Mike Ross, Patrick Huber, Seungwhan Moon, Zhaojiang Lin, Xin Luna Dong, Adithya Sagar, Xifeng Yan, and Paul A. Crook. 2024. Large language models as zero-shot dialogue state tracker through function calling. *Preprint*, arXiv:2402.10466.
- Andrea Madotto, Zhaojiang Lin, Genta Indra Winata, and Pascale Fung. 2021. Few-shot bot: Prompt-based learning for dialogue systems. *Preprint*, arXiv:2110.08118.
- 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.

Debjoy Saha, Bishal Santra, and Pawan Goyal. 2022. A study on prompt-based few-shot learning methods for belief state tracking in task-oriented dialog systems. *Preprint*, arXiv:2204.08167.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. *Toolformer: Language models can teach themselves to use tools*. In *Advances in Neural Information Processing Systems*, volume 36, pages 68539–68551. Curran Associates, Inc.

Armand Stricker and Patrick Paroubek. 2024. *Chitchat as interference: Adding user backstories to task-oriented dialogues*. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3203–3214, Torino, Italy. ELRA and ICCL.

Haipeng Sun, Junwei Bao, Youzheng Wu, and Xiaodong He. 2023. *Mars: Modeling context & state representations with contrastive learning for end-to-end task-oriented dialog*. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11139–11160, Toronto, Canada. Association for Computational Linguistics.

Zhao Tan, Xiping Liu, Qing Shu, Xi Li, Changxuan Wan, Dexi Liu, Qizhi Wan, and Guoqiong Liao. 2024. *Enhancing text-to-SQL capabilities of large language models through tailored promptings*. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 6091–6109, Torino, Italy. ELRA and ICCL.

Tom Young, Frank Xing, Vlad Pandealea, Jinjie Ni, and Erik Cambria. 2022. Fusing task-oriented and open-domain dialogues in conversational agents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11622–11629.

Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. *MultiWOZ 2.2 : A dialogue dataset with additional annotation corrections and state tracking baselines*. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 109–117, Online. Association for Computational Linguistics.

Xiaoying Zhang, Baolin Peng, Kun Li, Jingyan Zhou, and Helen Meng. 2023. *SGP-TOD: Building task bots effortlessly via schema-guided LLM prompting*. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13348–13369, Singapore. Association for Computational Linguistics.

## A Prompt Templates

We depict the various prompts created for our experiments. See Figure 4 for function calling DST, Figure 5 for SQL DST, Figure 6 for baseline DST, Figure 7 for response generation, Figure 8 for chitchat

detection, Figure 9 for task-oriented domain detection and Figure 10 for chitchat generation.

```

System Message
You are a task-oriented conversational AI assistant that helps users to book restaurants. Use the function definition below to create a function call with the correct arguments for the user's booking.

{"name": "find_book_restaurant",
 "description": "Find a restaurant and book a table",
 "parameters": {
  "type": "object",
  "properties": {
    "pricerange": {
      "type": "string",
      "description": "Price range of the restaurant",
      "possible_values": ["cheap", "moderate", "expensive"],
      "default_value": None,
    },
    "area": {
      "type": "string",
      "description": "Area where the restaurant is located",
      "possible_values": ["north", "east", "west", "south", "centre"],
      "default_value": None,
    },
  },
  "call_example": "find_book_restaurant(pricerange=None, area=centre, food='italian', name='pizza hut city centre', bookday='wednesday', booktimes='13:30', bookpeople=7)"
}

User Message
Output a function call with the correct function arguments given the customer's request. Make sure to follow the function definition. Focus only on the values mentioned in the last utterance.

-----Example 0:
Context:
Customer: I'm looking for an Italian restaurant for Saturday.
Assistant: <<function>>find_book_train(bookday='Saturday', food='Italian')

-----
Now complete the following example:
Context:
Customer: Hello, I am looking for a restaurant in Cambridge. I believe it is called Golden Wok.

```

Figure 4: Function calling DST prompt for the restaurant domain. Few-shot exemplars can be optionally added. We add 5 in all our few-shot experiments. We find that adding a *default\_value* field and one to several call examples helps improve performance.

## B End-to-end Training Details

We use the SimpleToD framework (Hosseini-Asl et al., 2020) to train an end-to-end TOD system, while relying only on a single language model. Each training example is composed of a concatenated text sequence which includes a dialogue context, the dialogue state for the last user turn, the database state (number of available options), response dialogue acts, and a delexicalized natural language response. We show an input example in Figure 11. Each training batch contains examples from each benchmark in equal proportion. During inference, we pass only the dialogue context to the model. Once the dialogue state is generated, we interrupt generation to fetch the database state, concatenate it to the current output and then continue generating.

We train the language model using LoRA (Hu et al., 2021), using a rank  $r$  of 64. We set the scaling  $\alpha$  to  $2r$ , and target the key, query, value and output weight matrices of the self-attention module. This amounts to roughly 1% of total parameters being trained. We use a learning rate of  $5e-5$ . We

**System Message**  
 You are a task-oriented conversational AI assistant that helps users to book restaurant. Using valid SQLite, answer the following multi-turn conversational questions for the table provided below.

```
CREATE TABLE restaurant(
name text,
food text,
pricerange text CHECK (pricerange IN (dontcare, cheap, moderate, expensive)),
area text CHECK (area IN (centre, east, north, south, west)),
booktime text,
bookday text,
bookpeople int
)
/*
5 example rows:
SELECT * FROM restaurant LIMIT 5;
name food pricerange area booktime bookday bookpeople
pizza hut city centre italian dontcare centre 13:30 wednesday 7
the missing sock international moderate east dontcare dontcare 2
golden wok chinese moderate north 17:11 friday 4
cambridge chop house dontcare expensive center 08:43 monday 5
darrys cookhouse and wine shop modern european expensive center 11:20
saturday 8
*/
```

**User Message**  
 Write a valid SQL query to extract the information from the Table given the customer's request. Make sure to end with a semicolon. Focus only on the values mentioned in the last utterance.

-----Example 0:  
 Context:  
 Customer: I'm looking for an Italian restaurant for Saturday.  
 Assistant: SELECT \* FROM restaurant WHERE bookday = saturday AND food = italian;

-----  
 Now complete the following example:  
 Context:  
 Customer: Hello, I am looking for a restaurant in Cambridge. I believe it is called Golden Wok.

Figure 5: SQL DST prompt for the restaurant domain. Few-shot exemplars can be optionally added. We add 5 in all our few-shot experiments. We follow (Hu et al., 2022) for this approach.

train for up to 2 epochs with early stopping on a single 80Gb A100.

## C Interactive Evaluation

We show our dialogue interface in Figure 12 and a few example interactions: Table 4, Table 5, Table 6, Table 7, show successful and unsuccessful interactions with the various models.

**User Message**  
 Capture entity values from last utterance of the conversation according to examples.  
 Focus only on the values mentioned in the last utterance.  
 Capture pair "entity:value" separated by colon and no spaces in between.  
 Separate entity:value pairs by hyphens.  
 Values that should be captured are:  
 - "pricerange" that specifies the price range of the restaurant (cheap/moderate/expensive)  
 - "area" that specifies the area where the restaurant is located (north/east/west/south/centre)  
 - "food" that specifies the type of food the restaurant serves  
 - "name" that specifies the name of the restaurant  
 - "bookday" that specifies the day of the booking  
 - "booktime" that specifies the time of the booking  
 - "bookpeople" that specifies for how many people is the booking made  
 Do not capture any other values!  
 If not specified, leave the value empty.

-----  
 -----Example 0:  
 Context:  
 Customer: I'm looking for an Italian restaurant for Saturday.  
 Assistant: bookday:'saturday'-food:italian'

-----  
 Now complete the following example:  
 Context:  
 Customer: Hello, I am looking for a restaurant in Cambridge. I believe it is called Golden Wok.

Figure 6: Baseline DST prompt for the restaurant domain. Few-shot exemplars can be optionally added. We add 5 in all our few-shot experiments. We follow (Hudeček and Dusek, 2023) for this approach. They do not use a system prompt in their implementation.

**User Message**  
 Definition: You are an assistant that helps people to book a restaurant. You can search for a restaurant by area, food, or pricerange. There is also a number of restaurants in the database currently corresponding to the user's request.  
 If multiple restaurants are available, the Assistant should ask for further preferences.  
 If you find a possible restaurant, the Assistant should provide [restaurant\_name], [restaurant\_address], [restaurant\_phone] or [restaurant\_postcode] if asked. Use these exact placeholders.  
 If no restaurants are available, the Assistant should ask for different preferences.  
 Before booking a table, the Assistant should ask for the time and the day of the booking and number of people. The Assistant should provide [reference] when the booking has been made. Use these exact placeholders.  
 \*\*Always act as if booking is available.\*\*  
 Write the Assistant response as a single line, based on the state and the database. Act friendly and engaging.

-----  
 -----Example 0:  
 Context:  
 Customer: I'm looking for an Italian restaurant for Saturday.  
 State: {restaurant:{food: italian, bookday: saturday}}  
 database:restaurant: 19  
 Assistant: <response> What area would you prefer ? </response>

-----  
 Now complete the following example:  
 Context:  
 Customer: Hello, I am looking for a restaurant in Cambridge. I believe it is called Golden Wok.  
 State: {restaurant: {name: Golden Wok }}  
 database:restaurant: 1  
 Assistant:

Figure 7: Response prompt for the restaurant domain. We use 5 few-shot examples for all experiments.

**System Message**  
 A user is using Cambridge's Towninfo Centre information assistant which can help users with information and bookings. These are tasks with certain specific domains.  
 Task domains include:  
 - train (booking train tickets)  
 - restaurant (finding and booking restaurants)  
 - hotel (finding and booking hotels)  
 - attraction (eg. "architecture", "sports", "entertainment", "cinema", "museum", "concert", "theatre"...)  
 - taxi (booking a taxi from one location to another)  
 You are an expert at determining if a User Turn contains task-related information or requests.

**User Message**  
 Use the Dialogue Context and the User Turn to classify the User Turn as "chitchat" or "task".  
 The criteria are:  
 If the User Turn contains:  
 - comments about personal life, opinions, or experiences  
 - casual comments about Cambridge or the domains (restaurants, trains, hotels, taxis, attractions)  
 then the User Turn is "chitchat".  
 If the User Turn contains  
 - an \*explicit\* request for information in a task domain (restaurant, train, hotel, attraction, taxi) or  
 - a request to perform an action in a task domain (restaurant, train, hotel, attraction, taxi) or  
 - an essential piece of information relevant to a task domain (restaurant, train, hotel, attraction, taxi)  
 then the User Turn is "task".

Dialogue Context:  
 User Turn:  
 Customer: I'm looking for an Italian restaurant for Saturday.

Respond with one word. Either "chitchat" or "task" between <label> and </label> tags.  
 Label:

Figure 8: Zero-shot chitchat detection prompt. This is meant to detect whether to perform chitchat or activate the few-shot TOD bot pipeline.

**System Message**  
 You are a conversational AI, capable of detecting the domain of a user request.  
 The only possible domains are :  
 - train (booking train tickets)  
 - restaurant (finding and booking restaurants)  
 - hotel (finding and booking hotels)  
 - attraction (eg. "architecture", "sports", "entertainment", "cinema", "museum", "concert", "theatre", "park", "church", "hotel", "nightclub", "swimming pool", "college", "concert hall", "boat", "historical landmark", "gallery", "shopping area", "nature preserve", "sports venue", "theme park")  
 - taxi (booking a taxi from one location to another)  
 Return only the domain name.

**User Message**  
 Select the appropriate domain given the customer's request. Respond with only one word: the domain name.  
 It is very important you focus on the customer's current request to make your decision.  
 Context:  
 Customer: I'm looking for an Italian restaurant for Saturday.  
 Domain:

Figure 9: Zero-shot domain selection prompt. This prompts the LLM to choose one of the possible MultiWOZ evaluation domains.

**System Message**  
 You are a friendly conversational AI assistant. Your goal is to engage with the user in a friendly conversation.

**User Message**  
 Respond to the user with a short response. Focus on being friendly and engaging. Write the response in between <response> and </response> tags.  
 Context:  
 Customer: I love to play video games !  
 Response:

Figure 10: Zero-shot chitchat prompt. This prompts the LLM to engage in a brief chitchat conversation with the user, in the context of a TOD interaction.

```
<context> (...) <user> I have always
wanted to try the food at a British
restaurant. <system> Yes, they always
look so fancy! <user> Are there any
expensive ones ? </context>
<belief> restaurant food british,
restaurant pricerange expensive
</belief> <db> [db_state4] </db>
<action> restaurant inform area,
restaurant inform area, restaurant inform
area, restaurant inform choice, restaurant
select </action> <response> Of course!
We have [choice] of them in the [area],
[area], and [area]. What area would you
prefer? </response>
```

Figure 11: Training example for our end-to-end fine-tuned approach, based on SimpleToD (Hosseini-Asl et al., 2020).

Figure 12: Dialogue interface used for our human evaluation. Participants use the seed turn to start the conversation and then follow the steps below. They then move to the QA tab to evaluate the dialogue.

User	I've been toying around with my new virtual reality headset lately.	<b>State</b>
Assistant	That sounds like a lot of fun! What kind of games have you been trying out on your VR headset?	{}
User	Well, do you know Mario? It's a fancy plumber from Japan, but he's Italian and speaks English.	
Assistant	Oh, I know Mario! He's a classic character. Have you been playing a Mario game on your VR headset?	{}
User	Well, duh, yeah. You've asked me specifically about VR. It's a bit odd in VR because it's a platformer, but it's a Mario game, still fun...	
Assistant	That's awesome! Mario in VR must be a whole new experience. It's great that you're enjoying it even in a different format. Have you found any cool new features in the VR version?	{}
User	In VR you get to hug Luigi, that's pretty good, you can't do that on a standard screen.	
Assistant	That's so cool that you get to hug Luigi in VR! It must feel more immersive and interactive. It's those little details that make the experience special. Have you discovered any other fun surprises in the VR version?	{}
User	Well, can't hug Luigi all day, I'm taking a break. I could need a taxi to get to ruskin gallery. I'm thinking about leaving from saffron brasserie after dinner with colleagues up to 07:15.	
Assistant	81418747547 is the contact number, I have you in a ford.	{"taxi": { "departure": "saffron brasserie", "destination": "ruskin gallery", "leaveat": "19:15"}}
User	Perfect, thank you.	
Assistant	It sounds like you had a great time with Mario in VR! Enjoy your dinner at Saffron Brasserie and have a fantastic time at the Ruskin Gallery. If you need any more help or information, feel free to ask!	{"taxi": { "departure": "saffron brasserie", "destination": "ruskin gallery", "leaveat": "19:15"}}

Table 4: An interaction with GPT-3.5-turbo where the assistant successfully handles chitchat and task-oriented user inputs. It also correctly tracks that *after dinner..07:15* should be mapped to 19:15 and closes the dialogue with a contextual response, that refers back to the initial chitchat.

User	I need a taxi asap ! I was eating at Pizza Hut City Centre when two men started fighting outside. The guards arrived and managed to break apart the fight. I think I'm ready to go home	<b>State</b>
Assistant	What a scary experience! Don't worry, I'm here to help. What's your destination, and I'll get a taxi booked for you ASAP?	{ "taxi": { "departure": "Pizza Hut City Centre" } }
User	I need to get to the gonville hotel	
Assistant	I've booked a taxi for you, an Executive will come for you. Contact number is 01223 312343. Do you need any other help?	{ "taxi": { "departure": "Pizza Hut City Centre", "destination": "gonville hotel", "leaveat": "asap" } }
User	That'll do	
Assistant	Ah, you must be relieved to be getting out of there! I hope the rest of your evening gets better. Your taxi should arrive soon, and you'll be safely back at the Gonville Hotel in no time. If you need anything else, don't hesitate to ask. Otherwise, have a great night and try to put that scary incident behind you!	{ "taxi": { "departure": "Pizza Hut City Centre", "destination": "gonville hotel", "leaveat": "asap" } }

Table 5: An interaction with Llama3-70B-Instruct where the assistant successfully handles an inter-mode turn (1st turn) by reacting empathetically while moving the task forward. Interestingly, regarding the dialogue state, we notice the model tracks "asap" as a time the taxi should leave, which makes sense given the context but may not work with a database expecting time values.

User	My grandmother thinks she can cook better and has never allowed me to eat at an Italian restaurant. But now she's away on a trip, so my whole family is going to an Italian restaurant to find out.	<b>State</b>
Assistant	Oh wow, that sounds like a fun adventure! I hope you all have a blast trying out different dishes and discovering your own favorite Italian flavors. And who knows, maybe your grandmother will be surprised by your culinary discoveries when she gets back! What's your favorite Italian dish so far?	{ }
User	I love lasagna, like Garfield the cat haha Could you suggest an Italian restaurant in the south ?	
Assistant	Ahaha, Garfield is a great reference! Lasagna is a classic choice, and there are so many amazing Italian restaurants in the South. I'd be happy to help you find one! Can you tell me which city or state you're in? That way, I can give you some personalized recommendations.	{ }
User	I need an italian restaurant in the south of Cambridge please	
Assistant	Ahaha, Garfield is a great reference! I've found a lovely Italian restaurant in the south of Cambridge that serves delicious lasagna. Bella Vita is a popular spot, located at G4 Cambridge Leisure Park Clifton Way Cherry Hinton and can be reached at 01223323737. Would you like to book a table?	{ "restaurant": "bookpeople": "1", "area": "south", "food": "italian" }
User	Yes please can you book for 7 people at 2 pm on Friday ? I know, we're quite a large group...	
Assistant	I've booked you in at Bella Vita. It's a lovely Italian restaurant in the south of Cambridge that serves delicious lasagna. The table will be reserved for 15 minutes. Reference number is : 00000004.	{ "restaurant": "bookpeople": "7", "area": "south", "food": "italian", "bookday": "fri- day", "booktime": "14:00" },
User	okay thank you bye	
Assistant	It was my pleasure to help! I hope you and your family have a fantastic time at Bella Vita and enjoy every bite of your lasagna. Don't forget to take plenty of photos and share them with your grandmother when she gets back. Have a great day and buon appetito!	{ "restaurant": { "bookpeople": "7", "area": "south", "food": "italian", "bookday": "fri- day", "booktime": "14:00" } }

Table 6: An example of an unsuccessful dialogue with Llama3-8B-Instruct. Although it seems fine on the surface, the Bella Vita does not exist in the database, it is a hallucination from the model. Other than that, the dialogue state tracking is accurate and the chitchat is well-handled, with the model finishing the dialogue with a few contextual remarks and *buon appetito*, which is a nice touch. In turn 2, the model classifies the user input as *chitchat*, which is a mistake. However, the response is contextual and still allows for the task to move forward.

		<b>State</b>
User	I love reading books. I find it interesting to read stories and get a peek into what the author is thinking.	
Assistant	Aww, that's wonderful! There's something so special about getting lost in a good book, isn't there? I'm curious, what kind of books do you usually enjoy reading? Are you more into fiction, non-fiction, or a mix of both?	{}
User	I like books on art and meditation. i read a good one by Rick Rubin recently, have you heard of him ?	
Assistant	Ah, Rick Rubin is amazing! I've heard great things about his book. I'm a big fan of his work with the Red Hot Chili Peppers and other artists. Meditation and art are such great topics to explore in books - it's wonderful that you're interested in those areas. Have you found that reading about meditation has helped you in your own practice?	{}
User	A bit, I'm not as consistent with it as I should be though	
Assistant	Ah, don't be too hard on yourself! Consistency can be tough, but it's great that you're acknowledging the importance of meditation in your life. Maybe we can brainstorm some ways to make it more enjoyable and sustainable for you? Do you have a favorite meditation spot or a specific technique that you enjoy? I'm all ears!	{}
User	I like sitting in the park usually and listning to the sounds around me	
Assistant	Ah, that sounds lovely! Being in nature can be so grounding and calming. I'm sure the sounds of the park help you focus and relax. Have you ever tried incorporating guided meditations into your practice, or do you prefer to just sit and listen to your surroundings?	{},
User	Yeah I'm looking apfr a book on that. But I have to go to Stevenage to get it. There's a bookstore there.	
Assistant	The byard art is the 1 museum in the south and the entrance fee is free.	{"attraction": { "type": "museum", "area": "south" }}
User	Hmm can you book me a train to Stevenage from Cambridge this Saturday, so I can go to the bookstore ?	
Assistant	What a lovely conversation we've had! I'd be delighted to help you book a train to Stevenage from Cambridge this Saturday. There is a train that will depart at 05:21 and arrive at 20:15. Would you like me to book that for you?	{"attraction": { "type": "museum", "area": "south" }, "train": { "bookpeople": "1", "arriveby": "20:15", "day": "saturday", "destination": "stevenage" }}
User	I need one that arrives by 7:15 pm.	
Assistant	There are 10 trains matching your query. Shall I book one?	{"attraction": { "type": "museum", "area": "south" }, "train": { "bookpeople": "1", "arriveby": "19:15", "day": "saturday", "destination": "Stevenage", "departure": "Cambridge" }}

Table 7: An extract from an interaction with Llama3-8B-Instruct, where a slightly ambiguous user typo leads the model to incorrectly try to find a museum and start tracking the *attraction* domain. We notice though that this is not a problem for continuing with the task of booking a train which the model tracks accurately. We also find the chitchat quite convincing, with the model displaying good knowledge of the music producer Rick Rubin.

# Exploration of Human Repair Initiation in Task-oriented Dialogue : A Linguistic Feature-based Approach

Anh Ngo<sup>1, 6</sup>, Dirk Heylen<sup>2</sup>, Nicolas Rollet<sup>3</sup>, Catherine Pelachaud<sup>4</sup>,  
Chloé Clavel<sup>1, 5</sup>

<sup>1</sup>ALMAnaCH, INRIA Paris, <sup>2</sup>University of Twente, The Netherlands,

<sup>3</sup>Télécom Paris, SES, Institut Polytechnique de Paris, I3-CNRS, <sup>4</sup>CNRS, ISIR, Sorbonne University,

<sup>5</sup>Télécom Paris, LTCI, Institut Polytechnique de Paris, <sup>6</sup>ISIR, Sorbonne University

Correspondence: [anh.ngo-ha@inria.fr](mailto:anh.ngo-ha@inria.fr)

## Abstract

In daily conversations, people often encounter problems prompting conversational repair to enhance mutual understanding. By employing an automatic coreference solver, alongside examining repetition, we identify various linguistic features that distinguish turns when the addressee initiates repair from those when they do not. Our findings reveal distinct patterns that characterize the repair sequence and each type of other-repair initiation.

## 1 Introduction

Human language complexities often expose flaws such as misunderstandings, misinterpretations, speech impediments, or social norm violations. Strategies people use in conversations to identify and address these problems, fostering mutual understanding, are called repair (Schegloff, 2007). Schegloff, 2007 distinguishes repair types based on who initiates and who provides the solution between the speaker and the addressee. This paper focuses on Other-initiated Self-repair, also called Other-initiated repair (OIR), where the addressee initiates repair for the speaker corrects, as highlighted by Dingemans and Enfield, 2024 as foundational for human language resilience, complexity, and flexibility.

Recent studies emphasize the need for Conversational Agents (CAs) to have repair mechanisms. Gehle et al., 2014 show that museum guide robots failing to promptly address issues led to visitor disengagement and conversation breakdowns, suggesting the importance of multimodal repair strategies. van Arkel et al., 2020 find that simple OIR mechanisms in agents improve communicative success and reduce computational and interaction costs for disambiguation in communication compared to pragmatic reasoning (interlocutors reason each other). Efforts to detect OIR in the literature are narrow. Purver et al., 2018 trained a supervised

classifier on four different datasets using turn-level features extracted from transcription, such as numbers of wh-words and fillers. The results indicate that challenging repairs are more common in task-oriented datasets. Besides, research on integrating OIR in CAs is limited and primarily relies on rule-based systems. For instance, Höhn, 2017 developed a rule-based chatbot with repair capabilities that recognize repair initiation in messaging conversations using conversational analysis rules, such as repetition, determiner and pronoun usage, and adjacent position.

**Example 1.** Sample of OIR sequence, annotated based on Dingemans and Enfield, 2015's coding schema. Data is in Dutch, English translation provided by DeepL<sup>1</sup>.

**TS SPEAKER:** en ik zie een uh  
ovaalvormig ding op het kopje (T-1)  
(and I see a uh oval-shaped thing on the  
cup)

**REPAIR INITIATOR:** op het platte kopje  
daarboven hè? (T0)  
(on the flat head up there huh?)

**TS SPEAKER:** ja (T+1)  
(yes)

A minimal OIR sequence comprises three components: trouble source (TS) turn (T-1), repair initiation (T0), and repair solution (T+1), depicted in Example 1. T-1 is where a potential communication problem arises, T0 is where the addressee signals a problem, and T+1 is where the speaker resolves the problem, completing the repair sequence. In addition, Dingemans and Enfield, 2015 categorized repair initiation into three types: *Open Request* (the least specific, no TS specified in T-1), *Restricted Request* (implying the TS location), and *Restricted Offer* (the most specific, proposing a candidate understanding). Our research aims to develop a CA

<sup>1</sup><https://www.deepl.com/>

system that can detect human repair initiation (T0) based on verbal and non-verbal cues and generate an appropriate repair solution (T+1). This work examines dialogue transcripts to identify linguistic features that distinguish OIR sequences from non-repair sequences and differentiate among the three types of OIR in task-oriented dialogues. Previous studies have identified various OIR practices, for instance, Schegloff et al., 1977 described five OIR formats in their study of American English conversation, while Dingemanse et al., 2014 find similarities in OIR formats across ten languages, such as question word "what?" or interjection "huh." The contributions of this paper are as follows: First, unlike previous studies that focused solely on references to trouble sources (TS), this work expands further by examining the acceptance of repair initiation by subsequent turn. Second, in addition to repetition, we incorporate an automatic coreference solver to see if repair initiators refer back to the TS and if the response acknowledges the repair initiation. Results show significant coreference involvement in *restricted request* and *restricted offer*.

## 2 Dataset

As repair occurs more often in task-oriented dialogue and is generally unaffected by familiarity or interaction mode (Colman and Healey, 2011), we employ dialogue transcripts from a Dutch multimodal task-oriented corpus (Rasenberg et al., 2022) within project CABB (Eijk et al., 2022), involving 20 dyads performing referential communication tasks to locate 16 stimulated geometrical objects called Fribbles. The data collection setup corresponds to the CABB dataset, described in (Eijk et al., 2022). Participants alternated between Director and Matcher roles to communicate and locate specified objects. Each participant’s speech was segmented into Turn Constructional Units (TCUs) and then orthographically transcribed based on standard spelling conventions of Dutch. The repair sequences were annotated following Dingemanse and Enfield, 2015’s coding schema, resulting in: 20 (*open request*), 32 (*restricted request*), and 255 (*restricted offer*) sequences, respectively.

We examine the interaction differences after a potential issue (turn labeled as TS) to compare instances when a person initiates repair versus when they do not identify trouble and request repair immediately. To do this, we selected all turns between T-1 and the repair initiation in T0, identifying 91

non-repair sequences. Appendix A provides sample data for each OIR type and details the non-repair selection method.

## 3 Feature Extraction

Based on the OIR coding schema (Example 1), T0 is considered a repair initiation if it (1) treats the prior turn containing trouble and (2) the subsequent turn T+1 acknowledges and responds to this request. To determine (1), we analyze the syntactic structure of the repair initiation turn (T0) regarding the potential TS in T-1 via coreferences. For (2), we examine how the TS speaker acknowledges the repair initiation by analyzing coreferences in T+1 that refer to the entity mentioned in T0 and the TS speaker’s self-repetition.

### 3.1 Feature extraction for repair initiation (T0) concerning prior turn (T-1)

**Part-of-Speech (POS) tagging and Lemmatization.** To investigate T0’s linguistic patterns across three OIR types compared to T0 in non-repair sequences, we leverage Stanza<sup>2</sup>, a multilingual NLP toolkit, for POS tagging and lemmatization. It enables us to comprehend the overall grammatical structure of the OIR turn and identify the most frequently used word types and their corresponding most common lemmas. The performance of Stanza’s pretrained model on Dutch is 94.97% for POS tagging and 95.33% for lemmatization. See Appendix C for the list of POS tags in Dutch.

**Coreference.** Coreference, a linguistic phenomenon in dialogue, involves referring to entities across turns using pronouns, demonstratives, or other expressions linked to previously mentioned nouns or concepts. Analyzing coreference patterns offers insights into the relationships between turns. By examining coreferences used by the repair initiator in T0, we investigate if T0 refers to an entity in T-1, potentially the TS. We utilized the coreference resolution model from the UTD\_NLP team (Li et al., 2022), which achieved the best performance at CODI-CRAC 2022 (Yu et al., 2022), with an average CONLL F1 score of 75.04 in resolving anaphora in dialogue. The coreference chain sample produced by the model is included in the Appendix B.

To analyze repair initiation structure and its grammatical ties to prior turns via coreferences,

<sup>2</sup><https://stanfordnlp.github.io/stanza/>

we used Seq2Pat<sup>3</sup>, a sequence-pattern-generation library (Kadioğlu et al., 2023). Each turn T0 after tokenization and POS tagging was fed into Seq2Pat to obtain a list of the most frequent sequential patterns. The instances of the coreference chain are tagged by [COREF]. Due to data imbalance among OIR types, different min\_frequency thresholds were set to each: min\_frequency = [5, 5, 30, 10] for *open request*, *restricted request*, *restricted offer*, and non-repair, respectively.

### 3.2 Feature extraction for repair acknowledgment in subsequent turn (T+1)

To determine if the TS speaker’s response in T+1 addresses the request made by the repair initiator in T0, we analyzed the coreferences initiated in T0 and used in T+1. We also examined the TS speaker’s self-repetition behavior when providing a repair solution, as these repetitions suggest the TS speaker’s language consistency and alignment with the trouble in T-1. To identify self-repetition, we used dialign<sup>4</sup>, a tool for measuring lexical alignment in human-agent interaction (Dubuisson Duplessis et al., 2017) (example in Appendix B).

## 4 Results and Discussion

### 4.1 Does T0 consider the prior turn T-1 as source of trouble?

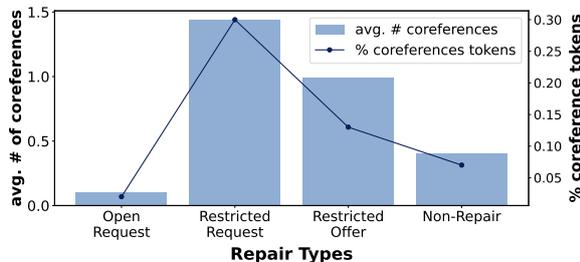


Figure 1: T0’s average number of coreferences and coreferences tokens proportion

Figure 1 shows the average number of coreferences in T0 (initiated in T-1) and the percentage of T0 tokens that are coreferences. *Restricted request* has the highest coreference usage (about 1.5 coreferences per T0, comprising approximately 30% of tokens), followed by *restricted offer* (around one coreference per T0, accounting for approximately 13% of tokens). Non-repair and *open request* show minimal coreference use, with about 0.5 coreferences (7% of tokens) and 0.1 coreferences (2%

<sup>3</sup><https://github.com/fidelity/seq2pat>

<sup>4</sup><https://github.com/GuillaumeDD/dialign>

of tokens) per T0, respectively. Both *restricted request* and *restricted offer* signal trouble in T-1, likely indicating dependence on coreferences for previously mentioned ambiguous entities. However, *restricted offer*, potentially introducing new entities and a longer T0 turn to propose candidate understanding, explains the lower coreference usage and proportion of coreference tokens compared to *restricted request*.

Figure 2 describes the most common sequential POS tag patterns for T0 across three OIR types and non-repair sequences. These patterns, displayed as bi-grams ["1<sup>st</sup> POS tag", "2<sup>nd</sup> POS-tag", frequency], are visualized on the y-axis, x-axis, and heatmap values, respectively. Additionally, Figure 3 depicts the top five most frequent POS tags and their corresponding lemmas.



Figure 2: T0 utterances POS tags Sequential Patterns

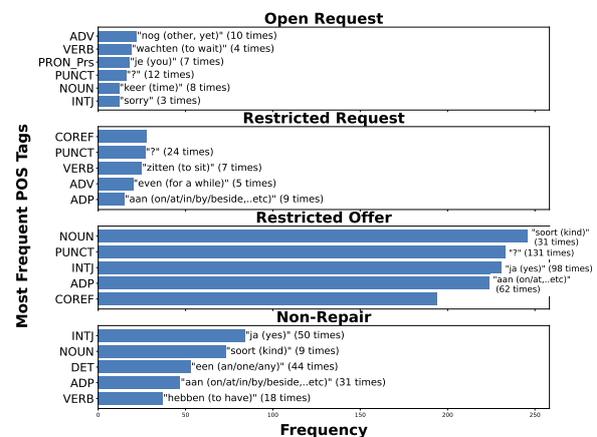


Figure 3: T0’s most frequent POS tags and its corresponding most frequent lemma

Regarding *open request*, frequent patterns in-

volve adverbs preceded by personal pronouns (10 times), verbs (7 times), interjections (6 times), and auxiliaries (5 times). The most common adverb lemma is "nog" (*yet*) expressing negation, while the personal pronoun "je" (*you*) suggests a request towards the prior turn's speaker. Compared to the other OIR types, *open request* uniquely involve auxiliaries, with the modal verb "kunnen" (*be able to/can/may*) being the most frequent, indicating the request for action from the prior speaker. Notably, verbs frequently found in these patterns, with the most common lemma being "wachten" (*to wait*), may indicate a request to slow down due to issues in the previous turn.

In *restricted request*, the notable correlation between coreferences and other word forms like verbs (16 times), adverbs (15 times), and prepositions (14 times) indicates heavy reliance on referring back to previously mentioned entities. Using interrogative pronouns (PRON\_Int) is a distinctive feature in this type, often followed by verbs (7 times, most frequently "zitten" (*to sit*)), prepositions (5 times, most commonly "aan" - equivalent to multiple English prepositions like *on, at, in, by, beside*), and personal pronouns (6 times, most frequently "je" (*you*)). Its most common lemma, "wat" (*what, which, any*), is used for asking questions, indicating a demand for clarification from the current speaker regarding what the prior speaker mentioned (potentially TS).

Considering *restricted offer*, the most frequent sequential patterns involve determiners followed by nouns (121 times) and coreferences (54 times). The sequences combining prepositions preceded or followed by determiners (60 or 80 times, respectively), nouns (58 or 97 times, respectively), or coreferences (60 or 67 times, respectively) are also common patterns. These patterns emphasize the scenario where the repair initiator is likely presenting or describing specific objects to offer the candidate understanding.

Non-repair sequences' T0 share similarities to *restricted offer* regarding the usage of noun phrases, determiners, and prepositions. However, the presence of adjectives sets it apart from all three OIR types, implying a focus on descriptive presentation. Especially in non-repair, there is a high occurrence of the combination of demonstrative pronouns with auxiliaries, determiners, and nouns, emphasizing the introduction of new entities, clarification, or stating existence rather than extensive reference back. Unlike *open request*, the auxiliary verb "zijn"

(*to be*) is the most frequent in non-repair sequences, often employed for demonstration.

Unique sequential POS tag patterns and particular behaviors in employing coreference chains to refer to entities from the preceding turn reveal that each OIR type initiates repair requests differently, setting them apart from non-repair sequences. Utilizing these extracted patterns could assist in creating a repair initiation detector for CA.

#### 4.2 Does subsequent turn T+1 acknowledge T0's request as repair initiation?

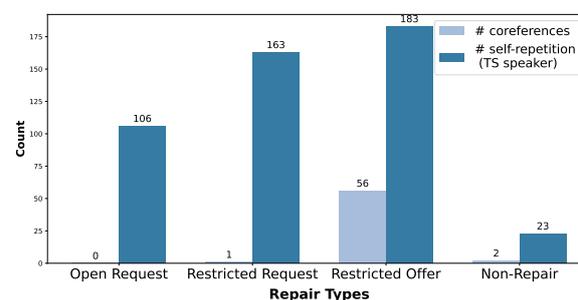


Figure 4: Distribution of Coreference (initiated in T0) and TS Speaker Self-repetition in T+1

Figure 4 examines T+1's acceptance of repair initiation from T0, showing the number of coreferences (initiated in T0) used in T+1 and the TS speaker's self-repetition (verb and noun) from T-1 to T+1. Regarding coreference, only in *restricted offer*, the TS speaker in T+1 uses several coreferences to refer to the entities initiated by the repair initiator in T0, unlike *open request* and *restricted request* where coreferences are rare. Since the *restricted offer* is the most specific repair initiation, it potentially prompts the TS speaker to use coreferences for confirming the proposed candidate.

In contrast, the high self-repetition across all OIR types suggests the TS speaker often repeats themselves (from T-1) to address the repair initiation request. Despite occasional similarities with repair initiation or a format resembling OIR, the infrequent use of coreferences and self-repetition in non-repair sequences suggests that the TS speaker potentially progressed the conversation without acknowledging it as a request for repair.

These patterns, particularly in *restricted offer*, could enhance repair solution generation models in CA by incorporating them with repair initiation sequential patterns.

## 5 Conclusion and Future Work

Utilizing Natural Language Processing approaches on dialogue transcripts, we identified linguistic and sequential patterns characterizing three types of OIR and non-repair sequences. The coreference chains used in T0 combined with sequential patterns of OIR structure are typical across OIR types and non-repair sequences, which reveal the grammatical structure of T0 and whether T0 treats the prior turn T-1 as containing trouble. Besides, the TS speaker's self-repetition and coreference chains (initiated by the repair initiator) used in T+1 show the behavior of the TS speaker in acceptance of the repair initiation from T0. Our future work will explore multimodalities like prosodic, facial and bodily cues, to develop a computational model for repair initiation detection and repair solution generation in the Conversational Agent.

## Acknowledgments

Data were provided (in part) by the Radboud University, Nijmegen, The Netherlands. This work has been supported by the Paris Île-de-France Région in the framework of DIM AI4IDF. This work was partially funded by the ANR-23-CE23-0033-01 SINNet project.

## References

- Marcus Colman and Patrick G. T. Healey. 2011. [The distribution of repair in dialogue](#). *Cognitive Science*, 33.
- Mark Dingemanse, Joe Blythe, and Tyko Dirksmeyer. 2014. [Formats for other-initiation of repair across languages: An exercise in pragmatic typology](#). *Studies in Language*, 38.
- Mark Dingemanse and N. J. Enfield. 2015. [Other-initiated repair across languages: Towards a typology of conversational structures](#).
- Mark Dingemanse and N. J. Enfield. 2024. [Interactive repair and the foundations of language](#).
- Guillaume Dubuisson Duplessis, Chloé Clavel, and Frédéric Landragin. 2017. [Automatic measures to characterise verbal alignment in human-agent interaction](#). In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 71–81, Saarbrücken, Germany. Association for Computational Linguistics.
- Lotte Eijk, Marlou Rasenberg, Flavia Arnese, Mark Blokpoel, Mark Dingemanse, Christian F. Doeller, Mirjam Ernestus, Judith Holler, Branka Milivojevic, Asli Özyürek, Wim Pouw, Iris van Rooij, Herbert Schriefers, Ivan Toni, James Trujillo, and Sara Bögels. 2022. [The cabb dataset: A multimodal corpus of communicative interactions for behavioural and neural analyses](#). *NeuroImage*, 264.
- Raphaela Gehle, Karola Pitsch, and Sebastian Benjamin Wrede. 2014. [Signaling trouble in robot-to-group interaction: emerging visitor dynamics with a museum guide robot](#). *Proceedings of the second international conference on Human-agent interaction*.
- Sviatlana Höhn. 2017. [A data-driven model of explanations for a chatbot that helps to practice conversation in a foreign language](#). In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 395–405, Saarbrücken, Germany. Association for Computational Linguistics.
- Serdar Kadioğlu, Xin Wang, Amin Hosseininasab, and Willem-Jan Hoeve. 2023. [Seq2pat: Sequence-to-pattern generation to bridge pattern mining with machine learning](#). *AI Magazine*, 44.
- Shengjie Li, Hideo Kobayashi, and Vincent Ng. 2022. [Neural anaphora resolution in dialogue revisited](#). In *Proceedings of the CODI-CRAC 2022 Shared Task on Anaphora, Bridging, and Discourse Deixis in Dialogue*, pages 32–47, Gyeongju, Republic of Korea. Association for Computational Linguistics.
- Matthew Purver, Julian Hough, and Christine Howes. 2018. [Computational models of miscommunication phenomena](#). *Topics in Cognitive Science*, 10(2):425–451.
- Marlou Rasenberg, Wim Pouw, Asli Özyürek, and Mark Dingemanse. 2022. [The multimodal nature of communicative efficiency in social interaction](#). *Scientific Reports*, 12.
- Emanuel A. Schegloff. 2007. *Sequence organization in interaction : a primer in conversation analysis I*. Cambridge University Press.
- Emanuel A. Schegloff, Gail Jefferson, and Harvey Sacks. 1977. [The preference for self-correction in the organization of repair in conversation](#). *Language*, 53:361.
- Jacqueline van Arkel, Marieke Woensdregt, Mark Dingemanse, and Mark Blokpoel. 2020. [A simple repair mechanism can alleviate computational demands of pragmatic reasoning: simulations and complexity analysis](#). In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 177–194, Online. Association for Computational Linguistics.
- Juntao Yu, Sopan Khosla, Ramesh Manuvinaurike, Lori Levin, Vincent Ng, Massimo Poesio, Michael Strube, and Carolyn Rosé. 2022. [The CODI-CRAC 2022 shared task on anaphora, bridging, and discourse deixis in dialogue](#). In *Proceedings of the*

## A Sample Data

### Example 2. Open request OIR sample

**TS SPEAKER:** op dat driehoek (T-1)  
(*on that triangle*)

**REPAIR INITIATOR:** wat zei je? (T0)  
(*what did you say?*)

**TS SPEAKER:** op die driehoek (T+1)  
(*on that triangle*)

### Example 3. Restricted request OIR sample

**TS SPEAKER:** deze heeft twee oren die aan de onderkant breder worden en een soort hanekam op zijn hoofd een kleintje (T-1)

(*this one has two ears that widen at the bottom and a sort of cock's comb on its head a little one*)

**REPAIR INITIATOR:** maar wat zei wat zei je in het begin? (T0)  
(*but what did you say at the beginning?*)

**TS SPEAKER:** een soort oren die aan de onderkant breder worden (T+1)  
(*a kind of ears that widen at the bottom*)

### Example 4. Restricted offer OIR sample

**TS SPEAKER:** waarbij je dus op de bovenkant zo'n zo'n mini uh kegeltje hebt (T-1)  
(*where you have one of those mini uh cones on the top*)

**REPAIR INITIATOR:** oh ja die zo scheef naar achter staat? (T0)  
(*oh yes which is so slanted backwards?*)

**TS SPEAKER:** ja precies (T+1)  
(*yes exactly*)

### Example 5. Non-repair sequence selection example.

**TS SPEAKER:** het is een een een een een kopje en aan de onderkant zit een uh ovale standaard zeg maar het kopje staat daarop -> **TS**  
(*it's a a a a a little cup and at the bottom there's an uh oval stand, let's say, the cup stands on it*)

**REPAIR INITIATOR:** zit er een driehoek op? -> **Non-repair**  
(*is there a triangle on it?*)

**TS SPEAKER:** nee er zit geen driehoek op -> **Non-repair**  
(*no there is no triangle on it*)

**REPAIR INITIATOR:** en het staat zeg maar op zo'n ovale ding ja zo'n pilvorm is het -> **OIR**  
(*and it stands on such an oval thing, yes such a pill shape is it*)

**TS SPEAKER:** ja het belangrijkste is dat het een soort van houder heeft waar het op staat zeg maar -> **Repair solution**  
(*yes the most important thing is that it has some kind of holder that it stands on, let's say*)

**Example 6.** Non-repair sequence sample (2). The second turn resembles an OIR format with the repair initiator repeating "drie bolletjes" (three balls). However, it is not considered OIR because the TS speaker continues with new information in the subsequent turn, indicating they saw it as acknowledgment rather than a repair request.

**TS SPEAKER:** oh ja deze heeft uh drie bolletjes telkens als armen  
(*oh yes this one has uh three balls each as arms*)

**REPAIR INITIATOR:** drie bolletjes  
(*three balls*)

**TS SPEAKER:** en staat op een groot vierkant  
(*and stands on a large square*)

## B Coreference Chain and Self-repetition Samples

**Example 7.** Coreference used by Repair Initiator in T0, initiated by TS Speaker in T-1

**TS SPEAKER:** um dit is de hoofdvorm met die ronde staaf aan de linkerkant die uitgesneden is met die punt erin (T-1)  
(*um this is the main shape with that round bar on the left cut out with that point in it*)

**REPAIR INITIATOR:** um je bedoelt met die schuine punt zo naar beneden? (T0)  
(*um you mean with that slant point so down?*)

**TS SPEAKER:** ja (T+1)  
(*yes*)

**Example 8. Coreference** used by TS Speaker in T+1, initiated by Repair Initiator in T0

TS SPEAKER: soort plakseltjes ofzo (T-1)

*(kind of sticky or something)*

REPAIR INITIATOR: ja lijkt een beetje op een stopcontact zou kunnen zo'n stekker? (T0)

*(yes looks a bit like a socket could be such a plug?)*

TS SPEAKER: ja je zou het in een stopcontact kunnen zetten (T+1)

*(yes you could put it in a socket)*

- DET - deteminers
- INTJ - interjection
- NOUN - noun
- PRON\_Dem - demonstrative pronouns
- PRON\_Int - interrogative pronouns
- PRON\_Prs - personal pronouns
- PUNCT - punctuations
- SYM - symbols
- VERB - verbs

**Example 9. Sample of TS speaker's self-repetition**

TS SPEAKER: dit is de hoofdvorm waarbij een yoghurtbakje links aan de hoofdvorm vastzit soort van klein staafje rechts en dan bovenop een rechthoekige staaf (T-1)

*(this is the main form where a yoghurt container on the left is attached to the main form kind of small bar on the right and then on top a rectangular bar)*

REPAIR INITIATOR: yoghurtbakje was? (T0)

*(yoghurt container was?)*

TS SPEAKER: ja yoghurtbakje op de kop links van de hoofdvorm zit er aan vastgeplakt (T+1)

*(yes yoghurt tray on the head left of the main form is stuck to it)*

## C POS tags List

- ADJ - adjectives
- ADP - prepositions and postpositions
- ADV - adverbs
- AUX - auxiliaries, including
  - perfect tense auxiliaries "hebben" (*to have*), "zijn" (*to be*)
  - passive tense auxiliaries "worden" (*to become*), "zijn" (*to be*), "krijgen" (*to get*)
  - modal verbs "kunnen" (*to be able, can*), "zullen" (*shall*), "moeten" (*must*), "mogen" (*to be allow*)
- CCONJ - coordinating conjunctions "en" (*and*), "of" (*or*)

# Comparing Pre-Trained Embeddings and Domain-Independent Features for Regression-Based Evaluation of Task-Oriented Dialogue Systems

Kallirroi Georgila

Institute for Creative Technologies, University of Southern California  
12015 Waterfront Drive, Los Angeles, CA 90094-2536, USA  
kgeorgila@ict.usc.edu

## Abstract

We use Gaussian Process Regression to predict different types of ratings provided by users after interacting with various task-oriented dialogue systems. We compare the performance of domain-independent dialogue features (e.g., duration, number of filled slots, number of confirmed slots, word error rate) with pre-trained dialogue embeddings. These pre-trained dialogue embeddings are computed by averaging over sentence embeddings in a dialogue. Sentence embeddings are created using various models based on sentence transformers (appearing on the Hugging Face Massive Text Embedding Benchmark leaderboard) or by averaging over BERT word embeddings (varying the BERT layers used). We also compare pre-trained embeddings extracted from human transcriptions with pre-trained embeddings extracted from speech recognition outputs, to determine the robustness of these models to errors. Our results show that overall, for most types of user satisfaction ratings and advanced/recent (or sometimes less advanced/recent) pre-trained embedding models, using only pre-trained embeddings outperforms using only domain-independent features. However, this pattern varies depending on the type of rating and the embedding model used. Also, pre-trained embeddings are found to be robust to speech recognition errors, more advanced/recent embedding models do not always perform better than less advanced/recent ones, and larger models do not necessarily outperform smaller ones. The best prediction performance is achieved by combining pre-trained embeddings with domain-independent features.

## 1 Introduction

The quality of a human-machine dialogue interaction can be influenced by various factors, such as the domain/genre of dialogue, the dialogue system capabilities, and the user expertise and expectations. This makes it very difficult to define what a

successful dialogue should look like, and evaluate system performance and predict user satisfaction. Thus, despite many years of research, dialogue evaluation still remains an unsolved problem.

In this paper, our focus is on task-oriented dialogue, and specifically on predicting user satisfaction after their interaction with the dialogue system. We use the Communicator corpus (Walker et al., 2001a, 2002) containing the logs of user interactions with 8 spoken dialogue systems. The user's task is to book a flight and in some cases also make hotel or car-rental arrangements. Each dialogue log is accompanied by user ratings after their interaction with the system. An example dialogue excerpt is shown in Figure 1 in the Appendix.

The original Communicator corpus contains system and user utterances (both human transcriptions and speech recognition outputs), timing information, and speech act and task annotations for the system's side of the conversation. An extended version of this corpus was developed by Georgila et al. (2005b, 2009) via automatic annotation. Georgila et al. (2005b, 2009) added speech act and task annotations for the user's side of the conversation, and dialogue context annotations, e.g., filled slots, filled slots values, grounded slots, speech acts history.

In this paper, we use Gaussian Process Regression for predicting user satisfaction ratings, because in our recent work (Georgila, 2022) it was shown to perform better than other regression methods, for this task and corpus. In our previous work (Georgila, 2022), we considered only domain-independent features (e.g., duration, number of filled slots, number of confirmed slots, word error rate). These features were domain-independent because they were just based on counts, and no lexical, semantic, or specific to the task information was used. Here, in addition to these domain-independent features, we also use pre-trained dialogue embeddings extracted from system and user utterances.

Our pre-trained dialogue embeddings are computed by averaging over sentence embeddings for each dialogue. Sentence embeddings are created using various models based on sentence transformers (Reimers and Gurevych, 2019) (appearing on the Hugging Face Massive Text Embedding Benchmark leaderboard), or by averaging over BERT word embeddings (Wieting et al., 2016; Coates and Bollegala, 2018) (varying the BERT layers used). By definition, these embeddings are domain-dependent because they encode lexical and semantic information about the domain. Also, we compare pre-trained embeddings extracted from human transcriptions versus pre-trained embeddings extracted from automatic speech recognition (ASR) outputs, to determine the robustness of these models to errors, which is an understudied research question (Mousavi et al., 2024). We investigate what level of performance can be achieved just by relying on the words of the system and user utterances from which we compute pre-trained dialogue embeddings, whether using only embeddings outperforms using only domain-independent features, and whether combining embeddings and domain-independent features can result in performance gains. We also examine the impact on performance of different feature combinations.

To our knowledge, our work is one of a few studies (if not the first) to compare such a large variety of pre-trained embeddings (including the most recent embedding models by OpenAI) under the same conditions, and the first study to do so for predicting user ratings in task-oriented dialogue. This is also the first work to compare all these different types of pre-trained embeddings with various domain-independent features for user ratings’ prediction in task-oriented dialogue. Last, but not least, this is one of a very limited number of studies comparing the performance of pre-trained embeddings on human transcriptions versus ASR outputs, and the first study to do so for user ratings’ prediction.

## 2 Related Work

Despite many years of research, dialogue evaluation still remains an unsolved problem (Hastie, 2012; Deriu et al., 2021; Mehri et al., 2022). For task-oriented dialogue there are subjective evaluation metrics, such as user satisfaction, computed using information from surveys (Hone and Graham, 2000; Paksima et al., 2009), and objective metrics, such as task completion and dialogue length, com-

puted using information from interaction logs.

PARADISE (Walker et al., 2000) is the most well-known framework for automatic evaluation of task-oriented dialogue. The goal of PARADISE is to optimize user satisfaction (or another desired quality) by formulating it as a linear combination of various factors, such as task success and dialogue cost (e.g., dialogue length, ASR errors). Weights calculated via linear regression determine the contribution of each factor. PARADISE can be used to predict user satisfaction at the end of the dialogue, but can also be applied to any point in the dialogue prior to completion. Generally it is useful to be able to evaluate on the fly how the dialogue is unfolding, so that appropriate measures can be taken (e.g., transfer to a human operator), if a dialogue is problematic. Based on this idea, much work has been done on estimating user satisfaction at the system-user exchange level rather than rating the whole dialogue (Engelbrecht et al., 2009; Higashinaka et al., 2010; Ultes and Minker, 2014; Schmitt and Ultes, 2015).

For chatbots and other non-task-oriented dialogue systems it is not clear what success means, and it is common to use subjective evaluations of system responses (e.g., coherence, engagingness) given some context, or use word-overlap similarity metrics (e.g., BLEU, ROUGE) even though such metrics do not correlate well with human judgments of dialogue quality (Liu et al., 2016). Recently, new evaluation metrics have been proposed for open-domain dialogue leveraging pre-trained language models such as BERT and DialoGPT (Mehri and Eskenazi, 2020a,b; Ghazarian et al., 2020), and commonsense knowledge bases (Ghazarian et al., 2023).

In this paper, we focus on predicting user satisfaction ratings for the whole dialogue. We use Gaussian Process Regression (GPR) for predicting user satisfaction ratings, because in our recent work (Georgila, 2022) it was shown to perform better than other regression methods, for this task and corpus. In our previous work (Georgila, 2022), we only used domain-independent features, but here we also use pre-trained dialogue embeddings extracted from system and user utterances.

Linear regression has been used before for dialogue evaluation (Walker et al., 2000, 2001b; Cervone et al., 2018; Georgila et al., 2019, 2020; Georgila, 2022). Also, Support Vector Regression has been used before for dialogue evaluation (Cervone et al., 2018; Georgila, 2022).

We use GPR for our experiments because modern regression methods are a natural evolution of the PARADISE framework. Furthermore, we do not have many data points for data-hungry methods such as neural networks. As we will see in section 3, we only have 500 data points in the training data and 506 data points in the test data.

### 3 Data and Features

We use the Communicator corpus (Walker et al., 2001a, 2002) because it has been used before for this task, but also because it is one of a few task-oriented dialogue corpora that include user ratings. Other popular corpora, such as MultiWOZ (Budzianowski et al., 2018; Eric et al., 2020), do not include user ratings or ASR outputs.

The original Communicator corpus contains system and user utterances (both human transcriptions and ASR outputs), timing information, and speech act and task annotations for the system’s side of the conversation based on the DATE scheme (Walker and Passoneau, 2001). An extended version of this corpus was developed by Georgila et al. (2005b, 2009) via automatic annotation. Based on the ASR outputs, speech act and task annotations for the user’s side of the conversation were added, as well as dialogue context annotations, e.g., filled slots, filled slots values. Basically these extended annotations are the kind of information one would get by deploying a dialogue system, but because the original corpus did not include such information, Georgila et al. (2005b, 2009) reconstructed it.

Georgila et al. (2009) verified the validity and reliability of these automatic annotations by evaluating them with respect to the task completion metrics of the original corpus and in comparison to manually annotated data. The utility of these extended annotations has been demonstrated by their use by various researchers for different purposes, such as learning dialogue policies (Henderson et al., 2005; Frampton and Lemon, 2006; Henderson et al., 2008; McLeod et al., 2019) and building simulated users (Schatzmann et al., 2005; Georgila et al., 2005a, 2006).

In the Appendix, Figure 1 shows an example dialogue excerpt including speech act and task annotations, and Figure 2 depicts an example dialogue state.

These extended dialogue context annotations are divided into two broad categories: logs of the current status of the slots (‘FilledSlotsStatus’, ‘Filled-

SlotsValuesStatus’, ‘GroundedSlotsStatus’), and logs containing information about how the status of the slots has changed over time through the dialogue (‘FilledSlotsHist’, ‘FilledSlotsValuesHist’, ‘GroundedSlotsHist’). The former inform us about the current status of the slots, and may only contain one instance per slot. The latter provide information about the order in which slots have been filled or confirmed, and may contain several instances of the same slot. The annotations also include the history of speech acts and tasks.

For our experiments we use the 2001 collection, which consists of 1,683 dialogues between human users and 8 dialogue systems. These systems vary in their dialogue policies, e.g., some of them request multiple pieces of information at the same time, others request explicit confirmation, others request implicit confirmation, etc. Overall there are 78,718 turns (39,419 system turns and 39,299 user turns). Similarly to Georgila (2022), for our experiments we only used dialogues for which all user ratings were available: ATT (157 dialogues), BBN (137 dialogues), CMU (69 dialogues), COLORADO (157 dialogues), IBM (77 dialogues), LUCENT (140 dialogues), MIT (166 dialogues), and SRI (103 dialogues). The first half of the dialogues from each system is used for training (500 dialogues in total) and the rest for testing (506 dialogues in total).

So our task is to predict the following user satisfaction ratings on a Likert scale (1-5, higher is better): ease of the tasks the user had to accomplish (‘Task-Ease’), whether it was easy or not to understand the system (‘System-Comprehend-Ease’), the user’s expertise (‘User-Expertise’), whether the system behaved as expected (‘System-Behaved-As-Expected’), and if the user would use the system again in the future (‘System-Future-Use’). We use the same domain-independent features as Georgila (2022), with the addition of the number of times the user requested a ‘start-over’. Our 17 domain-independent features are divided into 4 categories:

- **duration-related features (9):** overall duration, duration of the system talking part, duration of the user talking part, overall average duration per utterance, average duration per system utterance, average duration per user utterance, number of overall speech acts, number of system speech acts, number of user speech acts;
- **slots-related features (6):** number of filled

slots, number of filled slots without any ‘null’ values, number of grounded slots, number of filled slots in the dialogue history, number of filled slots without any ‘null’ values in the dialogue history, number of grounded slots in the dialogue history (all at the end of the dialogue) – we distinguish between slots filled with normal versus ‘null’ values as an extra piece of information;

- **word error rate (WER) (1)**: calculated as the edit distance between the ASR output and the transcription of the user utterance (this information was included in the original Communicator corpus);
- **start-over feature (1)**: number of ‘start-over’ requests by the user extracted from the human transcription or the ASR output.

All these features are automatically extracted from the data. Feature values are replaced with z-scores by subtracting from each feature value the mean for that feature and then dividing by the standard deviation for that feature. For each feature, the mean and standard deviation are calculated on the training data.

We use 4 variations of these feature combinations: ‘orig-man’ (original corpus with features from manual annotations such as human transcriptions of speech plus fully automatic annotations), ‘orig-auto’ (original corpus with fully automatic annotations), ‘ext-man’ (extended corpus with features from manual annotations plus fully automatic annotations), and ‘ext-auto’ (extended corpus with fully automatic annotations). So ‘ext-man’ is a super set of ‘orig-man’, and ‘ext-auto’ is a super set of ‘orig-auto’, because the extended corpus contains all the annotations of the original corpus plus new annotations (note that, as mentioned above, these new annotations are automatically generated). Also, ‘orig-man’ and ‘ext-man’ include both manual and automatic annotations, whereas ‘orig-auto’ and ‘ext-auto’ include only automatic annotations.

For duration, the number of user speech acts is only used in ‘ext-man’ and ‘ext-auto’, because (as discussed above) the original corpus did not include annotations of the user’s side of the conversation. Likewise, slots-related features are only part of the extended corpus (‘ext-man’ and ‘ext-auto’). Information about WER is only part of the manual annotations because it can be computed only when human transcriptions are available.

	orig -man	orig -auto	ext -man	ext- auto
<b>duration</b>	x	x	x	x
<b>slots</b>			x	x
<b>WER</b>	x		x	
<b>start-over</b>	x	x	x	x

Table 1: Categories of feature combinations; x means that a feature category is included.

For clarity, Table 1 shows exactly which features are used in each category.

We also compute pre-trained dialogue embeddings by averaging over sentence embeddings for each dialogue. Sentence embeddings are created using various models based on sentence transformers (appearing on the Hugging Face MTEB leaderboard), or by averaging over BERT word embeddings (varying the BERT layers used). We do not calculate z-scores for the embeddings.

We use the following types of embeddings from Hugging Face and OpenAI, and in parentheses we can see the sizes of the vectors they produce:

- ‘glove-6B-300d’ (300) (Pennington et al., 2014),
- ‘all-distilroberta-v1’ (768),
- ‘all-mpnet-base-v2’ (768),
- ‘all-MiniLM-L6’ (384),
- ‘all-MiniLM-L12’ (384),
- ‘e5-small-v2’ (384),
- ‘e5-base-v2’ (768),
- ‘e5-large-v2’ (1024) (Wang et al., 2024),
- ‘gte-small’ (384),
- ‘gte-base’ (768),
- ‘gte-large’ (1024),
- ‘bge-small-en-v1.5’ (384),
- ‘bge-base-en-v1.5’ (768),
- ‘bge-large-en-v1.5’ (1024),
- OpenAI’s ‘text-embedding-3-small’ (1536),
- OpenAI’s ‘text-embedding-3-large’ (3072).

The latest models of OpenAI have a new feature that allows selecting the size of the generated vector. According to OpenAI, this compressed vector retains its concept-representing properties. For ‘text-embedding-3-small’ we experimented with 3 vector sizes (50, 256, 1536) and for ‘text-embedding-3-large’ with 3 vector sizes (50, 512, 3072).

Because we only have 500 data points in the training data and 506 data points in the test data, and large vector sizes, we also applied Principal Component Analysis (PCA) for dimensionality reduction, with “whitening” to ensure that the resulting features are less correlated with each other. Huang et al. (2021) and Su et al. (2021) have found that “whitening” can enhance the isotropy of sentence embeddings, with the additional advantage of reducing their dimensionality.

We generated results with different numbers of PCA components, and we show results with a value of 50 which performed well for all models (better e.g., than 75 or 100). Of course, when we generated vectors of size 50 from OpenAI, we did not apply PCA. Note that we apply PCA only to the embedding vectors (the domain-independent features are not affected by PCA).

## 4 Experiments and Results

In our previous work (Georgila, 2022), we compared several state-of-the-art regression methods, and showed that GPR with an exponential kernel or a rational quadratic kernel performed the best. Thus, here we use GPR with an exponential kernel. Also, by performing more experiments, we verified again that GPR outperforms other regression methods, and that using an exponential kernel produces competitive results for different types of embeddings. For all GPR experiments we vary the length scale, and we report results for length scale equal to 1 (higher length scale values indicate smoother learned functions). Varying the length scale did not produce significant differences. GPR is considered as the state-of-the-art for regression, and has been used before in the NLP community for machine translation quality estimation (Cohn and Specia, 2013) and emotion prediction (Beck et al., 2014). For all our experiments we use the GPy library<sup>1</sup>, and GPR is applied after PCA.

To evaluate our models, for each of the 5 ratings, we calculate the Root Mean Square Error (RMSE). RMSE measures the average error be-

tween the model predictions and the ground truth (the ratings in the test data). Its value varies from 0 to 4, given that user ratings are on a scale from 1 to 5. Lower RMSE values are better.

### 4.1 Using Only Pre-Trained Embeddings

Table 2 shows results in terms of RMSE when using only our embedding models (not including domain-independent features), based on human transcriptions (‘man’) and ASR outputs (‘auto’).

For BERT, we experimented with various layer combinations, and we report the best results. We found that it helps to use the first layer (L1) together with the last layers (L10, L11, L12). Other researchers have also looked into the impact of different BERT layers, reporting that sometimes it is better not to use the last layer, as it is largely fine-tuned to the specific task (Li et al., 2020; Huang et al., 2021; Su et al., 2021). Although differences were small, the best layer combination was L1-10-11 which means that the vectors of layers L1, L10, and L11 were averaged. Layer L1 alone also produced competitive results. We hypothesize that layer L1 is important for our task because it encodes lexical information rather than semantic meaning, and for dialogue evaluation some words such as “start-over” or “no” can be quite predictive.

For each BERT layer, we also compared averaging of word embeddings versus using the output of the [CLS] token, and averaging performed better. Thus, here we only present results with averaging (see Table 5 in the Appendix for a comparison between averaging and using the output of the [CLS] token). Reimers and Gurevych (2019) showed that averaging of BERT word embeddings or using the output of the [CLS] token produces rather poor sentence embeddings, often worse than averaging GloVe word embeddings (even though BERT word embeddings are generally considered superior to GloVe word embeddings). However, this was not the case in our experiments where BERT most times (depending on the layer combination) worked better than GloVe.

Overall, ‘glove-6B-300d’, ‘all-distil-roberta-v1’, ‘all-mpnet-base-v2’, ‘all-MiniLM-L6’, and ‘all-MiniLM-L12’ did not perform well compared to the rest of the models. For ‘e5’, ‘gte’, and ‘bge’, the small versions performed well, and that was the case also for the large version of ‘bge’. This is interesting because it shows that larger models do not necessarily perform better than smaller models. The question that arises is whether larger models

<sup>1</sup><https://gpy.readthedocs.io/en/deploy/>

	TaskEase		SysComEase		UserExp		SysBehExp		SysFutUse	
	man	auto								
bert-L1	1.259	1.28	<b>1.163</b>	1.174	<b>1.287</b>	1.303	1.291	1.324	<b>1.363</b>	1.379
bert-L1-10-11	<b>1.254</b>	1.256	1.173	<b>1.167</b>	1.302	<b>1.289</b>	<b>1.283</b>	1.285	1.366	1.371
bert-L1-11-12	1.255	1.256	1.174	1.173	1.302	1.298	1.288	1.287	1.364	1.373
bert-L1-10-11-12	1.255	<b>1.253</b>	1.173	1.169	1.305	1.291	1.285	<b>1.28</b>	1.367	<b>1.37</b>
glove-6B-300d	1.294	1.286	1.177	1.176	1.296	1.294	1.33	1.321	1.385	1.389
all-distilroberta-v1	1.28	1.296	<b>1.171</b>	<b>1.174</b>	1.306	1.309	1.317	1.325	1.388	1.392
all-mpnet-base-v2	<b>1.271</b>	<b>1.278</b>	1.179	1.181	1.29	<b>1.289</b>	<b>1.305</b>	<b>1.317</b>	<b>1.377</b>	<b>1.379</b>
all-MiniLM-L6	1.276	1.282	1.172	<b>1.174</b>	<b>1.287</b>	1.295	1.312	1.318	1.381	1.387
all-MiniLM-L12	1.273	1.282	1.188	1.187	<b>1.287</b>	1.291	1.318	1.321	1.389	1.388
e5-small-v2	<b>1.254</b>	1.272	<b>1.175</b>	1.183	<b>1.285</b>	<b>1.286</b>	1.311	1.316	<b>1.376</b>	1.38
e5-base-v2	1.281	1.291	1.195	1.192	1.303	1.29	1.32	1.333	1.39	1.387
e5-large-v2	1.27	<b>1.271</b>	1.176	<b>1.181</b>	1.293	1.302	<b>1.299</b>	<b>1.308</b>	1.378	<b>1.375</b>
gte-small	<b>1.26</b>	<b>1.264</b>	<b>1.168</b>	<b>1.175</b>	<b>1.285</b>	<b>1.283</b>	1.303	1.306	<b>1.37</b>	1.373
gte-base	1.284	1.294	1.186	1.189	1.3	1.302	1.323	1.326	1.385	1.387
gte-large	1.265	1.27	1.172	1.178	1.29	1.286	<b>1.299</b>	<b>1.3</b>	1.374	<b>1.369</b>
bge-small-en-v1.5	<b>1.262</b>	<b>1.261</b>	1.181	1.185	1.296	1.29	1.299	1.31	1.374	1.372
bge-base-en-v1.5	1.281	1.283	1.186	1.186	1.308	1.304	1.315	1.311	1.392	1.386
bge-large-en-v1.5	<b>1.262</b>	1.271	<b>1.168</b>	<b>1.175</b>	<b>1.284</b>	<b>1.282</b>	<b>1.285</b>	<b>1.292</b>	<b>1.363</b>	<b>1.367</b>
openai-small-50	1.323	1.32	1.191	1.191	1.305	1.306	1.35	1.344	1.399	1.398
openai-small-256	1.31	1.297	1.183	1.183	1.302	1.299	1.347	1.326	1.396	1.387
openai-small-1536	1.263	1.268	1.161	1.163	1.285	<b>1.284</b>	1.313	1.312	1.38	1.374
openai-large-50	1.27	1.286	<b>1.149</b>	<b>1.158</b>	1.303	1.298	<b>1.3</b>	<b>1.308</b>	1.37	1.379
openai-large-512	<b>1.258</b>	<b>1.266</b>	1.161	1.165	<b>1.283</b>	1.285	1.303	1.312	1.362	1.361
openai-large-3072	1.264	1.268	1.164	1.166	1.297	1.29	1.31	1.317	<b>1.36</b>	<b>1.36</b>

Table 2: RMSE values when using only pre-trained embeddings (not including domain-independent features), based on the human transcriptions (‘man’) and the ASR outputs (‘auto’). For each block, the best value for each column is shown in a different color (specific to that block) and in bold. **The best value for each column across all blocks is shown in black and in bold.**

were negatively affected by being compressed more than smaller models, given that we used only 50 PCA components. However, as we see with the OpenAI models, this is not the case. The ‘openai-large-3072’ model was significantly compressed and yet performed well. When we experimented with different numbers of components the trends were the same, i.e., the small versions of ‘e5’, ‘gte’, and ‘bge’ still worked better than their base and large counterparts, with the exception of ‘bge’ where the large version also performed well.

For the OpenAI models, we can see that the models based on ‘text-embedding-3-large’ worked better than the models based on ‘text-embedding-3-small’. Interestingly, ‘openai-large-50’ works very well. Note that this is the model where the compression was done by OpenAI (not by our using of PCA). It is not clear what kind of dimensionality reduction algorithm OpenAI uses. For some ratings, we can see that applying PCA on ‘openai-

large-512’ and ‘openai-large-3072’ works better than ‘openai-large-50’.

Overall, differences in results across models are small, but there are trends:

- Larger models are not necessarily better than smaller models.
- More advanced/recent models do not always perform the best.
- Pre-trained embeddings are quite robust to ASR errors for our task, given that differences in RMSE values between corresponding ‘man’ and ‘auto’ models are small.

#### 4.2 Comparing Pre-Trained Embeddings and Domain-Independent Features

Table 3 shows the full results for the best performing embedding models from Table 2. So, for example, ‘orig-em-man’ means manual and automatic

	bert-L1	bert-L1-10-11	e5-small-v2	gte-small	bge-small-en-v1.5	bge-large-en-v1.5	openai-large-50	openai-large-512	openai-large-3072
<b>Task-Ease</b>									
orig-man: 1.292 ext-man: 1.276 orig-auto: 1.311 ext-auto: 1.284									
em-man	1.259	<b>1.254</b> <sup>†</sup>	<b>1.254</b> <sup>†</sup>	1.26	1.262	1.262	1.27	1.258	1.264
orig-em-man	1.236	1.242	<b>1.235</b> <sup>†</sup>	1.238	1.24	1.249	1.244	1.241	1.246
ext-em-man	1.235	1.241	<b>1.233</b> <sup>†</sup>	1.237	1.239	1.249	1.245	1.241	1.245
em-auto	1.28	<b>1.256</b> <sup>‡</sup>	1.272	1.264	1.261	1.271	1.286	1.266	1.268
orig-em-auto	1.253	1.245	1.248	<b>1.236</b> <sup>‡</sup>	1.237	1.258	1.256	1.248	1.25
ext-em-auto	1.252	1.244	1.244	<b>1.233</b> <sup>‡</sup>	1.235	1.257	1.256	1.247	1.248
<b>System-Comprehend-Ease</b>									
orig-man: 1.174 ext-man: 1.156 orig-auto: 1.178 ext-auto: 1.158									
em-man	1.163	1.173	1.175	1.168	1.181	1.168	<b>1.149</b> <sup>†</sup>	1.161	1.164
orig-em-man	1.138	1.152	1.15	1.141	1.152	1.149	<b>1.134</b> <sup>†</sup>	1.143	1.147
ext-em-man	1.136	1.151	1.148	1.14	1.15	1.148	<b>1.133</b> <sup>†</sup>	1.143	1.147
em-auto	1.174	1.167	1.183	1.175	1.185	1.175	<b>1.158</b> <sup>‡</sup>	1.165	1.166
orig-em-auto	1.148	1.15	1.16	1.149	1.154	1.155	<b>1.14</b> <sup>‡</sup>	1.149	1.15
ext-em-auto	1.147	1.149	1.157	1.148	1.152	1.153	<b>1.139</b> <sup>‡</sup>	1.15	1.15
<b>User-Expertise</b>									
orig-man: 1.286 ext-man: 1.295 orig-auto: 1.286 ext-auto: 1.293									
em-man	1.287	1.302	1.285	1.285	1.296	1.284	1.303	<b>1.283</b>	1.297
orig-em-man	1.276	1.296	1.274	<b>1.268</b> <sup>†</sup>	1.278	1.272	1.284	1.269	1.281
ext-em-man	1.282	1.305	1.278	<b>1.274</b> <sup>†</sup>	1.284	1.279	1.289	<b>1.274</b> <sup>†</sup>	1.286
em-auto	1.303	1.289	1.286	1.283	1.29	<b>1.282</b>	1.298	1.285	1.29
orig-em-auto	1.288	1.283	1.27	<b>1.262</b> <sup>‡</sup>	1.269	1.268	1.28	1.271	1.275
ext-em-auto	1.294	1.288	1.272	<b>1.266</b> <sup>‡</sup>	1.273	1.273	1.286	1.275	1.279
<b>System-Behaved-As-Expected</b>									
orig-man: 1.301 ext-man: 1.278 orig-auto: 1.33 ext-auto: 1.286									
em-man	1.291	<b>1.283</b>	1.311	1.303	1.299	1.285	1.3	1.303	1.31
orig-em-man	1.268	1.269	1.282	1.271	1.27	<b>1.267</b> <sup>†</sup>	<b>1.267</b> <sup>†</sup>	1.282	1.288
ext-em-man	1.262	1.262	1.273	1.264	1.263	<b>1.259</b> <sup>†</sup>	1.26	1.274	1.279
em-auto	1.324	<b>1.285</b> <sup>‡</sup>	1.316	1.306	1.31	1.292	1.308	1.312	1.317
orig-em-auto	1.291	<b>1.266</b> <sup>‡</sup>	1.287	1.274	1.277	1.273	1.273	1.29	1.294
ext-em-auto	1.282	<b>1.259</b> <sup>‡</sup>	1.278	1.265	1.267	1.265	1.266	1.281	1.284
<b>System-Future-Use</b>									
orig-man: 1.397 ext-man: 1.394 orig-auto: 1.41 ext-auto: 1.395									
em-man	1.363	1.366	1.376	1.37	1.374	1.363	1.37	1.362	<b>1.36</b> <sup>†</sup>
orig-em-man	<b>1.339</b> <sup>†</sup>	1.353	1.357	1.345	1.35	1.346	1.36	1.344	1.342
ext-em-man	<b>1.337</b> <sup>†</sup>	1.348	1.351	1.341	1.347	1.344	1.36	1.339	1.339
em-auto	1.379	1.371	1.38	1.373	1.372	1.367	1.379	1.361	<b>1.36</b> <sup>‡</sup>
orig-em-auto	1.358	1.355	1.36	1.346	1.347	1.351	1.362	1.344	<b>1.343</b> <sup>‡</sup>
ext-em-auto	1.354	1.35	1.354	1.341	<b>1.34</b> <sup>‡</sup>	1.347	1.362	<b>1.34</b> <sup>‡</sup>	<b>1.34</b> <sup>‡</sup>

Table 3: RMSE values for different combinations of embedding models and features. The best value for each row (i.e., best model) is shown in a color specific to that rating and in bold. **The best value for each rating is shown in black and in bold**; <sup>†</sup> means that ‘em-man’, ‘orig-em-man’, or ‘ext-em-man’ are significantly better than either ‘orig-man’ or ‘ext-man’ ( $p < 0.05$  or better); <sup>‡</sup> means that ‘em-auto’, ‘orig-em-auto’, or ‘ext-em-auto’ are significantly better than either ‘orig-auto’ or ‘ext-auto’ ( $p < 0.05$  or better). Also, ‘em-auto’ means only embeddings from ASR outputs, ‘ext-em-man’ means manual and automatic annotations from the extended corpus plus embeddings from human transcriptions, etc.

annotations from the original corpus ('orig-man') plus embeddings extracted from human transcriptions, 'ext-em-auto' means only automatic annotations from the extended corpus ('ext-auto') plus embeddings extracted from ASR outputs, 'em-man' means only embeddings extracted from human transcriptions, etc. Here, for each rating, we also see results using only the domain-independent features without embeddings; these results are slightly different from the results reported by Georgila (2022) because we additionally use the 'start-over' feature.

For all ratings, we measure statistical significance between the best values of 'em-man/orig-em-man/ext-em-man', and either 'orig-man' or 'ext-man'. Sometimes, the difference between 'em-man/orig-em-man/ext-em-man' and 'orig-man' is significant, but the difference between 'em-man/orig-em-man/ext-em-man' and 'ext-man' is not significant (or vice versa). In this case, we still mark the difference as significant in Table 3 (to avoid over-crowding Table 3 with too many different markings). We also measure statistical significance between the best values of 'em-auto/orig-em-auto/ext-em-auto', and either 'orig-auto' or 'ext-auto'. We mark differences as significant in the same way as explained above.

For all statistical significance calculations, for comparing models, we use the squared error values and the Wilcoxon signed-rank test with Holm-Bonferroni correction for repeated measures. We did not test for significance all combinations, but roughly a difference in the RMSE values of 0.02 or larger is likely to be significant at  $p < 0.05$  or better (depending on the variance of course).

For 'Task-Ease' the best models are 'bert-L1-10-11', 'e5-small-v2', and 'gte-small'. Using embeddings (with or without domain-independent features), based on human transcriptions ('man') or ASR outputs ('auto'), outperforms using only domain-independent features.

For 'System-Comprehend-Ease', the best model is 'openai-large-50' for all feature combinations. There are significant differences between the RMSE values of this model and the RMSE values of the domain-independent features.

For 'User-Expertise', the best models are 'gte-small', 'bge-large-en-v1.5', and 'openai-large-512'. Here differences between using only domain-independent features and using only embeddings are small and not significant, but they become significant once we combine domain-independent features and embeddings.

	Bas 3	Bas maj	BM
<b>Task-Ease</b>	1.471	1.721	<b>1.233</b>
<b>Sys-Compr-Ease</b>	1.421	1.285	<b>1.133</b>
<b>User-Expertise</b>	1.431	1.41	<b>1.262</b>
<b>Sys-Behave-Exp</b>	1.433	1.705	<b>1.259</b>
<b>Sys-Future-Use</b>	1.516	2.321	<b>1.337</b>

Table 4: RMSE values for the baseline always predicting score 3 (Bas 3) and the majority baseline (Bas maj), and the best of our models (BM). **The best value for each row is shown in bold.**

For 'System-Behaved-As-Expected', the best models are 'bert-L1-10-11', 'bge-large-en-v1.5', and 'openai-large-50'. For the best model ('bert-L1-10-11') and for 'auto' using only embeddings significantly outperforms using only domain-independent features. Combining domain-independent features and embeddings results in significant differences for both 'man' and 'auto'.

For 'System-Future-Use', the best models are 'bert-L1', 'bge-small-en-v1.5', 'openai-large-512', and 'openai-large-3072'. Using only embeddings performs much better than using only domain-independent features ( $p < 0.01$  for 'man' and  $p < 0.001$  for 'auto'). Combining domain-independent features and embeddings further improves performance ( $p < 0.001$  for both 'man' and 'auto').

Similarly to Georgila (2022), we also implemented simple baselines. Table 4 shows results for RMSE for each type of rating, for the baseline that always predicts score 3, the majority baseline, and the best result of our models taken from Table 3. As expected, our models significantly outperform the baselines ( $p < 0.001$ ).

Below we summarize our findings from comparing pre-trained embeddings with domain-independent features:

- For most types of user satisfaction ratings and advanced/recent pre-trained embedding models, using only pre-trained dialogue embeddings outperforms using only domain-independent features.
- Combining pre-trained embeddings and domain-independent features is better than just using pre-trained embeddings.
- Differences between corresponding 'man' and 'auto' models are small, and thus, we conclude that pre-trained dialogue embeddings are quite robust to ASR errors for our task.

- Using domain-independent features from the extended annotations sometimes helps, but overall, performance is similar to using features from the original annotations.

## 5 Conclusions and Discussion

We used GPR for predicting user satisfaction ratings. We used both domain-independent features and pre-trained dialogue embeddings extracted from system and user utterances. Our pre-trained dialogue embeddings were computed by averaging over sentence embeddings for each dialogue. Sentence embeddings were created using various models based on sentence transformers (appearing on the Hugging Face MTEB leaderboard) or by averaging over BERT word embeddings (varying the BERT layers used).

Our results showed that overall, for most types of user satisfaction ratings and advanced/recent pre-trained embedding models, using only pre-trained dialogue embeddings outperforms using only domain-independent features. This is very interesting, because it shows that we can do quite well relying only on information from words and system and user utterances, without any additional features. Combining embeddings and domain-independent features performed the best. This is also very interesting and could potentially revive interest in using domain-independent features. Although overall extracting domain-independent features from the extended annotations helped, performance was similar to using domain-independent features from the original annotations.

Interestingly, some simpler models (e.g., ‘bert-L1’) performed better than more complex and more recent models. Also, larger models did not necessarily outperform smaller ones. Because differences between corresponding ‘man’ and ‘auto’ models were small, we conclude that pre-trained embeddings are quite robust to ASR errors for our task. Overall, RMSE values ranged roughly from 1.1 to 1.4 depending on the model and feature combination.

Our overall contributions are as follows:

- To our knowledge, our work is one of a few studies (if not the first) to compare such a large variety of pre-trained embeddings (including the most recent embedding models by OpenAI) under the same conditions.
- Our work is the first study to compare such

a large variety of pre-trained embeddings (including the most recent embedding models by OpenAI) under the same conditions for predicting user ratings in task-oriented dialogue.

- Our work is also the first study to compare all these different types of pre-trained embeddings and various domain-independent features for user ratings’ prediction in task-oriented dialogue.
- Finally, this is one of a very limited number of studies comparing the performance of pre-trained embeddings on human transcriptions versus ASR outputs, and the first study to do so for user ratings’ prediction.

Throughout our experiments, to construct dialogue embeddings we used averaging (Wieting et al., 2016; Coates and Bollegala, 2018), but the problem with averaging is that it can result in loss of important conversational information (Reimers and Gurevych, 2019). For example, not all parts of a dialogue are of equal importance, and by trying to encode everything we may end up compressing too much information from parts that really matter.

Very little work has been done on constructing dialogue embeddings using techniques different from averaging. A notable recent attempt to construct dialogue embeddings is Dial2vec (Liu et al., 2022). Dial2vec uses self-guided contrastive learning (leveraging both positive and negative examples) and considers a dialogue as an information exchange process between two interlocutors. It learns embeddings for both interlocutors with the help of each other, and then the dialogue embedding is obtained by an aggregation of embeddings of the interlocutors. Dial2vec was used to construct dialogue embeddings for the tasks of domain categorization, semantic relatedness, and dialogue retrieval. Based on the idea of Dial2vec, an interesting future research direction would be to learn dialogue embeddings for the interlocutors (system and user) participating in successful versus unsuccessful dialogues, and by aggregating the embeddings of the interlocutors learn in turn dialogue embeddings for successful versus unsuccessful dialogues.

## Acknowledgments

This work was sponsored by the Army Research Laboratory under Cooperative Agreement Number W911NF-20-2-0053. Many thanks to the anonymous reviewers for their helpful comments.

## References

- Daniel Beck, Trevor Cohn, and Lucia Specia. 2014. Joint emotion analysis via multi-task Gaussian Processes. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1798–1803, Doha, Qatar.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5016–5026, Brussels, Belgium.
- A. Cervone, E. Gambi, G. Tortoreto, E. A. Stepanov, and G. Riccardi. 2018. Automatically predicting user ratings for conversational systems. In *Proc. of CLIC-It*.
- Joshua Coates and Danushka Bollegala. 2018. Frustratingly easy meta-embedding – computing meta-embeddings by averaging source word embeddings. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT), Volume 2 (Short Papers)*, pages 194–198, New Orleans, Louisiana, USA.
- Trevor Cohn and Lucia Specia. 2013. Modelling annotator bias with multi-task Gaussian Processes: An application to machine translation quality estimation. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 32–42, Sofia, Bulgaria.
- 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.
- Klaus-Peter Engelbrecht, Florian Gödde, Felix Hartard, Hamed Ketabdar, and Sebastian Möller. 2009. Modeling user satisfaction with hidden Markov models. In *Proc. of the Annual SIGdial Meeting on Discourse and Dialogue (SIGDIAL)*, pages 170–177, London, UK.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. MultiWOZ 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking base-lines. In *Proc. of the Language Resources and Evaluation Conference (LREC)*, pages 422–428, Marseille, France (Online).
- Matthew Frampton and Oliver Lemon. 2006. Learning more effective dialogue strategies using limited dialogue move features. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 185–192, Sydney, Australia.
- Kallirroi Georgila. 2022. Comparing regression methods for dialogue system evaluation on a richly annotated corpus. In *Proc. of the Workshop on the Semantics and Pragmatics of Dialogue (SemDial:DubDial)*, pages 81–93, Dublin, Ireland.
- Kallirroi Georgila, Carla Gordon, Hyungtak Choi, Jill Boberg, Heesik Jeon, and David Traum. 2019. Toward low-cost automated evaluation metrics for Internet of Things dialogues. In *Proc. of the International Workshop on Spoken Dialogue Systems Technology (IWSDS), Lecture Notes in Electrical Engineering 579*, pages 161–175, Singapore.
- Kallirroi Georgila, Carla Gordon, Volodymyr Yanov, and David Traum. 2020. Predicting ratings of real dialogue participants from artificial data and ratings of human dialogue observers. In *Proc. of the Language Resources and Evaluation Conference (LREC)*, pages 726–734, Marseille, France (Online).
- Kallirroi Georgila, James Henderson, and Oliver Lemon. 2005a. Learning user simulations for Information State Update dialogue systems. In *Proc. of Interspeech*, pages 893–896, Lisbon, Portugal.
- Kallirroi Georgila, James Henderson, and Oliver Lemon. 2006. User simulation for spoken dialogue systems: Learning and evaluation. In *Proc. of Interspeech*, pages 1065–1068, Pittsburgh, Pennsylvania, USA.
- Kallirroi Georgila, Oliver Lemon, and James Henderson. 2005b. Automatic annotation of COMMUNICATOR dialogue data for learning dialogue strategies and user simulations. In *Proc. of the Workshop on the Semantics and Pragmatics of Dialogue (SemDial:DIALOR)*, pages 61–68, Nancy, France.
- Kallirroi Georgila, Oliver Lemon, James Henderson, and Johanna D. Moore. 2009. Automatic annotation of context and speech acts for dialogue corpora. *Journal of Natural Language Engineering*, 15(3):315–353.
- Sarik Ghazarian, Yijia Shao, Rujun Han, Aram Galstyan, and Nanyun Peng. 2023. ACCENT: An automatic event commonsense evaluation metric for open-domain dialogue systems. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4398–4419, Toronto, Canada.
- Sarik Ghazarian, Ralph Weischedel, Aram Galstyan, and Nanyun Peng. 2020. Predictive engagement: An efficient metric for automatic evaluation of open-domain dialogue systems. In *Proc. of the AAAI Conference on Artificial Intelligence*, pages 7789–7796, New York, New York, USA.
- Helen Hastie. 2012. Metrics and evaluation of spoken dialogue systems. In Oliver Lemon and Olivier Pietquin, editors, *Data-Driven Methods for Adaptive Spoken Dialogue Systems*, pages 131–150. Springer.
- James Henderson, Oliver Lemon, and Kallirroi Georgila. 2005. Hybrid reinforcement/supervised learning for dialogue policies from COMMUNICATOR data. In

- Proc. of the IJCAI Workshop on Knowledge and Reasoning in Practical Dialogue Systems*, pages 68–75, Edinburgh, UK.
- James Henderson, Oliver Lemon, and Kallirroi Georgila. 2008. Hybrid reinforcement/supervised learning of dialogue policies from fixed datasets. *Computational Linguistics*, 34(4):487–511.
- Ryuichiro Higashinaka, Yasuhiro Minami, Kohji Dohsaka, and Toyomi Meguro. 2010. Issues in predicting user satisfaction transitions in dialogues: Individual differences, evaluation criteria, and prediction models. In *Proc. of the International Workshop on Dialogue Systems Technology (IWSDS), Lecture Notes in Computer Science 6392*, pages 48–60, Gotemba, Shizuoka, Japan.
- Kate S. Hone and Robert Graham. 2000. Towards a tool for the Subjective Assessment of Speech System Interfaces (SASSI). *Journal of Natural Language Engineering*, 6(3-4):287–303.
- Junjie Huang, Duyu Tang, Wanjun Zhong, Shuai Lu, Linjun Shou, Ming Gong, Daxin Jiang, and Nan Duan. 2021. WhiteningBERT: An easy unsupervised sentence embedding approach. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 238–244, Online and Punta Cana, Dominican Republic.
- Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. On the sentence embeddings from pre-trained language models. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9119–9130, Online.
- Che Liu, Rui Wang, Junfeng Jiang, Yongbin Li, and Fei Huang. 2022. Dial2vec: Self-guided contrastive learning of unsupervised dialogue embeddings. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7272–7282, Abu Dhabi, United Arab Emirates.
- Chia-Wei Liu, Ryan Lowe, Iulian V. Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2122–2132, Austin, Texas, USA.
- Sarah McLeod, Ivana Kruijff-Korbayová, and Bernd Kiefer. 2019. Multi-task learning of system dialogue act selection for supervised pretraining of goal-oriented dialogue policies. In *Proc. of the Annual SIGdial Meeting on Discourse and Dialogue (SIGDIAL)*, pages 411–417, Stockholm, Sweden.
- Shikib Mehri, Jinho Choi, Luis Fernando D’Haro, Jan Deriu, Maxine Eskenazi, Milica Gasic, Kallirroi Georgila, Dilek Hakkani-Tur, Zekang Li, Verena Rieser, Samira Shaikh, David Traum, Yi-Ting Yeh, Zhou Yu, Yizhe Zhang, and Chen Zhang. 2022. Report from the NSF Future Directions Workshop on Automatic Evaluation of Dialog: Research Directions and Challenges. In *arXiv preprint arXiv:2203.10012*.
- Shikib Mehri and Maxine Eskenazi. 2020a. Unsupervised evaluation of interactive dialog with DialoGPT. In *Proc. of the Annual SIGdial Meeting on Discourse and Dialogue (SIGDIAL)*, pages 225–235, Online.
- Shikib Mehri and Maxine Eskenazi. 2020b. USR: An unsupervised and reference free evaluation metric for dialog generation. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 681–707, Online.
- Seyed Mahed Mousavi, Gabriel Roccabruna, Simone Alghisi, Massimo Rizzoli, Mirco Ravanelli, and Giuseppe Riccardi. 2024. Are LLMs Robust for spoken dialogues? In *Proc. of the International Workshop on Dialogue Systems Technology (IWSDS)*, Sapporo, Japan.
- Taghi Paksima, Kallirroi Georgila, and Johanna D. Moore. 2009. Evaluating the effectiveness of information presentation in a full end-to-end dialogue system. In *Proc. of the Annual SIGdial Meeting on Discourse and Dialogue (SIGDIAL)*, pages 1–10, London, UK.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global vectors for word representation. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proc. of the Conference on Empirical Methods in Natural Language Processing and the International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China.
- Jost Schatzmann, Kallirroi Georgila, and Steve Young. 2005. Quantitative evaluation of user simulation techniques for spoken dialogue systems. In *Proc. of the Annual SIGdial Meeting on Discourse and Dialogue (SIGDIAL)*, pages 45–54, Lisbon, Portugal.
- Alexander Schmitt and Stefan Ultes. 2015. Interaction Quality: Assessing the quality of ongoing spoken dialog interaction by experts—And how it relates to user satisfaction. *Speech Communication*, 74:12–36.
- Jianlin Su, Jiarun Cao, Weijie Liu, and Yangyiwen Ou. 2021. Whitening sentence representations for better semantics and faster retrieval. In *arXiv preprint arXiv:2103.15316*.
- Stefan Ultes and Wolfgang Minker. 2014. Interaction quality estimation in spoken dialogue systems using hybrid-HMMs. In *Proc. of the Annual SIGdial Meeting on Discourse and Dialogue (SIGDIAL)*, pages 208–217, Philadelphia, Pennsylvania, USA.

- M. Walker, J. Aberdeen, J. Boland, E. Bratt, J. Garofolo, L. Hirschman, A. Le, S. Lee, S. Narayanan, K. Papineni, B. Pellom, J. Polifroni, A. Potamianos, P. Prabhhu, A. Rudnicky, G. Sanders, S. Seneff, D. Stallard, and S. Whittaker. 2001a. DARPA Communicator dialog travel planning systems: the June 2000 data collection. In *Proc. of Interspeech/Eurospeech*, pages 1371–1374, Aalborg, Denmark.
- Marilyn Walker, Candace Kamm, and Diane Litman. 2000. Towards developing general models of usability with PARADISE. *Journal of Natural Language Engineering*, 6(3-4):363–377.
- Marilyn Walker and Rebecca Passoneau. 2001. DATE: A dialogue act tagging scheme for evaluation of spoken dialogue systems. In *Proc. of the Human Language Technologies Conference*, pages 1–8, San Diego, California, USA.
- Marilyn Walker, Alex Rudnicky, John Aberdeen, Elizabeth Owen Bratt, John Garofolo, Helen Hastie, Audrey Le, Bryan Pellom, Alex Potamianos, Rebecca Passonneau, Rashmi Prasad, Salim Roukos, Greg Sanders, Stephanie Seneff, and Dave Stallard. 2002. DARPA communicator evaluation: Progress from 2000 to 2001. In *Proc. of Interspeech/ICSLP*, pages 273–276, Denver, Colorado, USA.
- Marilyn A. Walker, Rebecca Passoneau, and Julie E. Boland. 2001b. Quantitative and qualitative evaluation of Darpa Communicator spoken dialogue systems. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 515–522, Toulouse, France.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2024. Text embeddings by weakly-supervised contrastive pre-training. In *arXiv preprint arXiv:2212.03533*.
- John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2016. Towards universal paraphrastic sentence embeddings. In *Proc. of the International Conference on Learning Representations (ICLR)*, San Juan, Puerto Rico.

## Appendix

### SYSTEM UTTERANCE 1

*System output*

*System speech act*

*System task*

what trip would you like to take?

request\_info

top\_level\_trip

### USER UTTERANCE 1

*ASR output*

*User transcription*

*User speech act*

*User task*

i'd like to go from atlanta georgia to london england

i'd like to go from atlanta georgia to london england

[provide\_info,provide\_info]

[orig\_city,dest\_city]

### SYSTEM UTTERANCE 2

*System output*

*System speech act*

*System task*

flying from atlanta to london

implicit\_confirm

orig\_dest\_city

### SYSTEM UTTERANCE 3

*System output*

*System speech act*

*System task*

what day are you leaving atlanta?

request\_info

depart\_arrive\_date

### USER UTTERANCE 2

*ASR output*

*User transcription*

*User speech act*

*User task*

september twenty fourth

september twenty fourth

[provide\_info]

[depart\_date]

### SYSTEM UTTERANCE 4

*System output*

*System speech act*

*System task*

om monday, september twenty fourth

implicit\_confirm

depart\_arrive\_date

### SYSTEM UTTERANCE 5

*System output*

*System speech act*

*System task*

about what time do you want to leave?

request\_info

depart\_arrive\_time

### USER UTTERANCE 3

*ASR output*

*User transcription*

*User speech act*

*User task*

early evening

early evening

[provide\_info]

[depart\_time]

### SYSTEM UTTERANCE 6

*System output*

*System speech act*

*System task*

leaving in the early evening

implicit\_confirm

depart\_arrive\_time

### SYSTEM UTTERANCE 7

*System output*

*System speech act*

*System task*

please stand by while i look up the flight schedules on the web

status\_report

retrieval

### SYSTEM UTTERANCE 8

*System output*

*System speech act*

*System task*

ok, got them

status\_report

retrieval

### SYSTEM UTTERANCE 9

*System output*

*System speech act*

*System task*

i can only get you a flight in the early evening on united

present\_info

flight

Figure 1: An example dialogue excerpt showing both the ASR output and the transcription of the user's speech, and speech act and task annotations for both the system and the user. User speech act and task annotations, and dialogue context information are computed from the ASR output.

DIALOGUE LEVEL

**Turn:** user  
**TurnStartTime:** 990207520.875  
**TurnEndTime:** 990207525.562  
**TurnNumber:** 3  
**Speaker:** user  
**UtteranceStartTime:** 990207520.875  
**UtteranceEndTime:** 990207525.562  
**UtteranceNumber:** 3  
**DialogueActType:** user  
**ConvDomain:** about\_task  
**SpeechAct:** [provide\_info]  
**AsrInput:** <date\_time>early evening</date\_time>  
**TransInput:** <date\_time>early evening</date\_time>  
**Output:**

TASK LEVEL

**Task:** [depart\_time]  
**FilledSlot:** [depart\_time]  
**FilledSlotValue:** [early evening]  
**GroundedSlot:** [depart\_date]

LOW LEVEL

**WordErrorRatenois:** 0.00  
**WordErrorRate:** 0.00  
**SentenceErrorRate:** 0.00  
**KeywordErrorRate:** 0.0  
**ComputeErrorRatesReturn Value:** 0

HISTORY LEVEL

**FilledSlotsStatus:** [dest\_city],[orig\_city],[depart\_date],[depart\_time]  
**FilledSlotsValuesStatus:** [london england],[atlanta georgia],[september twenty fourth],[early evening]  
**GroundedSlotsStatus:** [],[dest\_city],[orig\_city],[depart\_date]  
**SpeechActsHist:** request\_info,[provide\_info,provide\_info],implicit\_confirm,request\_info,[provide\_info],implicit\_confirm,request\_info,[provide\_info]  
**TasksHist:** top\_level\_trip,[orig\_city,dest\_city],orig\_dest\_city,depart\_arrive\_date,[depart\_date],depart\_arrive\_date,depart\_arrive\_time,[depart\_time]  
**FilledSlotsHist:** [orig\_city,dest\_city],[depart\_date],[depart\_time]  
**FilledSlotsValuesHist:** [atlanta georgia,london england],[september twenty fourth],[early evening]  
**GroundedSlotsHist:** [],[orig\_city,dest\_city],[depart\_date]

Figure 2: An example dialogue state generated after user utterance 3 in Figure 1. Empty (‘[]’) values or ‘null’ values (not seen here) do not affect the accuracy of the slot values.

	TaskEase		SysComEase		UserExp		SysBehExp		SysFutUse	
	man	auto								
<b>Average of Word Embeddings</b>										
bert-L1	1.259	1.28	<b>1.163</b>	1.174	<b>1.287</b>	1.303	1.291	1.324	<b>1.363</b>	1.379
bert-L1-10-11	<b>1.254</b>	1.256	1.173	<b>1.167</b>	1.302	<b>1.289</b>	<b>1.283</b>	1.285	1.366	1.371
bert-L1-11-12	1.255	1.256	1.174	1.173	1.302	1.298	1.288	1.287	1.364	1.373
bert-L1-10-11-12	1.255	<b>1.253</b>	1.173	1.169	1.305	1.291	1.285	<b>1.28</b>	1.367	<b>1.37</b>
<b>Output of [CLS] Token</b>										
bert-L1	<b>1.276</b>	1.289	1.184	1.201	<b>1.305</b>	1.303	<b>1.296</b>	1.323	<b>1.375</b>	1.394
bert-L1-10-11	1.291	1.294	1.184	1.178	1.312	1.306	1.312	1.326	1.389	1.399
bert-L1-11-12	1.282	<b>1.287</b>	<b>1.178</b>	<b>1.176</b>	1.313	<b>1.302</b>	1.307	<b>1.304</b>	1.376	<b>1.384</b>
bert-L1-10-11-12	1.285	<b>1.287</b>	1.182	1.177	<b>1.305</b>	1.303	1.311	1.315	1.384	1.396

Table 5: RMSE values when calculating sentence embeddings as an average of BERT word embeddings versus using the output of the [CLS] token, based on the human transcriptions (‘man’) and the ASR outputs (‘auto’). Domain-independent features are not included. The best value for each column for the output of the [CLS] token is shown in red and in bold. The average of BERT word embeddings always outperforms the output of the [CLS] token. **The best value for each column across both types of models is shown in black and in bold.**

# Question Type Prediction in Natural Debate

Zlata Kikteva Alexander Trautsch Steffen Herbold Annette Hautli-Janisz

Faculty of Computer Science and Mathematics

University of Passau

firstname.lastname@uni-passau.de

## Abstract

In spontaneous natural debate, questions play a variety of crucial roles: they allow speakers to introduce new topics, seek other speakers' opinions or indeed confront them. A three-class question typology has previously been demonstrated to effectively capture details pertaining to the nature of questions and the different functions associated with them in a debate setting. We adopt this classification and investigate the performance of several machine learning approaches on this task by incorporating various sets of lexical, dialogical and argumentative features. We find that BERT demonstrates the best performance on the task, followed by a Random Forest model enriched with pragmatic features.

## 1 Introduction

Questions are at the core of human communication and can be used in a variety of ways, from simply eliciting information from the hearer to communicating a speaker's standpoint. They are also crucial in argumentation, where they make up around 5% of all speech acts and are rarely left ignored (Kikteva et al., 2022). However, question type prediction faces a number of challenges. First, the lexical surface does not correlate with the function of the question (e.g., 'Who would do that?' can either be a request for information on a set of entities or a rhetorical question communicating that no one would do that). Secondly, context is assumed to be crucial for their interpretation, however, it is not exactly clear what features in the context are indeed relevant. And lastly, a large majority of computational work assumes a bipartite distinction into information-seeking and rhetorical questions, a classification that does not capture the variety of functions that questions fulfil in debate. There is, in fact, a third category of questions referred to as assertive questions that has been theoretically motivated (Freed, 1994) and empirically tested (Visser

et al., 2020; Hautli-Janisz et al., 2022b). Such questions are characterised by the speaker's intention to express their opinion while still seeking information (as opposed to information-seeking questions the main purpose of which is to elicit a response, and the rhetorical ones which do not necessitate one).

With this paper, we examine the impact of different (1) context configurations and (2) combinations of carefully selected lexical, pragmatic and argumentative features for the question type prediction. We explore deep learning approaches that tend to effectively capture lexical features as well as statistical models that while still capable of representing lexical information, also benefit from categorical feature inputs. Our results indicate the introduction of the third question type drastically increases the complexity of the task when compared to the binary classification. We find that BERT achieves the highest scores when the input is enriched with lexical features of either the preceding material or response. Furthermore, we report our second-best results with a Random Forest model that makes use of a rich pragmatic feature set.

## 2 Previous work

There is a significant body of work on questions in computational linguistics in the context of question-answering systems, i.e., question classification based on the type of information expected as an answer, which focuses predominantly on factoid questions. Earlier approaches include extensive use of lexical and syntactic features combined with traditional statistical approaches (Zhang and Lee, 2003; Metzler and Croft, 2005; Huang et al., 2008; Silva et al., 2011; Loni, 2011; Tayyar Madabushi and Lee, 2016), while recently more work utilizes extensive capabilities of the deep learning models (Kim, 2014; Iyyer et al., 2014; Sun et al., 2019; Anhar et al., 2019; Yilmaz and Toklu, 2020)

Questions in natural communication, however, are often used to elicit more than just factual information from interlocutors and instead serve a variety of communicative purposes. Thus, to interpret the function of questions in discourse, researchers often adhere to a bipartite distinction: Harper et al. (2009) focus on identification of conversational versus factual questions, Bhattasali et al. (2015); Ranganath et al. (2016); Oraby et al. (2017) distinguish rhetorical questions from non-rhetorical ones; Kalouli et al. (2018, 2021) identifies information and non-information-seeking questions; Bagga et al. (2021) categorize questions into unpalatable and not unpalatable ones from the perspective of abusive language detection. They adopt a range of approaches such as the use of lexical features like n-grams, POS tags, speaker roles and word embeddings as well as the modelling of the context surrounding questions.

However, recent research in pragmatics suggests that question typology is a bit more complex than previously assumed. In their work, Hautli-Janisz et al. (2022a) and Kikteva et al. (2022) discuss a more fine-grained question typology that attempts to better capture the conversational functions that are fulfilled by questions in debate. In this work, we follow their distinction into information-seeking, rhetorical and assertive questions.

**Information-seeking questions** Also called pure questions (PQs), those are used to elicit information from an interlocutor. For instance, in Example 1, moderator Fiona Bruce seeks the views of the panel members on the matter of the voter ID in the UK.<sup>1</sup>

- (1) Fiona Bruce: *Will voter IDs protect the integrity of elections or just undermine the UK democracy?*

**Rhetorical questions** With RQs, speakers express their own opinion or standpoint, illustrated in Example 2 by Liz Saville-Roberts’s intention to communicate her dissatisfaction with the current state of the prison system in the UK.<sup>2</sup> The speaker poses the question without expecting to hear a response which is signalled by the fact that she continues talking.

- (2) Liz Saville-Roberts: *The black population in Wales is over-represented by five times within the prison population of Wales,*

*surely that is a desperate failure? That is an indication of the racism in our society in action.*

**Assertive questions** AQs serve the double purpose of communicating information and asking for confirmation/rejection from an interlocutor. In Example 3, Gillian Keegan expresses her frustration regarding the police having to inspect every package due to the new Brexit regulations, while at the same time expecting other panel members to agree with her opinion on the matter.<sup>3</sup>

- (3) Gillian Keegan: *The police probably have the legal right to open every packet and inspect every sausage. Isn’t that unreasonable?*

### 3 Data

The data underlying our investigation is taken from QT30 (Hautli-Janisz et al., 2022b), the largest ever dataset of broadcast political debate. The corpus comprises the transcriptions of 30 episodes of the UK’s talk show ‘Question Time’ (QT) between June 2020 and November 2021 and is manually annotated with Inference Anchoring Theory (IAT) (Budzynska et al., 2014, 2016), a framework that captures how argumentation unfolds and is reacted to in dialogue which allows us to extract questions of the three types as well as pragmatic features associated with them for the analysis.

The questions dataset used in the current work contains 2 867 questions, with the split into training and test given in Table 1. PQs make up almost 70% of all questions, both RQs and AQs are significantly less frequent and make up about 14% and 16% of the data, respectively. Questions extracted from QT30 are used for training; an additional 10 episodes of QT that were broadcast and analyzed in 2022 are used for testing. With this time split, we train on about 77% of the data and evaluate on the rest.

Table 1: Training and test split.

	PQ	RQ	AQ	Total
Training	1 555	306	343	2 204
Test	446	87	130	663
Total	2 001	393	473	2 867

<sup>1</sup><http://corpora.aifdb.org/qt30>, json ID: 21308

<sup>2</sup><http://corpora.aifdb.org/qt30>, json ID: 18464

<sup>3</sup><http://corpora.aifdb.org/qt30>, json ID: 23888

## 4 Question type prediction

### 4.1 Feature selection

**Lexical features** Lexical features include questions and corresponding preceding and response texts. We consider one locution, i.e., discourse unit, dialogically preceding the question to be the preceding context, while the response constitutes any number of locutions that are contributed by the same or different speaker than the question speaker, directly following the question until the end of that speaker’s turn. We represent them as n-grams for statistical models and as embeddings for deep learning models. For n-grams, we extract all available unigrams, as well as bi- and tri-grams that appear in at least two documents. We further process n-grams to identify the most relevant features by applying TF-IDF vectorization to the data. The vectorization is performed per question type allowing us to model feature representations for each question type separately.

**Pragmatic features in the response** We extract the following argumentative and pragmatic features from the response to the question:

- **Speaker roles** Information on whether the question and response material comes from the moderator, a panel member or an audience member and on whether the question and response speakers are the same or not.
- **Answers** Statements instantiating the content of the question.
- **Propositional relations** Inference (support between two statements), conflict (attack between two statements) and rephrase (reformulation or refinement of a previous statement).
- **Epistemic markers** Indicators of speaker commitment.

The number of answers, propositional relations, and epistemic markers is normalised on the level of the locution, i.e., we encode the relative frequency of a feature per locution in the response. Speaker roles are represented as categorical values.

### 4.2 Modeling

**Statistical models** In order to model both the pragmatic features of the response and lexical features of the question with its adjacent context, we

Table 2: Balanced accuracies for different context configurations. PREC, QU, and RESP stand for preceding, question, and response texts respectively.

Context	LSTM	BERT	RF	SVM
QU	0.37	0.47	0.40	0.40
PREC-QU	0.31	<b>0.48</b>	0.35	0.36
QU-RESP	0.42	<b>0.48</b>	0.35	0.35
PREC-RESP	0.34	0.40	0.34	0.36
PREC-QU-RESP	0.36	0.42	0.35	0.37

use Random Forest (RF) and Support Vector Machine (SVM) for classification. After hyperparameter tuning, we choose an entropy criterion and 200 estimators with a maximum depth of 8 and minimum sample split of 0.1 for RF; for SVM we use an RBF kernel with a one-vs-one multiclass classification strategy.

**Language models** To gain insight into how deep learning models compare to more traditional machine learning approaches, we use an LSTM model and an LLM. For LSTM we use softmax activation with categorical cross-entropy as a loss function and the Adam optimizer with a batch size of 64, a maximum sequence length of 400 and 100-dimensional embeddings trained over 6 epochs. For an LLM we use a cased, large variant of BERT (Devlin et al., 2018) with 336M parameters which we retrieved from the Huggingface Model Hub.<sup>4</sup> We adopt the same configuration as used by Huggingface to evaluate BERT on the GLUE benchmark (Wang et al., 2019). We train for three epochs, with a learning rate of 2e-05 and a batch size of 32 with a maximum sequence length of 400, which is sufficiently large for all inputs.

**Testing** In order to mitigate how unbalanced the dataset is, we resort to an oversampling technique for the training set by matching the number of underrepresented RQs and AQs to the number of PQs. For the same reason, for the evaluation of the models’ performance, we use balanced accuracy. The code is publicly available at <https://github.com/ZlataKikteva/sigdial2024-questions>.

## 5 Results

### 5.1 Context

We first model the impact of the lexical features in context on the multiclass question type prediction task as it has been observed to improve the

<sup>4</sup>‘bert-large-cased’

Table 3: Balanced accuracies for pragmatic features (in response) and lexical features (question text). The highest scores for each feature set are underlined; the highest score overall is in bold.

Model	Feature Set	Pragmatic Features				
		Speakers	Answers	Prop.rel.	Ep. markers	All
RF	Pragmatic Features only	0.41	<u>0.43</u>	0.38	0.37	0.40
	Pragmatic Features & Unigrams	0.41	<b>0.44</b>	0.43	0.42	0.42
	Pragmatic Features & Uni- and Bigrams	0.42	<u>0.43</u>	0.42	0.41	<u>0.43</u>
	Pragmatic Features & Uni-, Bi and Trigrams	0.41	0.42	0.41	0.41	<u>0.43</u>
SVM	Pragmatic Features only	0.42	<u>0.43</u>	0.38	0.37	0.37
	Pragmatic Features & Unigrams	0.40	0.35	0.36	0.34	<u>0.41</u>
	Pragmatic Features & Uni- and Bigrams	0.40	0.35	0.35	0.34	<u>0.42</u>
	Pragmatic Features & Uni-, Bi and Trigrams	<u>0.41</u>	0.36	0.36	0.35	<u>0.41</u>

predictions in a binary classification setting (Bhat-tasali et al., 2015; Kalouli et al., 2021). We model all possible context combinations of the question, preceding, and response material. The results are reported in Table 2. For this task, we employ both language models as well as RF and SVM models, for the latter unigrams are selected as input features the use of which results in better performance than bi-, trigrams, or any n-gram combination.

In this setting, BERT achieves the highest score overall with 0.48 for PREC-QU and QU-RESP context combinations confirming the positive effect of context on question type classification.<sup>5</sup> Notably, we see an even stronger impact of QU-RESP combination on LSTM performance with its score increasing from 0.37 to 0.42. However, with the statistical approach, the inclusion of context does not benefit either of the models.

From these results, we infer the following: (1) the use of context for multiclass question type prediction seems to be beneficial only in some settings; (2) the gap between the performance levels of an LLM and statistical models is not as large as it could be expected considering the disparity in the amounts of computational power required for using the latter. With this in mind, we explore in the next section the possibility of further improving RF and SVM results by incorporating pragmatic features.

## 5.2 Pragmatic and lexical features

The results of the statistical models using the pragmatic and lexical features are presented in Table 3. For this set of experiments, we use only question text for extracting relevant n-grams based on the results reported in Table 2. With this set of experiments, we find that the inclusion of pragmatic

<sup>5</sup>We tested other models which all yielded inferior results, including RoBERTa, DeBERTaV3, and a zero-shot setting with Vicuna.

features in addition to n-grams improves the performances of both, RF and SVM, with the former achieving the highest balanced accuracy score for statistical models of 0.44. Overall, the RF model has comparable results for all feature sets with the presence of an answer seemingly having a slight edge in terms of its impact on the performance. However, the SVM tends to rely more heavily on the speaker roles as well as the combination of all of the pragmatic features. Finally, a relatively high score of 0.43 can be observed in the setting with the pragmatic features only. However, after further inspection, we find that models in this setting predict only two of the classes, indicating that lexical features carry valuable information that cannot be overlooked.

We also note that when adopting a multiclass approach to question typology, the performance of the models drops considerably when compared to previous research focused on two classes of questions at a time. Kalouli et al. (2021) achieve accuracy of about 0.88 when identifying information and non-information-seeking questions; Ranganath et al. (2016) and Oraby et al. (2017) distinguish between rhetorical and non-rhetorical questions with F1-scores of about 0.64 and 0.76 respectively. Our results suggest that an introduction of the third category of questions increases the task complexity.

## 5.3 Error analysis

In order to further investigate the results we conduct an error analysis by examining confusion matrices for the best-performing models in both settings (see Appendices A and B). Unsurprisingly, we find that the unbalanced nature of the dataset with a higher number of PQs compared to AQs and RQs results in models demonstrating better performance in the case of the over-represented category. With respect to the BERT models that

take into account context, there is also a relatively high number of AOs and ROs misclassified as POs (about 40% for AOs and over 50% for ROs) and a considerably lower rate of misclassification between AOs and ROs (around 20%). As for the RF model that considers lexical features and answers, we observe that it is much better at identifying ROs than the BERT model (almost 40% of correct predictions compared to around 20%). In particular, it is more successful at distinguishing them from the POs which can be attributed to the fact that ROs are less likely to be answered because of their nature. Overall, the results of the error analysis indicate that in the case of both approaches, the models exhibit better performance when it comes to the POs while displaying varying degrees of difficulty with the other two question types.

## 6 Conclusion

In this paper, we introduce a task of question type prediction using a more fine-grained typology than the typically adopted bipartite distinction and find that the introduction of the third class increases the complexity of the task drastically. While questions can be tricky to interpret as speaker intention, which is notoriously hard to capture, is often one of the main indicators of the question type, the task complexity is further increased by a high class imbalance in naturally occurring data.

We tackle the task by adopting several approaches conventionally used for the binary question type prediction such as the use of lexical features and the incorporation of question-adjacent context as well as by using novel for this task pragmatic features including propositional relations and epistemic markers in responses to questions. We find that BERT exhibits the best performance with an RF model trained on a combination of pragmatic features and unigrams taking second place. However, neither of these results is truly satisfactory, suggesting that current machine learning approaches are not yet powerful enough to reason about the nature of a question if we adopt a finer granularity into three types.

## Acknowledgements

The work reported on in this paper was partially funded by the VolkswagenStiftung under grant Az. 98544 ‘Deliberation Laboratory’.

## References

- Refany Anhar, Teguh Bharata Adji, and Noor Akhmad Setiawan. 2019. Question classification on question-answer system using bidirectional-lstm. In *2019 5th International Conference on Science and Technology (ICST)*, volume 1, pages 1–5. IEEE.
- Sunyam Bagga, Andrew Piper, and Derek Ruths. 2021. “are you kidding me?”: Detecting unpalatable questions on reddit. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2083–2099.
- Shohini Bhattachali, Jeremy Cytryn, Elana Feldman, and Joonsuk Park. 2015. Automatic identification of rhetorical questions. 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 743–749.
- Katarzyna Budzynska, Mathilde Janier, Juyeon Kang, Chris Reed, Patrick Saint-Dizier, Manfred Stede, and Olena Yaskorska. 2014. Towards argument mining from dialogue. In *Proceedings of the Fifth International Conference on Computational Models of Argument (COMMA 2014)*, pages 185–196. IOS Press.
- Katarzyna Budzynska, Mathilde Janier, Chris Reed, and Patrick Saint Dizier. 2016. Theoretical foundations for illocutionary structure parsing. *Argument & Computation*, 7(1):91–108.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. *BERT: pre-training of deep bidirectional transformers for language understanding*. *CoRR*, abs/1810.04805.
- Alice F Freed. 1994. The form and function of questions in informal dyadic conversation. *Journal of pragmatics*, 21(6):621–644.
- F Maxwell Harper, Daniel Moy, and Joseph A Konstan. 2009. Facts or friends? distinguishing informational and conversational questions in social q&a sites. In *Proceedings of the sigchi conference on human factors in computing systems*, pages 759–768.
- Annette Hautli-Janisz, Katarzyna Budzynska, Conor McKillop, Brian Plüss, Valentin Gold, and Chris Reed. 2022a. Questions in argumentative dialogue. *Journal of Pragmatics*, 188:56–79.
- Annette Hautli-Janisz, Zlata Kikteva, Wassiliki Siskou, Kamila Gorska, Ray Becker, and Chris Reed. 2022b. *QT30: A corpus of argument and conflict in broadcast debate*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3291–3300, Marseille, France. European Language Resources Association.
- Zhiheng Huang, Marcus Thint, and Zengchang Qin. 2008. Question classification using head words and

- their hypernyms. In *Proceedings of the 2008 Conference on empirical methods in natural language processing*, pages 927–936.
- Mohit Iyyer, Jordan Boyd-Graber, Leonardo Claudino, Richard Socher, and Hal Daumé III. 2014. [A neural network for factoid question answering over paragraphs](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 633–644, Doha, Qatar. Association for Computational Linguistics.
- Aikaterini-Lida Kalouli, Katharina Kaiser, Annette Hautli-Janisz, Georg A Kaiser, and Miriam Butt. 2018. A multilingual approach to question classification. In *Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, pages 2715–2720.
- Aikaterini-Lida Kalouli, Rebecca Kehlbeck, Rita Sevastjanova, Oliver Deussen, Daniel A Keim, and Miriam Butt. 2021. Is that really a question?: Going beyond factoid questions in nlp. In *14th International Conference on Computational Semantics: IWCS 2021*, pages 132–143.
- Zlata Kikteva, Kamila Gorska, Wassiliki Siskou, Annette Hautli-Janisz, and Chris Reed. 2022. [The key-stone role played by questions in debate](#). In *Proceedings of the 3rd Workshop on Computational Approaches to Discourse*, pages 54–63, Gyeongju, Republic of Korea and Online. International Conference on Computational Linguistics.
- Yoon Kim. 2014. [Convolutional neural networks for sentence classification](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751, Doha, Qatar. Association for Computational Linguistics.
- Babak Loni. 2011. A survey of state-of-the-art methods on question classification. *Delft University of Technology, Tech. Rep.*, 55:57.
- Donald Metzler and W Bruce Croft. 2005. Analysis of statistical question classification for fact-based questions. *Information Retrieval*, 8:481–504.
- Shereen Oraby, Vrindavan Harrison, Amita Misra, Ellen Riloff, and Marilyn Walker. 2017. [Are you serious?: Rhetorical questions and sarcasm in social media dialog](#). In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 310–319, Saarbrücken, Germany. Association for Computational Linguistics.
- Suhas Ranganath, Xia Hu, Jiliang Tang, Suhang Wang, and Huan Liu. 2016. Identifying rhetorical questions in social media. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 10, pages 667–670.
- Joao Silva, Luísa Coheur, Ana Cristina Mendes, and Andreas Wichert. 2011. From symbolic to sub-symbolic information in question classification. *Artificial Intelligence Review*, 35:137–154.
- Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2019. How to fine-tune bert for text classification? In *Chinese computational linguistics: 18th China national conference, CCL 2019, Kunming, China, October 18–20, 2019, proceedings 18*, pages 194–206. Springer.
- Harish Tayyar Madabushi and Mark Lee. 2016. [High accuracy rule-based question classification using question syntax and semantics](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1220–1230, Osaka, Japan. The COLING 2016 Organizing Committee.
- Jacky Visser, Barbara Konat, Rory Duthie, Marcin Koszowy, Katarzyna Budzynska, and Chris Reed. 2020. Argumentation in the 2016 us presidential elections: annotated corpora of television debates and social media reaction. *Language Resources and Evaluation*, 54(1):123–154.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. [Glue: A multi-task benchmark and analysis platform for natural language understanding](#). *Preprint*, arXiv:1804.07461.
- Seyhmus Yilmaz and Sinan Toklu. 2020. A deep learning analysis on question classification task using word2vec representations. *Neural Computing and Applications*, 32(7):2909–2928.
- Dell Zhang and Wee Sun Lee. 2003. Question classification using support vector machines. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pages 26–32.

**A Confusion matrices for the best-performing models in the context setting**

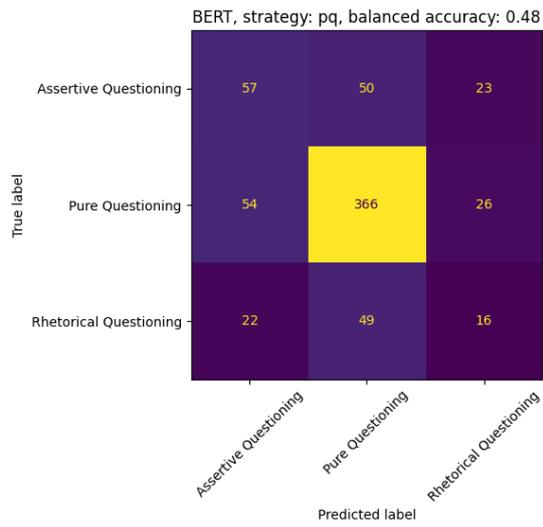


Figure 1: BERT in PREC-QU configuration

**B Confusion matrix for the best-performing model in the pragmatic and lexical feature setting**

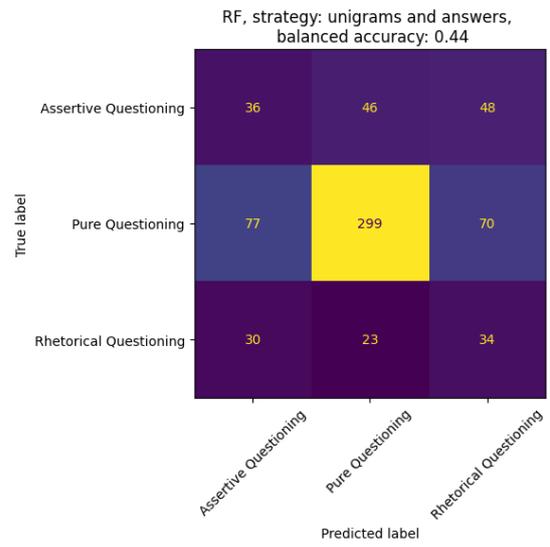


Figure 3: RF in Answers & Unigrams configuration

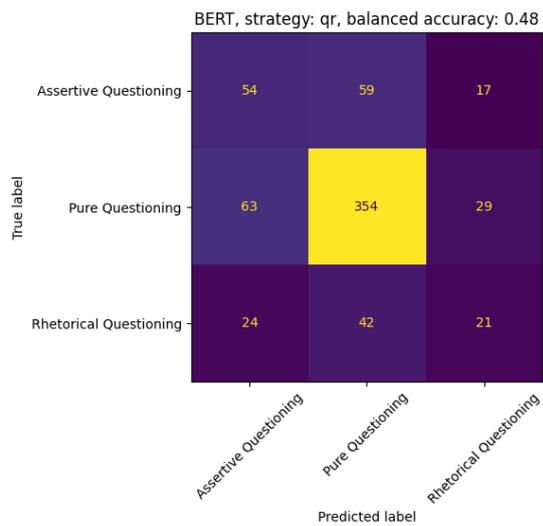


Figure 2: BERT in QU-RESP configuration

# MemeIntent: Benchmarking Intent Description Generation for Memes

Jeongsik Park\* Khoi P. N. Nguyen\* Terrence Li Suyesh Shrestha  
Megan Kim Vu Jerry Yining Wang Vincent Ng  
Human Language Technology Research Institute  
University of Texas at Dallas  
{jeongsik.park,khoi.nguyen6,vince}@utdallas.edu

## Abstract

While recent years have seen a surge of interest in the automatic processing of memes, much of the work in this area has focused on determining whether a meme contains malicious content. This paper proposes the new task of *intent description generation*: generating a description of the author's intentions when creating the meme. To stimulate future work on this task, we (1) annotated a corpus of memes with the intents being perceived by the reader as well as the background knowledge needed to infer the intents and (2) established baseline performance on the intent description generation task using state-of-the-art large language models. Our results suggest the importance of background knowledge retrieval in intent description generation for memes.

## 1 Introduction

*Memes*, which are "amateur media artifacts, extensively remixed and recirculated by different participants on social media networks" (Milner, 2012), are typically created with an intent to perform some "action" (Grundlingh, 2018). While many memes are intended to make a joke (where the author tries to make fun of a celebrity's weird accent, for instance), other memes may have malicious intentions. For instance, a meme author may seek to provoke fear (e.g., by conveying the message that vaccines contain microchips) or manipulate public opinions (e.g., by portraying Hillary Clinton as a corrupt politician before the 2016 presidential campaign with the goal of garnering support for Donald Trump). The core task in automated meme understanding, therefore, involves identifying the intent behind the creation of a meme.

In this paper, an *intent* is defined as an action that the meme author does via the meme. For example, the meme in Figure 1a intends to "mock



Figure 1: Example memes from Dimitrov et al. (2021a) (left) and Sharma et al. (2023) (right).

Justin Trudeau as a communist for being similar to Fidel Castro", while Figure 1b "makes fun of Donald Trump's hypocrisy regarding his view on the severity of COVID-19". As shown in these examples, intents are best represented in textual form. Therefore, intent identification is naturally cast as a generation task, hence will be called *intent description generation*. Automatically generating the intent description of a meme is by no means an easy task, for at least two reasons:

First, background knowledge is often needed for proper interpretation of a meme. *Background knowledge* refers to the knowledge that is not present in the meme but is needed to recognize the intent when combined with the information that is explicitly stated in the meme. There are different kinds of background knowledge, including historical knowledge (e.g., "Make America Great Again" is the slogan used by Trump in his presidential campaigns'), general political ideologies (e.g., 'progressives favor stricter gun control policies'), or knowledge of the meme culture (e.g., 'the meme template Drakeposting<sup>1</sup> funnily expresses an objection and an approval'), etc. For example, Figure 1a requires the knowledge that Castro is a communist

\*These authors contributed equally to this work.

<sup>1</sup><https://knowyourmeme.com/memes/drakeposting>

leader and that there has been a fear of communism in the West since the Cold War. Figure 1b, on the other hand, does not require any special knowledge as all of the information necessary to understand the meme is presented at face value.

Second, in order to derive the intent, complex *inference mechanisms* may be needed to combine background knowledge with different pieces of information extracted from the image and text portions of the meme. Those "combination" steps reflect how a human thinks, such as based on logos (i.e., logical reasoning), ethos (i.e., speaker's authority), or pathos (i.e., emotional appeal). Some of these steps are not about logical reasoning, thus harder to automate (Mondorf and Plank, 2024).

As an example of such inference mechanisms, consider Figure 1a again. To arrive at the final intent, we first have to recognize that the person on the left is Justin Trudeau (Canada's prime minister) and the person on the right is Fidel Castro (Cuba's former leader). In addition, Trudeau has his mouth open whereas Castro has his mouth closed, which signifies that Trudeau is speaking and Castro is listening. When combining this information with the text "Happy Father's Day", one can infer that Trudeau either admits that Castro is his father or simply likes Castro enough to send his greetings to him. Then, combining the background knowledge that Castro was a staunch communist with the fact derived earlier that Trudeau admires Castro, the meme poster is trying to transfer the communist nature of Castro to Trudeau to damage Trudeau's reputation. Given the negative sentiment towards communism in the Western public, the final intent is thus "mocking Trudeau as a communist for being similar to Castro".

Intent description generation, though challenging, is a task whose solution has both practical and theoretical significance. From a practical perspective, knowledge of the intent being perceived through the meme could be useful for other meme-related processing tasks. For instance, knowing what the intent is could facilitate the determination of whether a meme contains harmful content (Pranick et al., 2021a) or the detection of persuasion techniques (Dimitrov et al., 2021a). Theoretically speaking, being able to generate intents like humans requires that a machine read between the lines and achieve a deeper level of understanding of perceptual input, enabling machine perception to get one step closer to human perception.

Our contributions in this paper are four-fold.

First, we propose the new task of intent description generation. Second, we construct the first benchmark for intent description generation, called *MemeIntent*, which shows the background knowledge required for each meme and its final intent(s). Third, we produce preliminary results on MemeIntent from two state-of-the-art language and vision-language models. Finally, based on the experimental results, we justify the need for more careful treatments of background knowledge in meme processing. To stimulate future research in intent description generation for memes, MemeIntent has been made publicly available<sup>2</sup>.

The rest of this paper is organized as follows. Section 2 provides an overview of related work on automated meme processing. In Section 3, we describe our intent description generation benchmark, MemeIntent. To get an idea of how challenging intent description generation is, we conduct experiments on MemeIntent, discussing our experimental setup in Section 4 and showing preliminary evaluation results of two state-of-the-art large language models on MemeIntent in Section 5. Finally, we present our conclusions in Section 6.

## 2 Related Work

There has been a recent surge of interest in meme processing. Table 1 summarizes our survey on the tasks that have been proposed up to date for meme processing. These tasks can be classified into three groups: categorization, interpretation, and explanation, which will be described next.

### 2.1 Categorization Tasks

A growing effort has been made to assemble internet memes and categorically label them along various dimensions. These tasks can be broadly categorized into two groups.

The first group is composed of tasks that ask to classify malicious memes, including the offensive (Suryawanshi et al., 2020a), trolling (Suryawanshi et al., 2020b), hateful (Kiela et al., 2020), anti-semitic (Chandra et al., 2021), harmful (Pranick et al., 2021b,c), and misogynous (Fersini et al., 2022). The second group is composed of tasks about detecting other aspects of memes such as persuasion techniques (Dimitrov et al., 2021b), figurative language (e.g., allusion, irony, sarcasm, contrast, etc.) (Liu et al., 2022), people's roles (e.g.,

<sup>2</sup><https://github.com/JeongSikPark1998/MemeIntent>

Task	Dataset name	Topics	Size
Offensiveness Identification	MultiOFF (Suryawanshi et al., 2020a)	US Election	743
Troll Classification	TamilMemes (Suryawanshi et al., 2020b)	Tamil Memes	2,969
Hatefulness Detection	HatefulMemes (Kiela et al., 2020)	N/A	10K
Antisemitism	Jewtocracy (Chandra et al., 2021)	Antisemitism	3,102+3,509
Harm Detection	HarMeme (Pramanick et al., 2021b)	Covid	3544
Harm Detection	HARM-C&P (Pramanick et al., 2021c)	Covid, Politics	3,544; 3,552
Persuasion Technique Detection	SemEval-2021-T6 (Dimitrov et al., 2021b)	Mixed	950
Emotion Classification	Memotion (Sharma et al., 2020)	N/A	10K
Fine-grained Hatefulness Detection	WOAH-5 (Mathias et al., 2021)	N/A	10K
Misogyny Identification	MAMI (Fersini et al., 2022)	N/A	15K
Figurative Language Detection	FigMemes (Liu et al., 2022)	Politics	5,141
Role Labelling of Entities	HVVMemes (Sharma et al., 2022)	Covid, Politics	7K
Explaining Hate	HatReD (Hee et al., 2023)	N/A	3,228
Explaining Role of Entities	ExHVV (Sharma et al., 2023)	Covid, Politics	4,680
Meme Captioning	MemeCap (Hwang and Shwartz, 2023)	No offensive/sexual	6,387
<b>Intent Description Generation</b>	<b>MemeIntent</b>	<b>Mixed</b>	<b>950</b>

Table 1: **Tasks related to memes processing and associated benchmarks.** *Mixed* means *politics, vaccines, COVID-19, gender equality*. The three groups (separated by horizontal lines) are about categorization, explanation, and interpretation tasks, respectively.

hero, villain, or victim) (Sharma et al., 2022), emotion (e.g., humor, sarcasm, motivation, or offensiveness) (Sharma et al., 2020), and attacked targets (e.g., religion, race, sex, nationality, or disability) (Mathias et al., 2021).

## 2.2 Interpretation Tasks

The second category of work on meme processing involves the relatively new task of meme interpretation, which involves generating text that captures the final meaning of a meme. Because intent description generation is a meme interpretation task, this category is the central interest of this paper.

To the best of our knowledge, meme interpretation has only been studied by Hwang and Shwartz (2023), who proposed the task of *meme captioning*, which means "describing the meaning of the meme". MemeCap, the dataset they produced as part of their work, contains memes images from Reddit. For each meme, they manually annotated the meme captions, the *literal captions* (i.e., the caption of the image excluding the text), and the *visual metaphors* (i.e., associations between entities on the meme and its actual target).

Intent description generation, while being a meme interpretation task, can be seen as the next level of meme captioning. Grundlingh (2018), a linguist studying memes, has argued that a meme, like an utterance, has both *illocutionary* and *perlocutionary* acts. In other words, a meme *says something to do something*. As such, while meme captioning is about what the meme is *saying* (the illocutionary act), intent description generation is



Figure 2: (a) A meme from MemeCap (Hwang and Shwartz, 2023), with title "Simpsons predicted it yet again". (b) A meme from Dimitrov et al. (2021a).

concerned with what the meme is *doing* (the perlocutionary act).

For example, for the meme in Figure 2a, the caption from MemeCap is "The Simpsons was correct about its use of Trump and Greta Thurnberg." However, the intent requires one reasoning step further to show that "the meme insults Greta Thurnberg as a pushy kid".

## 2.3 Explanation Tasks

The third category of work, like the second category, also involves generating text, but the focus here is generating a textual *explanation* of the mes-

sage conveyed in a meme, as described below.

Sharma et al. (2023) defined the task of generating an explanation of *why* an entity plays the given role in the meme, where the role can be one of "hero", "villain", and "victim". Hee et al. (2023), on the other hand, addressed the task of explaining the reason why a meme is hateful toward a general target group.

Note that these explanation tasks are different from the interpretation tasks. The explanation tasks can be viewed as *constrained* interpretation tasks: in Sharma et al.'s task, both the entity and the role are given, whereas in Hee et al.'s task, the general target is given. In contrast, such constraints are not present in the interpretation task. As an example, consider the meme in Figure 1b again. The final intent that we would have produced for this meme (as the output of interpretation) is "The meme poster makes fun of Trump for the change in his recognition of the severity of the Coronavirus". However, when the target is constrained to be "the Democratic Party", the explanation would be "The Democratic Party is portrayed as a victim of false allegations", which is entirely different in meaning.

### 3 Benchmark Creation

In this section, we will show details about MemeIntent.

#### 3.1 SemEval 2021 Task 6

We chose to annotate the meme collection of SemEval 2021 Task 6 (Dimitrov et al., 2021b). This dataset contains 950 memes, each of which has the image, the text extracted from the image, and the persuasion techniques used. Based on these memes, we built the MemeIntent benchmark. This dataset is favored due to its wide range of opinionated topics, including politics, vaccines, COVID-19, and gender equality. Moreover, each meme in this dataset cannot be properly interpreted without relying on both the visual cues and the textual cues. Therefore, the dataset asks for a 'true' multimodal processing ability in the solutions, as well as the capacity to retrieve relevant world knowledge to interpret contents on such topics.

#### 3.2 Annotation Scheme

Our annotation scheme and procedure is illustrated in Figure 3, while further details are shown in Appendix B. For each meme, we annotate two fields:

- **Intent:** what the author might be trying to do

---

#### Annotation guideline

---

**Intent:** Write one sentence about what the author ultimately wants to do with the meme, as perceived by the annotator. This must be written in good English (complete sentence, with a period at the end)

**BK:** Write the additional knowledge, besides the visible image and text, you needed to use to derive the intent. Examples are information about a public figure or an explanation for a related event.

---

Table 2: Annotation Guidelines.

through the meme (e.g., "The meme encourages people to get vaccinated because they are safe"). A meme can have multiple intents, representing its multiple meanings.

- **Background knowledge (BK):** the knowledge that is not present in the meme, but is needed to recognize the intent when combined with the information that is explicitly stated in the meme. That includes historical knowledge, general political ideologies, or knowledge of the meme culture, etc.

Note that, we allowed multiple intents in a meme to respect the subjectivity of meme interpretation. Built on top of theories from Bach and Harnish (1984), Grundlingh (2018) argued that a meme, like an utterance, could have more than one inference, which depends on the context of communication. Therefore, it is necessary to collect different intents perceived by different annotators, which is a natural consequence of the difference in their backgrounds and personalities. For example, consider the meme in Figure 2b. Depending on how one thinks about gun use, they may interpret the intent of the meme as "accusing Trump of being violent" or "praising Trump for his policies".

Additionally, the annotations include BK to provide extra guidance for learning algorithms in generating intents. As memes usually require a high level of cultural understanding (Milner, 2012), learning systems should be able to store and appropriately retrieve truthful knowledge about the world. The BK was collected to support testing such capabilities.

#### 3.3 Annotation Procedure

Now, we seek to design an annotation procedure to label high quality intents and BKs. To control quality, dataset creators typically maintain *inter-annotator agreement* scores – the higher the score, the more reliable the dataset is (Artstein, 2017).

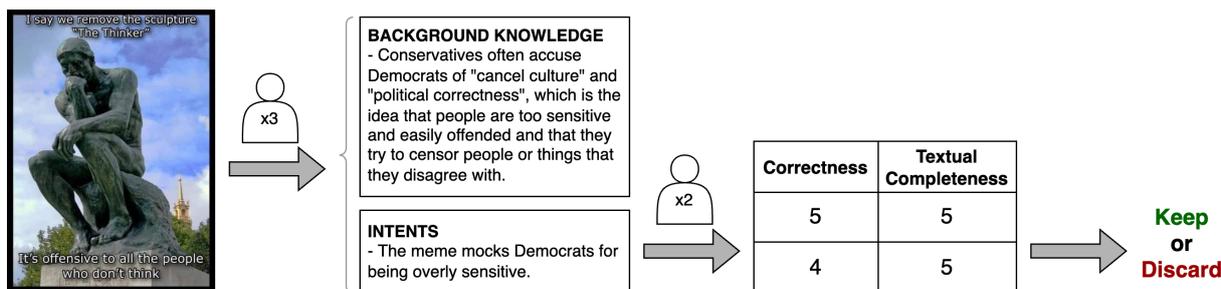


Figure 3: Annotation scheme and procedure of MemeIntent.

In order to obtain such scores, datasets must only involve *categorical* labels, which makes it easy to determine if two annotations agree. However, MemeIntent contains unstructured text annotations where there is no trivial way to check if two sentences (e.g., two intents) are identical. For reference, Hwang and Schwartz (2023) created the MemeCap dataset with only one round of annotations.

We questioned what a reliable procedure to annotate free-text data looks like. Towards that goal, we referred to the work of Wiegrefe and Marasovic (2021), who surveyed 65 papers that produce datasets for explainable NLP. For improving annotation quality, the authors advocate for "a two-stage COLLECT-AND-EDIT" approach, where annotations are first collected (stage 1), and then edited by a new annotator (stage 2). This approach is recommended due to its potential "to increase linguistic diversity via multiple annotators per-instance, reduce individual annotator biases, and perform quality control", and thus has been employed in constructing various free-text datasets (Parikh et al., 2020; Do et al., 2021; Li et al., 2018). Although COLLECT-AND-EDIT does not return any scores at the end, it has been shown to yield high quality annotations without further steps. Therefore, we used COLLECT-AND-EDIT as our annotation procedure.

We recruited five students in computer science, all of whom are native speakers of English, to label the dataset. All annotators went through roughly two hours of initial training and received regular feedback to adhere to our annotation guidelines. Each meme was annotated by three annotators in a sequential COLLECT-AND-EDIT manner: each of the three annotators, given the annotations of the previous person (which was initially empty), could add new intents, add new background knowledge sentences, or modify existing ones, based on their own interpretation of the meme. At the end of this

stage, each meme had one or more unique intents, along with a list of BK sentences that is relevant to the understanding of the meme.

To control quality, we asked two reviewers to review each intent. To avoid biases, those reviewers were made sure to review memes that they did not annotate. The reviews were recorded as answers in the 5-point Likert scale<sup>3</sup> to two questions:

- *Correctness*: How much do you agree that this is the author's intent?
- *Textual Completeness*: How much do you agree that this sentence has complete English writing with good grammar?

We removed all intents that received at least one correctness score lower than 4 from any of the reviewers. If no intents remained for a meme, we would restore the intent(s) with the highest average Correctness score.

Overall, 11.4% of memes in MemeIntent have more than one intent. The mean number of words in the intents is  $10.6 \pm 4.8$ . For background knowledge, the list of BK for each meme has an average of  $1.7 \pm 1.3$  sentences. The mean correctness scores of the intents are  $4.76 \pm 0.37$  (on the 1-5 Likert scale), while the mean Textual Completeness is  $4.54 \pm 0.71$ . These statistics provide suggestive evidence that the COLLECT-AND-EDIT procedure indeed produces high-quality annotations. Appendix A shows further qualitative analysis of the memes in MemeIntent, while Appendix C presents our ethics statement regarding the dataset.

## 4 Experimental Setup

With MemeIntent constructed, we now evaluate the performance of state-of-the-art models on intent generation for memes. The experiments were set up to answer the following research questions (RQs):

<sup>3</sup>The 5-point Likert scale is a numerical scale for recording agreement level, going from 1 (strongly disagree) to 5 (strongly agree).

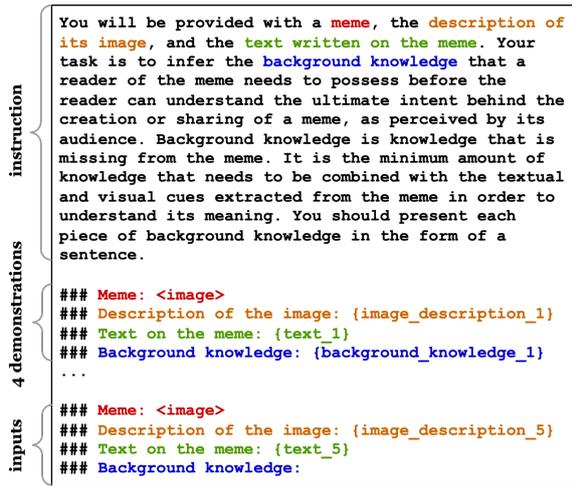


Figure 4: Prompt template used for BK generation in Llava few-shot learning setup. For zero-shot learning, the demonstrations are omitted.

**RQ1:** What is the effect of adding background knowledge to the input on models’ performance in intent description generation for memes?

**RQ2:** How do state-of-the-art models perform in intent description generation?

To that end, we designed a 3-factor experimental setup, consisting of 2 models  $\times$  2 learning setups  $\times$  3 input types, resulting in 12 settings. Finally, models’ outputs from all settings are evaluated automatically and by humans. This section describes those factors and the evaluation metrics.

#### 4.1 Three Input Types

We designed three input types that vary only in the treatment of background knowledge in the input.

In the first type, NoBK, only the *surface* information of the meme is fed to the model, including the meme itself<sup>4</sup>, the extracted text on the meme, an automatically generated caption of the image without the text.

In the second type, AutoBK, we introduced automatically generated BK into the process. In particular, the BK is generated from a different model (BK generation model) in the same setting as the intent (i.e., the same model type and learning setup), using the prompt template in Figure 4. The BK is then fed to the intent description generation model along with surface information to generate the intents.

Finally, the HumanBK type replaces the generated BK with the BK annotated by humans. The

<sup>4</sup>Note that the meme is ignored by LMMs because it does not take images as input.

prompt template for this input type is shown in Figure 5. This input setting is to gauge the upper bound on performance improvement given the human-annotated BK.

#### 4.2 Two Models

Next, we selected two of the best open-sourced models for experiments.

**Vision Language Model (Llava 1.6)** Because intent description generation is a vision-language task, it is natural to use a vision-language model (VLM) to generate intent descriptions. In our experiments, we used Llava 1.6 (Liu et al., 2023), one of the most popular open-source vision language models with state-of-the-art performance in many visual reasoning tasks. We chose the variant llava-v1.6-mistral-7b-hf<sup>5</sup> for its superior performance among the Llava-Next variants with model size no more than 10B. It contains Mistral-7B-Instruct-v0.2<sup>6</sup> as the base language model and CLIP-ViT-L-336px (Radford et al., 2021) as the vision encoder.

**Aided Large Language Model (Llava 1.6 + Llama 3)** We also experimented with a pure large language model (LLM) with the aid of an image captioner. In other words, we employed a two-staged pipeline including (1) image captioning and (2) text-based intent description generation. For image captioning, we again leveraged Llava 1.6 to generate the captions for the memes. These captions, which describe the images themselves, would act as a proxy for the actual image to the LLM<sup>7</sup>. We then used Llama 3<sup>8</sup> to generate intents from the caption and other textual inputs. Llama 3, the most capable open-source LLM as of now (May 2024), has a decoder-only transformer architecture and was trained on 15 trillion tokens from public data. We used the variant Meta-Llama-3-8B-Instruct<sup>9</sup>.

In all experiments, we kept the hyperparameters of the models the same with the default values and only tuned the `max_new_tokens`, setting its final value to 100 for intent description generation and

<sup>5</sup><https://huggingface.co/llava-hf/llava-v1.6-mistral-7b-hf>

<sup>6</sup><https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2>

<sup>7</sup>To be fair with the VLM setting, we also feed the image caption to the VLM model.

<sup>8</sup><https://ai.meta.com/blog/meta-llama-3/>

<sup>9</sup><https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

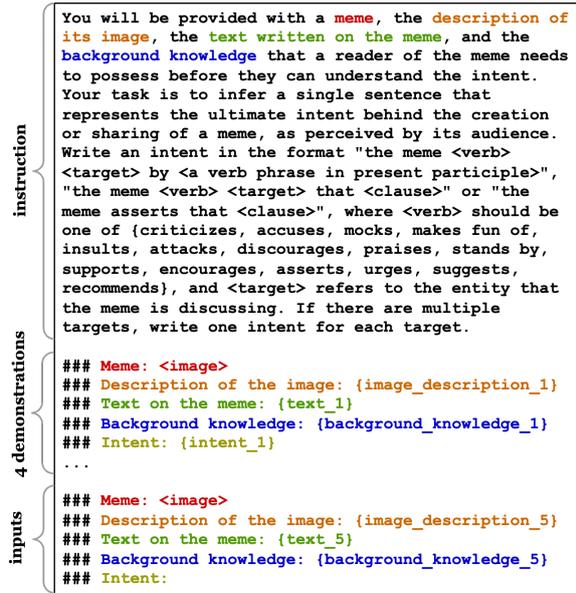


Figure 5: Prompt template used for intent description generation in Llava, with HumanBK input type, in few-shot learning setup. For zero-shot learning, the demonstrations are omitted.

500 for background knowledge generation in both models.

### 4.3 Two Learning Setups

For each of the input types and models, we further experimented with two learning setups: zero-shot and few-shot learning (Mann et al., 2020)<sup>10</sup>. Overall, the two setups differ in the existence of the demonstrations. In the **zero-shot** setup, the *prompt* to the model includes an *instruction* and the *inputs* for the current meme. Meanwhile, in the **few-shot** setup, the prompt also includes 4 *demonstrations*, which are carefully crafted examples of input-output for 4 randomly chosen memes from MemeCap’s test set. We illustrated the prompt used in few-shot learning in Figure 5.

### 4.4 Evaluation Metrics

To automatically evaluate model-outputted intents, we employed four metrics that are commonly used for text generation tasks, namely BLEU-4 (Papineni et al., 2002), ROUGE-L (Lin, 2004), BERT-F1 (Zhang et al., 2020), and SelfCheckGPT-NLI (Manakul et al., 2023). When there are multiple ground-truth intents, we took the maximum (i.e., best) of the scores when comparing the generated

<sup>10</sup>We attempted to fine-tune Llama model on the training data of MemeCap. However, the results turned out to be not as good as zero-shot and few-shot learning. Therefore, we omitted the result in this paper.

Model	Setup	BK	Metrics			
			BLEU	ROU.	BERT	Self.
Llama	zero-shot	No	0.011	0.243	0.89	0.354
		Auto	0.006	0.214	0.884	0.321
		Human	0.014	0.232	0.887	<u>0.475</u>
	few-shot	No	0.015	<u>0.287</u>	<u>0.899</u>	0.339
		Auto	0.013	0.282	0.899	0.34
		Human	<b>0.024</b>	<b>0.312</b>	<b>0.904</b>	0.439
Llava	zero-shot	No	0.006	0.231	0.885	0.352
		Auto	0.004	0.21	0.88	0.405
		Human	0.011	0.255	0.891	<b>0.506</b>
	few-shot	No	0.003	0.214	0.88	0.248
		Auto	0.002	0.134	0.867	0.459
		Human	0.003	0.225	0.883	0.313

Table 3: Automatic Evaluation Results on Intent Description Generation. For each metric, the overall best results are in **bold**, while the second best results are underlined. Abbreviations: Setup (learning setup), BK (input types), No (NoBK), Auto (AutoBK), Human (HumanBK), ROU (ROUGE-L), BERT (BERT-F1), Self (SelfCheckGPT-NLI).

intent with the ground truths<sup>11</sup>. When making comparisons between settings on a metric, we performed the two-sided T-test with significant level  $\alpha = 0.05$ . Finally, we conducted human evaluation on the outputs of some selected settings to verify observations from automatic evaluation.

## 5 Results and Discussion

### 5.1 Automatic Evaluation

Table 3 shows the automatic evaluation results on the generated intents of the two models across all learning setups and input types. Meanwhile, Table 4 reports the corresponding results for the generated BK in AutoBK settings by calculating similarity scores with the human-annotated background knowledge.

**Input types (RQ1).** The central observation from our experiment is that background knowledge is crucial to the performance of intent description generation. Specifically, for most settings, HumanBK input type gave the statistically highest performance across metrics. More interestingly, NoBK usually gave better performances than AutoBK, except in Llama few-shot, where there is no statistical significance. There are a few exceptions to this rule: SelfCheckGPT-NLI gave higher scores

<sup>11</sup>For SelfCheckGPT-NLI, we assigned  $score \leftarrow 1 - score$  to be consistent with the other metrics that the higher the score is, the closer the two pieces of text are.

Model	Setup	Metrics			
		BLEU	ROU.	BERT	Self.
Llama	zero-shot	0.003	0.073	0.827	0.331
	few-shot	<b>0.008</b>	<b>0.127</b>	<b>0.844</b>	0.294
Llava	zero-shot	<u>0.003</u>	<u>0.086</u>	0.83	0.384
	few-shot	<u>0.002</u>	0.072	0.821	<b>0.392</b>

Table 4: **Automatic Evaluation Results on Background Knowledge Generation in AutoBK setting.** For each metric, the overall best result is in **bold**, and the second best is underlined. Abbreviations: Setup (learning setup), ROU (ROUGE-L), BERT (BERT-F1), Self (SelfCheckGPT-NLI).

for AutoBK than NoBK where Llava was used, and NoBK sometimes outperformed HumanBK in Llama zero-shot experiments (on ROUGE-L and BERT-F1).

These results show that using human-annotated BK during the process produces better performance than no BK or auto-generated BK. We will further investigate these results via human evaluation (Section 5.2).

**Performance on Background knowledge generation (RQ1).** We take a closer look at the performance of models in BK generation. On BLEU, ROUGE-L, and BERT-F1, few-shot is better than zero-shot for Llama, and the opposite happens for Llava. However, SelfCheckGPT-NLI flips those results for both models.

We can connect these results with intent generation performance in AutoBK settings. In fact, among the AutoBK settings in Table 3, those with the best BK generation scores also score the highest on intent description generation. This suggests a correlation between the performance of BK generation and that of intent description generation across settings.

**Models (RQ2).** On BLEU, ROUGE-L, and BERT-F1, Llama (aided by Llava’s image captions) outperformed Llava alone for most of the settings – across input types and learning setups. However, three over four metrics<sup>12</sup> gave a statistically higher score for Llava in experiments where zero-shot and NoBK input were used. Besides, SelfCheckGPT-NLI gave statistically higher scores for Llava when where AutoBK was used. Therefore, none of these models can entirely outperform the other across settings.

<sup>12</sup>except BLEU which did not show statistical significance

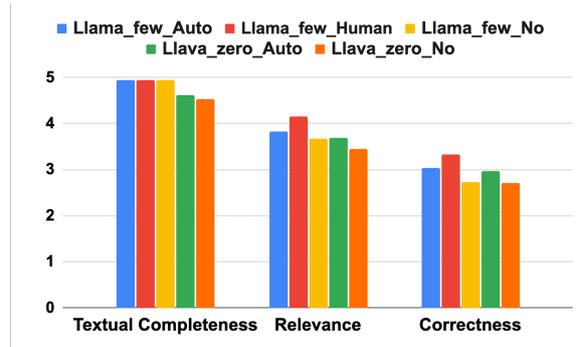


Figure 6: **Human Evaluation Results on Intent Description Generation.**

**Learning setup (RQ2).** In general, when Llama was used, few-shot was better than zero-shot. Conversely, the opposite happened when Llava was used. The superior performance of few-shot learning in Llama is aligned with the intuition that having demonstrations is useful. Meanwhile, Llava’s inferior few-shot performance has been discussed by its authors that Llava was not explicitly trained to take multiple images as input<sup>13</sup>.

The general trend above does have a few exceptions: SelfCheckGPT gave statistically higher scores for few-shot learning in Llava AutoBK, and zero-shot learning in Llava HumanBK.

## 5.2 Human Evaluation

For human evaluation, we evaluated the model outputs on 30 randomly chosen memes. Two annotators evaluated the outputs along three dimensions: **Textual Completeness** (i.e., How much do you agree that this sentence has complete English writing with good grammar?), **Relevance** (i.e., How relevant the sentence is to the topic of the meme?), and **Correctness** (i.e., How much do you agree that this is the author’s intent?). Answers were recorded in the 5-point Likert scale, where higher scores indicate better quality.

To select settings for evaluation, we first focused on the effect of leveraging background knowledge to enhance the prediction of a meme’s intent. Noticing the superior performance of Llama few-shot in most metrics, we selected its three settings – NoBK, AutoBK, and HumanBK – for human evaluation.

Next, the automatic evaluation results showed that in Llava NoBK and AutoBK settings, NoBK scored higher on BLEU, ROUGE, and BERT metrics; however, AutoBK scored higher on

<sup>13</sup>[https://huggingface.co/docs/transformers/en/model\\_doc/llava#usage-tips](https://huggingface.co/docs/transformers/en/model_doc/llava#usage-tips)

SelfCheckGPT-NLI. We know that NLI measures were used to assess the faithfulness of summarization, focusing on analyzing textual entailment between the context and the summary (Maynez et al., 2020). Given the contradiction between SelfCheckGPT-NLI and other metrics, we evaluated Llava zero-shot to determine whether SelfCheckGPT-NLI accurately captures the correctness between two sentences. We selected this setting since it demonstrated better scores among the two Llava settings.

Results are shown in Figure 6. Firstly, among the Llama few-shot settings, HumanBK significantly outperformed all other settings across all three metrics, which agrees with the automatic evaluation. Furthermore, while there was no statistical significance between AutoBK and NoBK in automatic evaluation, the human evaluation showed that AutoBK exhibits a higher performance than NoBK. These further demonstrate that a more sophisticated BK can influence the performance of intent generation.

Secondly, upon examining two outputs from the Llava model, it is observed that the performance of AutoBK surpasses that of NoBK across all three metrics in human evaluation. This is consistent with the SelfCheckGPT-NLI score, indicating that this metric effectively captures the correctness between the two sentences in our experiments.

## 6 Conclusion and Future Work

We examined the novel task of generating the description of intents in memes, specifically by (1) constructing MemeIntent, a benchmark of memes with intents, and background knowledge and (2) producing baseline results on our dataset against which future models can be compared. Our key findings suggest the importance of background knowledge treatments in intent description generation. To stimulate research on this task, we make our annotations publicly available.

Regarding future work, the experimental results w.r.t. the models and zero-shot vs. few-shot are inconclusive. Therefore, more experimentation is needed to get a clearer picture. Besides, we attempted to fine-tune Llama on the training set of MemeCap and test on MemeIntent, but the result was not good. This seems to be a failure to generalize from one meme interpretation dataset to another. Therefore, more efforts should be put into looking at the discrepancies between current datasets.

## References

- Ron Artstein. 2017. [Inter-annotator Agreement](#). In Nancy Ide and James Pustejovsky, editors, *Handbook of Linguistic Annotation*, pages 297–313. Springer Netherlands, Dordrecht.
- Kent Bach and Robert M. Harnish. 1984. *Linguistic communication and speech acts*, 1.ed., 2. print edition. MIT Press, Cambridge, Mass.
- Mohit Chandra, Dheeraj Pailla, Himanshu Bhatia, Aadilmehdi Sanchawala, Manish Gupta, Manish Shrivastava, and Ponnurangam Kumaraguru. 2021. [“Subverting the Jewtocracy”: Online Antisemitism Detection Using Multimodal Deep Learning](#). In *Proceedings of the 13th ACM Web Science Conference 2021, WebSci ’21*, pages 148–157, New York, NY, USA. Association for Computing Machinery.
- Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021a. [SemEval-2021 task 6: Detection of persuasion techniques in texts and images](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 70–98, Online. Association for Computational Linguistics.
- Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021b. [SemEval-2021 Task 6: Detection of Persuasion Techniques in Texts and Images](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 70–98, Online. Association for Computational Linguistics.
- Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. 2021. [e-snli-ve: Corrected visual-textual entailment with natural language explanations](#).
- Elisabetta Fersini, Francesca Gasparini, Giulia Rizzi, Aurora Saibene, Berta Chulvi, Paolo Rosso, Alyssa Lees, and Jeffrey Sorensen. 2022. [SemEval-2022 Task 5: Multimedia Automatic Misogyny Identification](#). In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 533–549, Seattle, United States. Association for Computational Linguistics.
- Lezandra Grundlingh. 2018. [Memes as speech acts](#). *Social Semiotics*, 28(2):147–168.
- Ming Shan Hee, Wen-Haw Chong, and Roy Ka-Wei Lee. 2023. [Decoding the Underlying Meaning of Multimodal Hateful Memes](#). In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pages 5995–6003, Macau, SAR China. International Joint Conferences on Artificial Intelligence Organization.
- EunJeong Hwang and Vered Shwartz. 2023. [MemeCap: A dataset for captioning and interpreting memes](#).

- In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1433–1445, Singapore. Association for Computational Linguistics.
- Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. [The Hateful Memes Challenge: Detecting Hate Speech in Multimodal Memes](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 2611–2624. Curran Associates, Inc.
- Qing Li, Qingyi Tao, Shafiq Joty, Jianfei Cai, and Jiebo Luo. 2018. [Vqa-e: Explaining, elaborating, and enhancing your answers for visual questions](#).
- 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.
- Chen Liu, Gregor Geigle, Robin Krebs, and Iryna Gurevych. 2022. [FigMemes: A Dataset for Figurative Language Identification in Politically-Opinionated Memes](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7069–7086, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. [Improved baselines with visual instruction tuning](#).
- Potsawee Manakul, Adian Liusie, and Mark Gales. 2023. [SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017, Singapore. Association for Computational Linguistics.
- Ben Mann, N Ryder, M Subbiah, J Kaplan, P Dhariwal, A Neelakantan, P Shyam, G Sastry, A Askell, S Agarwal, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 1.
- Lambert Mathias, Shaoliang Nie, Aida Mostafazadeh Davani, Douwe Kiela, Vinodkumar Prabhakaran, Bertie Vidgen, and Zeerak Waseem. 2021. [Findings of the WOAHS 5 Shared Task on Fine Grained Hateful Memes Detection](#). In *Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021)*, pages 201–206, Online. Association for Computational Linguistics.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. [On faithfulness and factuality in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Ryan M. Milner. 2012. *The World Made Meme: Discourse and Identity in Participatory Media*. Ph.D. thesis, University of Kansas.
- Philipp Mondorf and Barbara Plank. 2024. [Beyond accuracy: Evaluating the reasoning behavior of large language models – a survey](#).
- 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.
- Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. [ToTTo: A controlled table-to-text generation dataset](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1173–1186, Online. Association for Computational Linguistics.
- Shraman Pramanick, Dimitar Dimitrov, Rituparna Mukherjee, Shivam Sharma, Md. Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2021a. [Detecting harmful memes and their targets](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2783–2796, Online. Association for Computational Linguistics.
- Shraman Pramanick, Dimitar Dimitrov, Rituparna Mukherjee, Shivam Sharma, Md. Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2021b. [Detecting Harmful Memes and Their Targets](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2783–2796, Online. Association for Computational Linguistics.
- Shraman Pramanick, Shivam Sharma, Dimitar Dimitrov, Md. Shad Akhtar, Preslav Nakov, and Tanmoy Chakraborty. 2021c. [MOMENTA: A Multimodal Framework for Detecting Harmful Memes and Their Targets](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4439–4455, Punta Cana, Dominican Republic. 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.
- Chhavi Sharma, Deepesh Bhageria, William Scott, Srinivas PYKL, Amitava Das, Tanmoy Chakraborty, Viswanath Pulabaigari, and Björn Gambäck. 2020. [SemEval-2020 task 8: Memotion analysis- the visuo-lingual metaphor!](#) In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 759–773, Barcelona (online). International Committee for Computational Linguistics.
- Shivam Sharma, Siddhant Agarwal, Tharun Suresh, Preslav Nakov, Md. Shad Akhtar, and Tanmoy

Chakraborty. 2023. [What Do You MEME? Generating Explanations for Visual Semantic Role Labelling in Memes](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(8):9763–9771.

Shivam Sharma, Tharun Suresh, Atharva Kulkarni, Himanshi Mathur, Preslav Nakov, Md. Shad Akhtar, and Tanmoy Chakraborty. 2022. [Findings of the CONSTRAINT 2022 Shared Task on Detecting the Hero, the Villain, and the Victim in Memes](#). In *Proceedings of the Workshop on Combating Online Hostile Posts in Regional Languages during Emergency Situations*, pages 1–11, Dublin, Ireland. Association for Computational Linguistics.

Shardul Suryawanshi, Bharathi Raja Chakravarthi, Mihael Arcan, and Paul Buitelaar. 2020a. [Multimodal meme dataset \(MultiOFF\) for identifying offensive content in image and text](#). In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*, pages 32–41, Marseille, France. European Language Resources Association (ELRA).

Shardul Suryawanshi, Bharathi Raja Chakravarthi, Pranav Verma, Mihael Arcan, John Philip McCrae, and Paul Buitelaar. 2020b. [A Dataset for Troll Classification of TamilMemes](#). In *Proceedings of the WILDRE5– 5th Workshop on Indian Language Data: Resources and Evaluation*, pages 7–13, Marseille, France. European Language Resources Association (ELRA).

Sarah Wiegrefe and Ana Marasovic. 2021. [Teach Me to Explain: A Review of Datasets for Explainable Natural Language Processing](#). 1.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [BERTScore: Evaluating Text Generation with BERT](#). ArXiv:1904.09675 [cs].

## A Challenges in Interpretation

We conduct a manual analysis of two memes taken from our dataset, with the goal of understanding the challenge of interpreting memes.

In Figure 7a, one first sees a person in a colorful outfit (via the jacket, the glasses, the hair). Then they may infer that this is an LGBT person. After that, they read the text saying ‘*Trump scares me*’. To connect that with the image, they further interpret the emotion of the LGBT person and recognize that they are apparently nonchalant. This is a word-face contrast, which suggests there is something wrong with one of the two. If the facial expression is ‘wrong’, one knows that the LGBT person may have a problem expressing fear, and the intent is to criticize Trump for being a scary person. But the other interpretation is more probable – that the words are wrong. Then, the LGBT person is actually not scared, thus being over-sensitive. This



Figure 7: Some example memes from the SemEval-2021 Task 6 dataset.

interpretation path might be triggered by the unconventional outlook of the person, which typically *scares* people, thus making them think the person in the meme is a bad person. Finally, because the liberals in the US support LGBT rights, this line of interpreting leads to a more significant intent, that is *mocking the liberals as over-sensitive and scary*.

Consider another example in Figure 7b. A reader first sees a lion biting the zebra, about to kill it; then a lion looking at a hedgehog with upright quills, not sure if it is attacking or not. Then they read the first text saying ‘*unarmed victim*’. *Victim* refers to the zebra, and *unarmed* is a word for humans, so this is a metaphor for unarmed people being attacked. This line of thought triggers the reader’s knowledge about the constant debate over gun control policies in the US. Next, the phrase ‘*armed victim*’ with the word *victim* crossed suggests that the hedgehog, or metaphorically the gun owner, is safe. Finally, the rhetorical question ‘*Any questions?*’ conveys that this is clear evidence so that *gun use should be allowed with no doubt*. Along this line of reasoning, the fact that the zebra is violently bitten provokes fear in the reader, which urges them to become the hedgehog and get a gun for self-defense.

In both of these examples, sophisticated logical (*logos*) and emotional (*pathos*) processes have triggered each other, forming the most probable interpretation path that leads to the recognition of the intent. This is a tricky task that only humans with appropriate knowledge and experience can perform. In fact, logical reasoning requires sufficient *background knowledge* to have the right facts to start with (e.g., how an LGBT person usually looks like, that liberals support LGBT rights, gun control is debated). Moreover, humans are also easily triggered by emotional stimuli (e.g., a strange look is scary, and safety is important). Those emotions are two-fold – they ‘disambiguate’ multiple possible

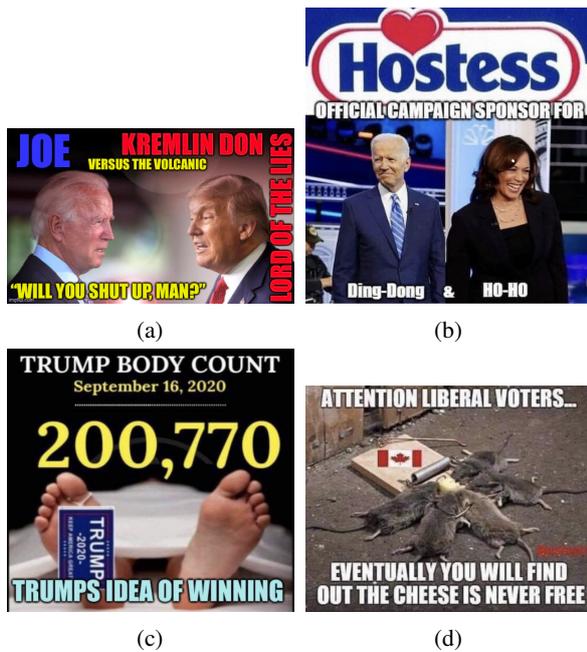


Figure 8: Example memes used in our annotation guidelines

interpretations via psychological biases, while also reinforcing the intent through pathos. This orchestrated effect can form a logical fallacy in disguise to achieve the final intent.

## B Annotation Details

This section shows details about our annotation procedure.

### B.1 Guidelines For Intents

This task introduces the notion of *intent*. An intent of a meme is what the meme author ultimately wants to do with the meme, perceived by the audience.

For example, the meme in Figure 8a has one final intent, which is: [The meme] praises Biden for being a better leader than Trump.

#### B.1.1 Frequently Asked Questions

**Can I write new Intents?** Yes, you should! If you think some intent is missing, add it.

**How to write the Intent?** Write an intent in format "[<verb1> <target1> <etc.>] x n", so that this sentence when being prefixed by "The meme" will form a grammatically correct sentence. For example, write "insults Trudeau for lying and insults anyone who believes in him as stupid".

**What to do with multiple intents?** Rank them by decreasing order of preference (i.e., from what you believe the most to the least).

### B.2 Guidelines For Background Knowledge

For BK, summarize the key background information (skipping trivial knowledge). Also, write down (1) what you don't know but seem important, and (2) what you are not sure if it's right. Use question marks for those, e.g., "The place is in Cuba?".

For example:

- (Figure 8b) Hostess is a company that has products named Ding-Dong and Ho-Ho. Ding-Dong is also used to refer to someone who is slow. Ho-ho is also used to referred to someone who is lustful.
- (Figure 8c) A lot of Americans died during Covid 19, when Trump was presiding over the United States.
- (Figure 8d) This is about the Canadian election. Canada has 2 parties, one of which is the Liberal party.

## C Ethics Statement

**Broader implications.** As mentioned before, the solution to the intent description generation task is of practical significance. From a practical perspective, knowledge of the message being conveyed in a meme could be useful for other meme-related processing tasks. For instance, knowing what the message is could facilitate the determination of whether a meme contains harmful content. Theoretically speaking, being able to generate messages like humans requires that a machine read between the lines and achieve a deeper level of understanding of perceptual input, enabling machine perception to get one step closer to human perception.

**Ethical considerations.** Having said that, we are all aware that some memes contain harmful content, so when our models are applied to these harmful memes, they will make an intent that is harmful explicitly. The resulting message could have a negative psychological impact on the users, especially if they are the target of the harmful content. Therefore, as with many other AI/NLP technologies, our models should be used with care. We should emphasize that our intent is to build models for generating the messages conveyed in memes, hoping that readers of memes will be less likely to

be manipulated after understanding the messages being conveyed.

**Human annotator information.** All annotators were hired during Fall 2023 - Spring 2024 as student workers (15-20 hours/week) with full consent. All of the annotators were undergraduate and graduate students in computer science aged around 18-24. The group comprised both male and female students with members from Asian ethnicity, with fluent to native English level.

**Steps taken to protect annotators from harmful content.** All annotators were provided with a thorough instructional training session in which they were instructed on how to annotate the data and how to go about the whole task. During training, annotators were shown the types of memes that they would work with so that they have an idea of the dataset's nature. The annotators have full autonomy to withdraw from the project at their own judgment.

**Terms of use.** This dataset is consistent with the terms of use and the intellectual property and privacy rights of people. There is nothing about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses.

**Data distribution.** We have open-sourced the data produced from this work. It is released on a GitHub repository with the MIT license.

# Automating PTSD Diagnostics in Clinical Interviews: Leveraging Large Language Models for Trauma Assessments

Sichang Tu,<sup>1</sup> Abigail Powers,<sup>1</sup> Natalie Merrill,<sup>1</sup> Negar Fani,<sup>1</sup>

Sierra Carter,<sup>2</sup> Stephen Doogan,<sup>3</sup> Jinho D. Choi<sup>1</sup>

<sup>1</sup>Emory University, Atlanta, GA, USA

<sup>2</sup>Georgia State University, Atlanta, GA, USA

<sup>3</sup>Doogood Foundation, New York, NY, USA

{sichang.tu, abigail.d.powers, natalie.merrill, nfani,jinho.choi}@emory.edu

scarter66@gsu.edu, sdoogan@rlsciences.com

## Abstract

The shortage of clinical workforce presents significant challenges in mental healthcare, limiting access to formal diagnostics and services. We aim to tackle this shortage by integrating a customized large language model (LLM) into the workflow, thus promoting equity in mental healthcare for the general population. Although LLMs have showcased their capability in clinical decision-making, their adaptation to severe conditions like Post-traumatic Stress Disorder (PTSD) remains largely unexplored. Therefore, we collect 411 clinician-administered diagnostic interviews and devise a novel approach to obtain high-quality data. Moreover, we build a comprehensive framework to automate PTSD diagnostic assessments based on interview contents by leveraging two state-of-the-art LLMs, GPT-4 and Llama-2, with potential for broader clinical diagnoses. Our results illustrate strong promise for LLMs, tested on our dataset, to aid clinicians in diagnostic validation. To the best of our knowledge, this is the first AI system that fully automates assessments for mental illness based on clinician-administered interviews.

## 1 Introduction

Mental health has become a vital element of overall well-being. The prevalence of mental illness poses, however, a critical challenge to healthcare, underscoring the urgent need for an increased capacity of mental health services. Only 29% of people with psychosis receive formal care, leaving a significant portion completely untreated (WHO: [World Health Organization \(2021\)](#)). Aside from obstacles such as high costs, limited awareness, and stigma surrounding mental health, the shortage of the mental health workforce has been a major factor exacerbating this gap. According to WHO, the average ratio of mental health workers per 100,000 population was 13, making it difficult for people to access reliable and readily administered mental health diagnostics, as well as subsequent support and interventions.

The emergence of Large Language Models (LLMs) has suggested innovative solutions to this challenge. Several studies have explored LLM applications in mental health for condition detection ([Zhang et al., 2022](#)), support and counseling ([Ma et al., 2023b](#)) as well as clinical decision-making ([Fu et al., 2023](#)), and shown the feasibility for LLMs to enhance the workforce of mental healthcare ([Hua et al., 2024](#)). By harnessing LLMs' ability to interpret languages that involve high expertise, it is possible to mitigate the service gap in the healthcare ecosystem through the automation of condition detection and diagnosis without the need of training so many professionals, which is both costly and time-consuming.

Despite these advancements, notable limitations persist in the current research on automatic diagnosis for mental health. Most studies have focused on prevalent conditions like stress ([Lamichhane, 2023](#)) and depression ([Qin et al., 2023](#)), with scant attention to less common but more severe conditions like Post-traumatic Stress Disorder (PTSD). Moreover, while prior studies have leveraged data from social media, clinical notes, and electronic health records, very few have utilized clinical interviews, and even in those cases, they rely on basic self-administered scales estimated in dialogues between computers and patients ([Galatzer-Levy et al., 2023](#)). No work has employed diagnostic interviews between real clinicians and patients that are systematically conducted, resulting in a dearth of practical research on the automatic diagnosis of mental illness.

In this paper, we present an LLM-based system that listens to hour-long conversations between clinicians and patients and performs diagnostic assessments for PTSD. Our final model is evaluated by clinicians specialized in PTSD, suggesting a great potential for LLMs while highlighting certain limitations (Section 6). Our primary contributions are:<sup>1</sup>

<sup>1</sup>Our final model is publicly available through our open-source project at <https://github.com/emorynlp/TraumaNLP>.

- A new dataset comprising over 700 hours of interviews between clinicians and patients is created. Every interview consists of multiple diagnostic sections, featuring a series of questions and corresponding assessments from clinicians based on the interview contents (Section 3).
- A novel and comprehensive pipeline is developed to process the interview dataset, so it can be used to build automatic assessment models on PTSD, which can be easily adapted to a broad range of diagnostic interviews (Section 4).
- Assessment models achieving promising results are developed using two state-of-the-art LLMs, showcasing LLMs' ability to answer diagnostic questions through information extraction and text summarization on the interviews (Section 5).

To the best of our knowledge, this is the inaugural system designed to conduct diagnostic assessments on mental health while interpreting real-world interviews administered by clinicians. We believe that this work will foster clinical collaboration between human experts and Artificial Intelligence, thus promoting equitable access to appropriate care for all populations affected by mental illness.

## 2 Related Work

Pre-trained language models have been widely applied in many healthcare tasks (Enghardt et al., 2023; Hu et al., 2023; Peng et al., 2023; Ma et al., 2023a; Liu et al., 2023a). The emergence of LLMs has introduced new capabilities and innovations in healthcare to this domain (Nori et al., 2023; Cascella et al., 2023). This section introduces the related research of LLMs and their applications in healthcare, particularly in mental health.

### 2.1 LLMs in Mental Health

The advent of LLMs like GPT (OpenAI, 2023), Llama (Touvron et al., 2023), and PaLM (Chowdhery et al., 2022) has sparked research into their applications in mental health (Ji et al., 2023). One key area is using conversational agents for mental health support and counseling, where LLMs excel at generating empathetic responses (Lai et al., 2023; Ma et al., 2023b; Loh and Raamkumar, 2023), highlighting their potential as digital companions or on-demand service providers. Additionally, the research on decision-support systems for novice counselors underscores their potential to enhance mental healthcare provision (Fu et al., 2023).

Research has also explored LLMs in disease detection and diagnosis (Zhang et al., 2022), focusing on issues like depression (Qin et al., 2023), stress (Lamichhane, 2023), and suicidality (Bhaumik et al., 2023). Closer to our work, Bartal et al. (2023) use text-based narratives from new mothers to assess childbirth-related PTSD with GPT and neural network models. Although GPT showed moderate performance, it holds promise for clinical diagnosis with further refinement. These studies typically use zero/few-shot prompting for binary or multi-label classification, demonstrating LLMs' capabilities in detecting mental health issues without fine-tuning, despite challenges like unstable responses, potential bias, and interpretation inaccuracies.

Some research has pivoted towards fine-tuning LLMs for domain-specific performance enhancement. Xu et al. (2023) present two fine-tuned models, Mental-Alpaca and Mental-FLAN-T5, outperforming GPT-3.5 and GPT-4 in multiple mental health prediction tasks. Based on Llama-2, Yang et al. (2023) train MentaLLaMA on 105K social media data enhanced by GPT. The model performance is on par with other state-of-the-art methods, while providing interpretable analysis.

### 2.2 LLMs in Clinical Interview and Diagnosis

Research on using LLMs on clinical interview data and diagnosis is limited. Wu et al. (2023) utilize GPT to augment the Extended Distress Analysis Interview Corpus by generating a new dataset from provided profile and rephrasing existing data. The augmented data outperforms the original imbalanced data in PTSD diagnosis. Galatzer-Levy et al. (2023) adopt Med-PaLM-2 to predict Major Depression Disorder (MDD) and PTSD on eight item Patient Health Questionnaire and PTSD Checklist-Civilian version ratings.

## 3 PTSD Interview Data

This study utilizes data from diagnostic interviews administered as part of a larger study on risk and resiliency to the PTSD development in a population seeking medical care (Gluck et al., 2021). Participants were recruited from waiting rooms in primary care, gynecology and obstetrics, and diabetes medical clinics at a publicly funded, safety-net hospital. Data were collected from 2012 to 2023, and inclusion criteria were ages between 18 and 65 with the capacity to provide informed consent. The parent study was conducted according to the latest version

Section	Questions	Variables	Example Question	Example Variable
LBI	31	15	What has been your primary source of income over the past month?	lbi_a1
THH	39	20	In the past, have you been treated for any emotional or mental health problems with therapy or hospitalization?	thh_tx_yesno
CRA	17	20	What would you say is the one that has been most impactful where you are still noticing it affecting you?	critaprobenotes
CAP	241	92	In the past month, have you had any unwanted memories of the [Event] while you were awake, so not counting dreams?	dsm5capscribtb01 trauma1_distress

Table 1: Statistics and examples for each of the four sections employed in this study.

of the Declaration of Helsinki (World Medical Association, 2013), and consent from the participants was obtained after explaining the procedures. The informed consent was approved by our Institutional Review Board and Research Oversight Committee.

### 3.1 Participants

Participants were paid \$60.00 for this interview and underwent semi-structured diagnostic interviews conducted by doctoral-level clinicians or doctoral students supervised by a licensed clinical psychologist on staff. A total of 411 interviews were conducted with 336 unique participants, some of whom had follow-up interviews after >1 month. 93.4% of the participants were women and 79.5% were Black or African American ( $M_{age} = 31.4$ ), where 38.7% had a high school education or less and 57.9% reported a monthly household income of < \$1,000.

### 3.2 Interview Procedures

The diagnostic interview begins with a section of the Longitudinal Interval Follow-Up Evaluation to assess global adaptive functioning across various psychosocial domains, including work, household, relationship as well as general functioning, and life satisfaction in the past month (Keller et al., 1987). Videos of the interviews are recorded using online conferencing software such as Zoom and Microsoft Teams. Each interview lasts 1.5 hours on average, involving the participant and 1-2 interviewers.

### 3.3 Psychiatric Diagnoses and Treatment

A total of 10 sections are applied during the interview. Among them, 4 sections are administered to the majority of participants; thus, this study focuses on those 4 sections. The first two sections, the Life Base Interview (LBI) and the Treatment History & Health (THH), are internally designed to assess the history of psychiatric diagnoses and treatment, as well as the presence of suicidality. The other two sections, the Criterion A (CRA) and the Clinician-

Administered PTSD Scale for DSM-5 (CAP), follow the standard diagnostic criteria for PTSD outlined in the Diagnostic and Statistical Manual of Mental Disorders (DSM-5; Weathers et al. (2018)). Every section is accompanied by a set of questions, linked to variables that store pertinent values derived from the corresponding answers. Table 1 shows statistics and examples for each of the 4 sections.<sup>2</sup>

**LBI** It assesses the participant’s functioning over the past month, addressing topics such as daily life, work, relationships with friends and family, and overall life satisfaction.

**THH** It covers the participant’s treatment/health history, including past physical and mental conditions as well as treatments received, such as medication and therapeutic services.

**CRA** It assesses whether the participant has been exposed to (threatened) death, serious injury, or sexual violence, with a focus on potential traumatic experiences the participant might have endured.

**CAP** It centers on issues the participant may have encountered due to traumatic events, including distress, avoidance of trauma-related stimuli, negative thoughts and feelings, and trauma-related arousal.

## 4 Data Processing

Every video is converted into an MP3 audio file and *transcribed* by two automatic speech recognizers, whose results are *aligned* to produce a high-quality transcript. The transcript is *segmented* into multiple sections based on the relevant questions, and each question is *paired* with its assessment result.

### 4.1 Transcription

Two commercial tools, Rev AI<sup>3</sup> and Azure Speech-to-Text<sup>4</sup>, and an open-source tool, OpenAI Whis-

<sup>2</sup>Descriptions of all 10 sections are provided in Appendix A.

<sup>3</sup>Rev AI: <https://www.rev.ai>

<sup>4</sup>Azure Speech-to-Text: <https://bit.ly/42r24pA>

per (Radford et al., 2023), are tested for automatic speech recognition (ASR) on our dataset. Whisper gives the lowest Word Error Rate (WER; Klakow and Peters (2002)) of 0.13, compared to 0.21 and 0.16 from Rev AI and Azure, respectively. Whisper also exhibits better performance in handling noisy environments and numbers that Azure often misses or inaccurately transcribes (Table 2). Despite its superior ASR performance, Whisper does not identify speakers, a feature found in the others. Thus, both Azure and Whisper are run on all audios and their results are combined to obtain the best outcomes.

Tool	Examples
Azure	(1) I got 2020 on the 24 with three. Three will be 3 is turning 2116, one 15211. (2) They happened in 2017 and I'll be 60 next month, so 5556 something like that.
Whisper	(1) I got two to be 20 on the 24th, well, three, three is turning 20, one 16, one 15, two 11. (2) That happened in 2017 and I'll be 60 next month, so. 55, 56, something like that.

Table 2: Comparisons between Azure and Whisper transcripts, with equivalent tokens coded in matching colors.

## 4.2 Alignment

To map the speaker diarization (SD) output from Azure to the Whisper output, Align4D<sup>5</sup> is used such that the first and last words of every utterance in the Azure output are aligned to their corresponding words in the Whisper transcript with speaker info, and form a speaker turn spanning all words between those words. Some words in the Whisper transcript may get left out from this mapping, which are combined with either preceding or following adjacent utterances using heuristics.

Text-based Diarization Error Rate (TDER; Gong et al. (2023)) is used, more suitable than traditional metrics like WER or Diarization Error Rate (DER; Fiscus et al. (2006)), for evaluating text-based SD. Transcripts from 29 audios produced by Microsoft Teams are used as the gold-standard, where Teams identifies speakers via different audio channels with near-perfect SD. Our aligned method achieves a TDER of 0.56, a significant improvement over the TDER of 0.62 achieved by Azure alone.

## 4.3 Segmentation

Each interview is conducted through multiple sections comprising a series of questions (Section 3.3), yet recorded as one continuous video. It is crucial to segment the video into sections, each of which is

<sup>5</sup>Align4D: <https://github.com/emorynlp/align4d>

split into sessions, where a session contains content relevant to a specific question. Here, a session is defined as a list of utterances where the first utterance includes the corresponding question, and it is followed by another session whose first utterance includes the next question (if it exists). Algorithm 1 describes how a section is matched in the transcript.

---

### Algorithm 1: section\_match( $U, Q^c$ )

---

**Input:**  $U$ : a list of utterances,  $Q$ : a list of questions.

**Output:** An ordered list of tuples comprising utterance IDs and their matching scores.

```

1  $S \leftarrow \text{similarity\_matrix}(U, Q^c)$ ;
2  $T \leftarrow [\max(S_{*,i}) : 1 \leq i \leq |Q^c|]$ ;
3 if average( $T$ ) > 0.6 and
4   ( $|\text{select}(T, 0.8)| \geq 3$  or  $|\text{select}(T, 0.9)| \geq 2$ )
5   then return sequence_alignment( $S$ );
6 return  $\emptyset$ 
```

---

Let  $U$  be a list of utterances, and  $Q^c$  a list of core questions for a specific section.<sup>6</sup>  $S \in \mathbb{R}^{|U| \times |Q^c|}$  is created, where  $S_{i,j}$  is a similarity score between  $u_i \in U$  and  $q_j \in Q^c$  (L1).  $T \in \mathbb{R}^{|Q^c|}$  is then created by selecting the maximum similarity score for every question (L2). Given a function  $\text{select}(T, s)$  that returns a list of scores in  $T$  greater than  $s$ , the section is matched if  $T$ 's average score is > 0.6 (L3) and if there exist at least 3 or 2 questions whose matching scores are > 0.8 or 0.9, respectively (L4). If the section is matched, Gong et al. (2023)'s sequence alignment algorithm is applied to  $S$ , which returns an ordered list of utterance IDs and their matching scores for questions in  $Q^c$ ; otherwise, it returns an empty list (L5). In our case, Sentence Transformer is used to create embeddings for utterances & questions (Reimers and Gurevych, 2019), and cosine similarity is used to estimate the scores. Overlap between spans of two sections may occur due to incorrect matching. Algorithm 2 shows how to remove such overlaps. Let  $Q_i^c$  be a list of core questions for the  $i$ 'th section, and  $R_i = \text{sm}(U, Q_i^c)$  (sm: section\_match). Given  $(R_1, R_2)$ ,  $R_1'$  is created by taking a subset of  $R_1$  whose utterance IDs exist in  $R_2$  (L1), and  $R_2'$  is created similarity (L2). If  $R_1'$  contains more questions with scores > 0.6 than  $R_2'$ , implying  $Q_1^c$  is more likely matched to the overlapped span than  $Q_2^c$ ,  $R_2'$  is removed from  $R_2$  (L4); otherwise,  $R_1'$  is removed from  $R_1$  (L5).

Finally, Algorithm 3 shows how session spans are found for a specific section.  $C^e$  is a list of tuples

<sup>6</sup>Core questions are required for retrieving essential information, while optional questions depend on the answers to the core questions, so are often skipped during the interview.

---

**Algorithm 2:** `remove_overlap( $R_1, R_2$ )`

---

**Input:**  $R_1, R_2$ : ordered lists of tuples comprising utterance IDs and their matching scores for the first and second sections, respectively.

**Output:** ( $R_1, R_2$ ): updated lists without overlaps.

```
1  $R'_1 \leftarrow [(i, s) : \forall (i, s) \in R_1 \wedge (i, *) \in R_2]$ ;
2  $R'_2 \leftarrow [(i, s) : \forall (i, s) \in R_2 \wedge (i, *) \in R_1]$ ;
3 if  $|\text{select}(R'_1, 0.6)| > |\text{select}(R'_2, 0.6)|$  then
4   return ( $R_1, R_2 \setminus R'_2$ )
5 return ( $R_1 \setminus R'_1, R_2$ )
```

---

comprising utterance IDs and their scores for the  $k$ 'th section created by Algorithms 1 and 2 (L1) (`ro`: `remove_overlap`).  $C^\ell$  is created in the same manner, except adapting the Levenshtein Distance (LD) as the similarity metric (L2) (Levenshtein, 1966).  $\text{sel}(C, s)$  returns a list of tuples comprising utterance IDs and their matched question IDs, where the scores  $> s$ .  $\text{last}(U, Q_*)$  returns the first utterance ID of the  $(k + 1)$ 'th section if exists; otherwise, it returns the last utterance ID of  $U$ .  $C$  is created by taking the intersection of  $C^e$  and  $C^\ell$  whose scores  $> 0.8$  and  $0.7$ , and the last utterance ID (L3).<sup>7</sup>

For each span  $U'$  of utterances between  $C_i$  and  $C_{i+1}$  (exclusive for both ends), a list  $Q'$  of optional questions related to  $C_i$  is created (L5-7).  $T^e$  is a list of tuples comprising utterance IDs in  $U'$  and their matched question IDs in  $Q'$  with scores  $> 0.8$ , and  $T^\ell$  is created using LD (L8-9). The intersection of  $T^e$  and  $T^\ell$  is appended to a list  $O$  (L10), which is then merged with  $C$  and sorted to produce  $V$  (L11).

For each span  $U''$  between  $V_i$  and  $V_{i+1}$ , a list  $Q''$  of any questions have not been matched in that span is created (L14). Bipartite matching bw.  $U''$  and  $Q''$  are performed to find matches optimizing several criteria in Appendix B.1 (L15), accumulated, merged, and sorted to produce the final list (L16-17).

#### 4.4 Assessment Pairing

Answers to the questions are used to determine the values of the variables (Table 1), resulting in many-to-many relations between questions and variables (many-questions to one-variable is the most common case). Our data comprises five variable types. (1) *Scale* assesses on an ordinal scale with ratings for intensity, severity, or likeness. (2) *Category* selects among binary choices or distinct class labels. (3) *Measure* captures various units such as duration, frequencies, and ages. (4) *Notes* are summarized

<sup>7</sup>Any section not matched by Algorithm 1 is considered absent.

---

**Algorithm 3:** `session_match( $U, Q_{1..4}^c, Q_{1..4}^o, k$ )`

---

**Input:**  $U$ : a list of utterances,  $Q_{1..4}^{c/o}$ : lists of core|optional questions for the 1..4'th sections,  $k$ : the section index to segment sessions in.

**Output:** ( $R_1, R_2$ ): updated lists without overlaps.

```
1  $C^e \leftarrow \text{ro}(\text{sm}^e(U, Q_k^c), \text{sm}^e(U, Q_{\forall j \neq k}^c))$ ;
2  $C^\ell \leftarrow \text{ro}(\text{sm}^\ell(U, Q_k^c), \text{sm}^\ell(U, Q_{\forall j \neq k}^c))$ ;
3  $C \leftarrow (\text{sel}(C^e, 0.8) \cap \text{sel}(C^\ell, 0.7)) \cup \text{last}(U, Q_k^c)$ ;
4  $O \leftarrow \emptyset$ ;
5 for  $i \leftarrow 1$  to  $(|C| - 1)$  do
6    $U' \leftarrow$  a list of utterances between  $C_i$  and  $C_{i+1}$ ;
7    $Q' \leftarrow$  a list of questions in  $Q_k^c$  related to  $C_i$ ;
8    $T^e = \text{sel}(\text{sm}^e(U', Q'), 0.8)$ ;
9    $T^\ell = \text{sel}(\text{sm}^\ell(U', Q'), 0.7)$ ;
10   $O \leftarrow O \cup (T^e \cap T^\ell)$ ;
11  $(V, W) = (\text{sorted}(C \cup O), \emptyset)$ ;
12 for  $i \leftarrow 1$  to  $(|V| - 1)$  do
13    $U'' \leftarrow$  a list of utterances between  $V_i$  and  $V_{i+1}$ ;
14    $Q'' \leftarrow$  a list of questions in  $Q_k^c \cup Q_k^o$  that are
      between  $V_i$  and  $V_{i+1}$ ;
15    $T \leftarrow$  the best bipartite matching results between
       $U''$  and  $Q''$  optimizing several criteria in B.1;
16    $W \leftarrow W \cup T$ 
17 return  $\text{sorted}(V \cup W)$ 
```

---

texts documented by the interviewers. (5) *Rule* is calculated based on predefined rules derived from the other variable types. Table 3 shows the statistics of all variables for each section in our dataset.

Type	Variables					Count
	LBI	THH	CRA	CAP	Total	
Scale	7	1	0	40	48	9,722
Category	4	9	15	3	31	4,258
Measure	2	0	1	24	27	3,482
Notes	1	10	3	0	14	1,146
Rule	1	0	1	25	27	6,326

Table 3: Statistics of the five types of variables. Examples of these variables are provided in Appendix B.2.

## 5 Experiments

### 5.1 Dataset

The original data contains 411 interviews (Sec. 3). Whisper tends to generate irrelevant or repetitive sequences when prolonged silences occur, rendering about  $\approx 20\%$  of the resulting transcripts unusable. To address this issue, silence removal and noise cancellation techniques are applied, recovering  $\approx 80\%$  of them. Among the 393 successful transcripts, 322

VT	Template
<b>S&amp;C</b>	[INTRO]. Based on the patient’s interview history, please determine $\{keywords\}$ that the patient $\{symptom\}$ . [RETURN]. [REASON]. The "answer" should be in the range $\{range\}$ . $\{attributes\}$
<b>M</b>	[INTRO]. Based on the patient’s interview history, please calculate $\{keywords\}$ that the patient have $\{symptom\}$ . [RETURN]. [REASON]. The "answer" should be $\{type\}$ .
<b>N</b>	[INTRO]. Based on the formatted data from patient’s interview, please determine whether or not the formatted data includes this specified information $\{single\_slot\}$ . [RETURN]. The "reason" gives a brief explanation on whether the formatted data includes or omits the information. The "answer" should be either "yes" or "no", indicating the presence or absence of the information in formatted data.

Table 4: Instruction templates for **Scale**, **Category**, **Measure**, and **Notes** variables. VT: Variable type, [INTRO]: Imagine you are a professional clinician, [RETURN]: Return the answer as a JSON object with "reason" and "answer" as the keys, [REASON]: The "reason" should provide a brief justification or explanation for the answer.

of them have human assessments (§4.4), which are used to evaluate our approach (Table 5).

	Audios	Hours	Turns	Tokens
Original	411	703	116,501	6,035,027
Transcribe	393	651	90,174	5,499,662
Evaluation	322	515	71,412	4,335,977

Table 5: Statistics of our PTSD interview dataset.

Compared to other interview datasets<sup>8</sup>, our dataset is the largest in the mental health domain. While existing datasets often involve human-machine dialogues or crowdworker simulations, ours consists of formal diagnostic interviews conducted entirely by clinicians, making it the first clinician-administered interview dataset. Additionally, our dataset aims to generate comprehensive diagnostic reports rather than just single scores, providing more detailed resource for clinical practice.

## 5.2 Large Language Models (LLMs)

The state-of-the-art commercial and open-source large language models, GPT-4 and Llama-2 (Touvron et al., 2023), are adapted for our experiments.<sup>9</sup> For each question, a model takes all sessions related to the variable to which the question pertains (§4.4), and an instruction to provide the answer and explanation. Table 4 shows our templates including replaceable patterns to generate the instruction for each variable type. For **Scale**,  $\{keywords\}$  can be replaced with "how severe", and  $\{symptom\}$  with "have unwanted dreams in the past month". For **Category**,  $\{keywords\}$  can be replaced with "which of the following categories best describes", and  $\{symptom\}$  with "usual employment status". To constrain the answer generated by the model, details such as the answer  $\{range\}$  for **S&C**, and the

value  $\{type\}$  for **Measure** are incorporated. **S** has a special pattern  $\{attributes\}$ , directing the model to return a particular score under certain conditions.

Assessing model performance for **Notes** poses a challenge as they must be compared against text summarized by interviewers. Given the complexity of this task, it is decomposed into multiple subtasks of binary classifications, information extraction, and categorization by adapting Chain-of-Thought (Wei et al., 2023). First, GPT is asked to generate a list of slots for each **N** variable, based on a batch of summary notes from interviewers. Because many of these slots have similar meanings, albeit varying in naming, GPT is again asked to cluster them. The clusters generated by GPT are manually refined, resulting in final grouped slots that cover 95+% of the initial generation. For each of these slots, an LLM is tasked with determining if relevant content for the slot is present in the provided sessions.<sup>10</sup>

## 5.3 Zero-shot V.S. Few-shot Settings

Zero-shot and few-shot settings are tested across all variable types<sup>11</sup>. For **Scale**, two few-shot settings are explored: one including an example for a single scale point, and the other covering examples for all scale points. For the GPT model, few-shot settings mostly outperform zero-shot settings in predicting **Category**, **Measure**, and **Notes** variables. For **Scale**, the few-shot setting with a single example results in the lowest performance. On the other hand, the few-shot setting including examples for all scale points shows a slight improvement in model performance. Thus, few-shot settings are used for all experiments with GPT. In contrast, the Llama model consistently yields inferior outcomes with few-shot settings compared to zero-shot settings, leading us to adopt zero-shot settings for all Llama experiments.

<sup>8</sup>Statistics of the comparison is provided in Appendix C.1.

<sup>9</sup>Specific versions, parameters, and costs for these large language models are provided in Appendix C.3 and C.4.

<sup>10</sup>Appendix C.5 gives slot examples for **Notes** variables.

<sup>11</sup>Appendix C.2 gives details on zero/few-shot settings.

Type	Count	Accuracy		RMSE		Bias		Recall	
		GPT-4	Llama-2	GPT-4	Llama-2	GPT-4	Llama-2	GPT-4	Llama-2
Scale	9,722	58.9	46.7	1.10	1.63	-0.04	0.51	-	-
Scale <sub>g</sub>	9,722	67.3	59.0	0.85	1.01	-0.04	0.51	-	-
Category	4,258	77.2	63.6	-	-	-	-	-	-
Measure	3,482	64.4	56.5	-	-	-0.34	-0.004	-	-
Notes	1,146	-	-	-	-	-	-	48.1	52.7
Rule	6,326	68.4	59.8	0.80	0.92	-0.15	0.44	-	-

Table 6: Model performance on all variable types (§4.4) using four evaluation metrics (§5.4).

## 5.4 Evaluation Metrics

Since each variable type is uniquely defined, different evaluation metrics are employed accordingly. Accuracy is computed for all types except *Notes*. For *Notes*, since the model identifies the presence of information in the provided sessions based on predefined slots, Recall is used as the primary metric to gauge the coverage of relevant information detected by the model. For *Scale*, the Root Mean Square Error (RMSE) and Bias evaluation are used. RMSE quantifies the magnitude of errors, whereas Bias evaluation calculates the proportion of positive and negative residuals, thereby revealing any directional bias in the model predictions.

## 5.5 Results

Table 6 gives the results for each variable type. For *Scale*, additional evaluation is conducted for CAP whose original scaling ranges from 0 to 4 where 0 indicates the absence of symptoms, 1 denotes minimal symptoms, and 2+ are considered symptoms that meet or exceed the threshold for clinical significance. To reflect this clinical demarcation, scale points are categorized into three scale groups, 0, 1, and 2+, and evaluated as *Scale<sub>g</sub>*.<sup>12</sup>

GPT consistently shows significantly higher accuracy, averaging 10.5% more across all types than Llama, and reaches an accuracy of 68.4% for *Rule* accumulating outcomes of other types. Regarding RMSE, GPT exhibits an error rate of 0.8 for *Rule* using results of *Scale*, implying that it is less than one scale off from human judgment on average. In terms of Bias, ranging from -1: *completely biased to negative* to 1: *completely biased to positive*, GPT displays a marginal bias toward negative for *Scale*, while Llama shows a strong positive bias, implying that GPT is a bit conservative in predicting a higher scale, whereas Llama tends to overestimate. GPT underestimates more than Llama for *Measure*, how-

ever, showing a slight negative bias of 0.15 for *Rule*. For *Notes*, Llama exhibits better performance with a recall of 52.7% than GPT, suggesting that Llama is more effective in retrieving relevant information. Considering that these models are not fine-tuned on our data, this level of performance is very promising, as we can achieve a robust model for practical use with further training.

## 6 Error Analysis

A thorough error analysis is conducted by proportionally sampling 100+ examples per variable type. Six types of major errors are identified (Table 7), with only two attributed to LLMs and the remainder caused by external factors, implying that the true LLM performance may be even higher.

**Misaligned Reasoning** One predominant error type occurs when models deviate from instructions of the rating scheme, presenting seemingly logical reasoning, although it ultimately leads to incorrect conclusions. In Table 7, both models fail to align the key term provided by the participant, *extremely*, with the definition of score 4 - “Extreme, dramatic physical reactivity”. Llama tends to deviate further than GPT, resulting in a higher RMSE.

**False Negatives** is a major error type caused by:

1. Inaccurate assessments by clinicians. In Table 7, the participant reports *five times a week*, yet the clinician incorrectly records the frequency of monthly basis as 5, which should have been 20 times a month.
2. Ambiguity in *Scale* where answers may fall between two scales, resulting in potentially valid model predictions being marked incorrect.
3. The model’s inability to recognize paraphrased information in *Notes*, mistakenly indicating the absence of slot information. This issue particularly affects GPT’s performance due to its strict interpretation of wording variations.

<sup>12</sup>Appendix C.6/C.7 presents results for each section/variable.

Type	History	Gold	Auto		Ext
			GPT	LM	
MR	Have you had any physical reactions when something reminded you of what happened? ... I had a horrible headache. ... How many times in the past month has that happened? ... Those two times. ... How long did it take you to sort of feel back to normal? I swear. It took me a minute. I got up. I got a glass of water. It took me about. I say two to three hours. ... So how bad was that Headache? Do you think there are any other symptoms? It was extremely. I never had. I had it like that.	4	3	2	
FN	... can you think about like how often that might happen in the last month about? I feel like about like five times a week.	5	20	20	✓
EI	... when did those start for you? ... So, since around age 12, at least yeah yeah because it took me a long time to really trust my stepfather.	480	NA	108	✓
TE	... how satisfied and fulfilled have you felt about your life, with zero being like not at all, couldn't have a worse life, and 10 being perfect, couldn't have a better life? I would say a C, because it's a lot more things that I want to do to be at a 10.	2	3	3	✓
SM	So how many times in the past month would you say some things made you upset that reminded you of it? Rarely, maybe like two, three times? Very rarely.	2	1	1	✓
CR	... thinking about your work in the past month, how have you been doing? ... It's a normal, consistent, um, it's a normal, consistent routine where I do the same thing, do the same thing every day.	40	NA	40	

Table 7: Examples of the six error types. MR: Misaligned Reasoning, FN: False Negative, EI: External Information, TE: Transcription Error, SM: Session Mismatching, CR: Commonsense Reasoning. Gold: clinician’s answers, Auto: model-predicted answers. Ext: errors caused by external factors, not LLMs. NA: the model predicts None. Clinician’s questions are highlighted in blue. Patient’s key information to the questions are highlighted in red.

**External Information** One common issue is the absence of external information, such as the prior knowledge about the patient (e.g., medical histories, demographics) or the content of previous interview questions. In Table 7, although both models see the onset of symptoms at age 12, they fail to provide an accurate response of the total symptom duration in months because the patient’s current age (that is 52) is not provided in the transcript. In this case, GPT tends to generate a None answer, while Llama tends to hallucinate the patient’s age, and thus produces an answer based on an arbitrary assumption.

**Transcription Error** Transcription errors from automatic speech recognizers often cause LLMs to incorrectly interpret the answers, especially with short responses (e.g., *yes*, *no*, single digits like 6), medical terminologies, or non-verbal cues such as nodding. In Table 7, the number ‘6’ is incorrectly transcribed as ‘C’ in the participant’s response.

**Session Mismatching** A question can be mismatched with the transcript, especially when the clinician extensively paraphrases it. In such cases, the segmented session may or may not contain all the necessary information to answer the question. In Table 7, both models correctly answer based on the patient’s response (1: Minimal). However, due to the mismatch, the session is missing a part where the patient also indicates 2 (clearly present but still manageable), which is recorded as gold.

**Commonsense Reasoning** The models’ limitations extend to inferring basic human experiences. Unable to deduce standard working hours from a *normal, consistent routine* in Table 7, the models fall short of clinician-like assumptions of a typical 40-hour workweek, showcasing a gap in applying commonsense logic to the assessment.

## 7 Conclusion

In this study, we undertook the task of automating PTSD diagnostics using 411 clinician-administered interviews. To ensure the data quality, we develop an end-to-end pipeline streamlining transcription, alignment, segmentation, and assessment pairing. We also construct a pioneering framework for this task by leveraging two state-of-the-art LLMs. Our findings reveal the substantial potential of LLMs in assisting clinicians with diagnostic validation and decision-making processes. Our error analysis suggests future directions for improvement, such as incorporating external information or commonsense knowledge to engineer more comprehensive instructions. We envision that this framework holds promise for addressing a broader spectrum of mental health conditions and offers novel insights into LLM applications within the mental health domain. We plan to collect more data and train a custom LLM to better preserve patients’ privacy, and develop a dialogue system to conduct the interviews.

## Limitations

Although the experiment results prove the capability of LLMs to automate PTSD diagnosis, their applications in real-world unsupervised clinical settings are premature. To avoid the possible negative influence of model errors on the patients, we recommend using this framework as a supportive tool for clinicians in diagnostics and decision-making.

It should be noted that the clinician annotated gold assessment data is not perfect, which may affect evaluation accuracy. However, this framework makes it easier to identify and refine the inaccuracies in the gold assessment data and thus improve its overall validity. We leave the data augmentation as the next step of our future work.

In addition, the experiments in this paper utilize LLMs without fine-tuning. One limitation is that we have little control over the model predictions. The models, especially Llama-2, generate unexpected outputs that violate the instructions. Furthermore, data privacy concerns restrict the use of models like GPT for clinical data. To address these issues and enhance framework adaptability, future work will focus on developing more controllable, open-sourced models that guarantee data protection in line with clinical domain restrictions.

Due to strict Institutional Review Board (IRB) regulations concerning the confidentiality of real patient information, we are unable to release the dataset, even in an anonymized format. However, recognizing the importance of contributing to the research community, we are pleased to announce that we will release the framework utilized in our study. This, we believe, will facilitate further research and innovation, as our methodology is versatile and can be adapted to a wide array of mental health conditions, provided the requisite interview question sets and video/transcripts are available.

## Ethical Considerations

The diagnostic interview data used in this paper was collected with informed consent approved by the Institutional Review Board (IRB) and Research Oversight Committee. The authors and clinicians involved in the research have passed Research, Ethics, Compliance, and Safety Training through Collaborative Institutional Training Initiative<sup>13</sup> (CITI Program). For the use of LLMs, this study exclusively employs anonymized interviews,

<sup>13</sup><https://about.citiprogram.org>

ensuring the confidentiality and privacy of all participants. All practices in this research adhere to the ACL Code of Ethics.

## Acknowledgements

We gratefully acknowledge the support of the DooGood Foundation. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the DooGood Foundation.

## References

- Alon Bartal, Kathleen Jagodnik, Sabrina Chan, and Sharon Dekel. 2023. [Chatgpt Demonstrates Potential for Identifying Psychiatric Disorders: Application to Childbirth-Related Post-Traumatic Stress Disorder](#).
- Runa Bhaumik, Vineet Srivastava, Arash Jalali, Shanta Ghosh, and Ranganathan Chandrasekaran. 2023. Mindwatch: A smart cloud-based ai solution for suicide ideation detection leveraging large language models. *medRxiv*, pages 2023–09.
- Marco Cascella, Jonathan Montomoli, Valentina Bellini, and Elena Bignami. 2023. [Evaluating the Feasibility of Chatgpt in Healthcare: An Analysis of Multiple Clinical and Research Scenarios](#). *Journal of Medical Systems*, 47(1).
- 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 Sankaranarayanan Pilla, 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](#).
- Zachary Enghardt, Chengqian Ma, Margaret E. Morris, Xuhai "Orson" Xu, Chun-Cheng Chang, Lianhui Qin, Daniel McDuff, Xin Liu, Shwetak Patel, and Vikram Iyer. 2023. [From Classification to Clinical Insights: Towards Analyzing and Reasoning About Mobile and Behavioral Health Data With Large Language Models](#).

- Jonathan G. Fiscus, Jerome Ajot, Martial Michel, and John S. Garofolo. 2006. The Rich Transcription 2006 Spring Meeting Recognition Evaluation. In *Proceedings of International Workshop on Machine Learning and Multimodal Interaction*, pages 309–322.
- Guanghui Fu, Qing Zhao, Jianqiang Li, Dan Luo, Changwei Song, Wei Zhai, Shuo Liu, Fan Wang, Yan Wang, Lijuan Cheng, Juan Zhang, and Bing Xiang Yang. 2023. [Enhancing Psychological Counseling with Large Language Model: A Multifaceted Decision-Support System for Non-Professionals.](#)
- Isaac R. Galatzer-Levy, Daniel McDuff, Vivek Nataraajan, Alan Karthikesalingam, and Matteo Malgaroli. 2023. [The Capability of Large Language Models to Measure Psychiatric Functioning.](#)
- Rachel L. Gluck, Georgina E. Hartzell, Hayley D. Dixon, Vasiliki Michopoulos, Abigail Powers, Jennifer S. Stevens, Negar Fani, Sierra Carter, Ann C. Schwartz, Tanja Jovanovic, Kerry J. Ressler, Bekh Bradley, and Charles F. Gillespie. 2021. [Trauma exposure and stress-related disorders in a large, urban, predominantly african-american, female sample.](#) *Archives of Women’s Mental Health*, 24(6):893–901.
- Chen Gong, Peilin Wu, and Jinho D. Choi. 2023. [Aligning Speakers: Evaluating and Visualizing Text-based Speaker Diarization Using Efficient Multiple Sequence Alignment.](#) In *Proceedings of the 35th IEEE International Conference on Tools with Artificial Intelligence*, ICTAI’23.
- Jonathan Gratch, Ron Artstein, Gale Lucas, Giota Strattou, Stefan Scherer, Angela Nazarian, Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, David Traum, Skip Rizzo, and Louis-Philippe Morency. 2014. [The distress analysis interview corpus of human and computer interviews.](#) In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 3123–3128, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Yan Hu, Qingyu Chen, Jingcheng Du, Xueqing Peng, Vipina Kuttichi Keloth, Xu Zuo, Yujia Zhou, Zehan Li, Xiaoqian Jiang, Zhiyong Lu, Kirk Roberts, and Hua Xu. 2023. [Improving Large Language Models for Clinical Named Entity Recognition via Prompt Engineering.](#)
- Yining Hua, Fenglin Liu, Kailai Yang, Zehan Li, Yi-han Sheu, Peilin Zhou, Lauren V. Moran, Sophia Ananiadou, and Andrew Beam. 2024. [Large Language Models in Mental Health Care: a Scoping Review.](#)
- Shaoxiong Ji, Tianlin Zhang, Kailai Yang, Sophia Ananiadou, and Erik Cambria. 2023. [Rethinking Large Language Models in Mental Health Applications.](#)
- Martin B. Keller, Philip W. Lavori, Barbara Friedman, Eileen Nielsen, Jean Endicott, Pat McDonald-Scott, and Nancy C. Andreasen. 1987. [The Longitudinal Interval Follow-up Evaluation. A comprehensive method for assessing outcome in prospective longitudinal studies.](#) *Archives Of General Psychiatry*, 44(6):540–548.
- Dietrich Klakow and Jochen Peters. 2002. [Testing the Correlation of Word Error Rate and Perplexity.](#) *Speech Communication*, 38(1):19–28.
- Tin Lai, Yukun Shi, Zicong Du, Jiajie Wu, Ken Fu, Yichao Dou, and Ziqi Wang. 2023. [Psy-Llm: Scaling up Global Mental Health Psychological Services with Ai-based Large Language Models.](#)
- Bishal Lamichhane. 2023. [Evaluation of Chatgpt for Nlp-based Mental Health Applications.](#)
- Vladimir I Levenshtein. 1966. [Binary Codes Capable of Correcting Deletions, Insertions, and Reversals.](#) *Soviet Physics Doklady*, 10(8):707–710.
- Jialin Liu, Changyu Wang, and Siru Liu. 2023a. [Utility of Chatgpt in Clinical Practice.](#) *Journal of Medical Internet Research*, 25:e48568.
- June M. Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi Liao, and Jiamin Wu. 2023b. [Chatcounselor: A Large Language Models for Mental Health Support.](#)
- Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. [Towards emotional support dialog systems.](#)
- Siyuan Brandon Loh and Aravind Sesagiri Raamkumar. 2023. [Harnessing Large Language Models’ Empathetic Response Generation Capabilities for Online Mental Health Counselling Support.](#) *arXiv preprint arXiv:2310.08017*.
- Chong Ma, Zihao Wu, Jiaqi Wang, Shaochen Xu, Yaonai Wei, Zhengliang Liu, Xi Jiang, Lei Guo, Xiaoyan Cai, Shu Zhang, Tuo Zhang, Dajiang Zhu, Dinggang Shen, Tianming Liu, and Xiang Li. 2023a. [Impressiongpt: An Iterative Optimizing Framework for Radiology Report Summarization with Chatgpt.](#)
- Zilin Ma, Yiyang Mei, and Zhaoyuan Su. 2023b. [Understanding the benefits and challenges of using large language model-based conversational agents for mental well-being support.](#) In *AMIA Annual Symposium Proceedings*, volume 2023, page 1105. American Medical Informatics Association.
- Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. 2023. [Capabilities of Gpt-4 on Medical Challenge Problems.](#)
- OpenAI. 2023. [Gpt-4 Technical Report.](#)
- Cheng Peng, Xi Yang, Kaleb E Smith, Zehao Yu, Aokun Chen, Jiang Bian, and Yonghui Wu. 2023. [Model Tuning or Prompt Tuning? A Study of Large Language Models for Clinical Concept and Relation Extraction.](#)

- Wei Qin, Zetong Chen, Lei Wang, Yunshi Lan, Weijieying Ren, and Richang Hong. 2023. [Read, diagnose and chat: Towards explainable and interactive llms-augmented depression detection in social media.](#)
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. [Robust Speech Recognition via Large-scale Weak Supervision.](#) In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*, pages 28492–28518.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-Bert: Sentence Embeddings using Siamese Bert-Networks.](#)
- David V Sheehan, Yves Lecrubier, K Harnett Sheehan, Patricia Amorim, Juris Janavs, Emmanuelle Weiller, Thierry Hergueta, Roxy Baker, Geoffrey C Dunbar, et al. 1998. The mini-international neuropsychiatric interview (mini): the development and validation of a structured diagnostic psychiatric interview for dsm-iv and icd-10. *Journal of clinical psychiatry*, 59(20):22–33.
- Ying Shen, Huiyu Yang, and Lin Lin. 2022. [Automatic depression detection: An emotional audio-textual corpus and a gru/bilstm-based model.](#)
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal 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.](#)
- Michel Valstar, Björn Schuller, Kirsty Smith, Timur Al-maev, Florian Eyben, Jarek Krajewski, Roddy Cowie, and Maja Pantic. 2014. [Avec 2014: 3d dimensional affect and depression recognition challenge.](#) In *Proceedings of the 4th International Workshop on Audio/Visual Emotion Challenge, AVEC '14*, page 3–10, New York, NY, USA. Association for Computing Machinery.
- Frank W. Weathers, Michelle J. Bovin, Daniel J. Lee, Denise M. Sloan, Paula P. Schnurr, Danny G. Kaloupek, Terence M Keane, and Brian P. Marx. 2018. [The Clinician-Administered PTSD Scale for DSM-5 \(CAPS-5\): Development and initial psychometric evaluation in military veterans.](#) *Psychological Assessment*, 30(3):383–395.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.](#)
- World Health Organization. 2021. *Mental health atlas 2020.* World Health Organization.
- World Medical Association. 2013. [World Medical Association Declaration of Helsinki: Ethical Principles for Medical Research Involving Human Subjects.](#) *The Journal of the American Medical Association*, 310(20):2191–2194.
- Yuqi Wu, Jie Chen, Kaining Mao, and Yanbo Zhang. 2023. [Automatic Post-Traumatic Stress Disorder Diagnosis via Clinical Transcripts: A Novel Text Augmentation with Large Language Models.](#) In *2023 IEEE Biomedical Circuits and Systems Conference (BioCAS)*, pages 1–5.
- Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia Gabriel, Hong Yu, James Hendler, Marzyeh Ghassemi, Anind K. Dey, and Dakuo Wang. 2023. [Mental-Llm: Leveraging Large Language Models for Mental Health Prediction via Online Text Data.](#)
- Kailai Yang, Tianlin Zhang, Ziyang Kuang, Qianqian Xie, Sophia Ananiadou, and Jimin Huang. 2023. [Mental-lama: Interpretable Mental Health Analysis on Social Media with Large Language Models.](#)
- Binwei Yao, Chao Shi, Likai Zou, Lingfeng Dai, Mengyue Wu, Lu Chen, Zhen Wang, and Kai Yu. 2022. [D4: a chinese dialogue dataset for depression-diagnosis-oriented chat.](#)
- Tianlin Zhang, Annika M. Schoene, Shaoxiong Ji, and Sophia Ananiadou. 2022. [Natural language processing applied to mental illness detection: a narrative review.](#) *npj Digital Medicine*, 5(1).

## A Section Details

Table 8 - 11 give examples for 4 core sections. Each example includes the standard interview Question, the Variable that the question belongs to, and the example Sessions between the Clinician and the Participant.

The Mini International Neuropsychiatric Interview (MINI) is a brief, structured diagnostic interview for diagnosing 17 major psychiatric disorders (Sheehan et al., 1998). We adopt 6 modules from MINI to assess conditions such as Major Depressive Episode (MDE), Mania & Hypomania (MH), PTSD (past incidents), Psychosis Symptoms (PS), Substance Use Disorder (SUD), and Alcohol Use Disorder (AUD). Table 12 provides an example from the MDE module.

<b>Q</b>	What has been your primary source of income over the past month?
<b>V</b>	lbi_a1
<b>S</b>	<b>C:</b> You got to do it all over again. Are you working full time? <b>P:</b> Yes.
<b>Q</b>	How would you rate your overall satisfaction on a scale of 1 to 10, with 1 being the best and 10 being the worst?
<b>V</b>	lbi_e1
<b>S</b>	<b>C:</b> In the past month, like how satisfied have you felt with your life? If we were doing like a scale of one to 10, one is like, it's the worst. This is the worst I've ever had in my life. 10 being like, this is, I'm living my best life. Living my life like it's golden. <b>P:</b> I actually feel like that now. I actually do. Cause until January 1st of this year, I had been unemployed the last two years.

Table 8: Two examples of the LBI section.

<b>Q</b>	Do you have any current physical health conditions?
<b>V</b>	thh_medicalcond
<b>S</b>	<b>C:</b> OK, so now we're going to move on to talking about your health and treatment history. Do you currently have, do you have any current physical health conditions? Did you say no? OK, I couldn't hear what you were saying. Go ahead. <b>P:</b> I have a skin condition called eczema.
<b>Q</b>	In the past, have you been treated for any emotional/mental health problems with therapy or hospitalization?
<b>V</b>	thh_tx_yesno
<b>S</b>	<b>C:</b> In the past, have you been treated for any emotional or mental health problem with therapy or hospitalization? <b>P:</b> No. Yes.

Table 9: Two examples of the THH section.

<b>Q</b>	Tell me a little bit more about what happened.
<b>V</b>	trauma1whathappened
<b>S</b>	<b>C:</b> OK, and what would that be? <b>P:</b> My mom worked at the airport here in xxx. It was the food catering place. They put the food, made the food for the planes. When I was a child, every year, they would sponsor a day at xxx. They would go out there and barbecue. We took over the whole picnic area. You had free entrance to the park, plus tickets to do all the little fair games and all that good stuff. Having a good time. My mom asked my stepfather to go with us because he had a car. He said he didn't wanna go and he wasn't going nowhere. So my mom put us all on the bus. We drove the bus out there. When we came home, it was like 11 o'clock. Of course, we living in xxx. You know that bus ride was long. It was dark, dark when we got home and she had all three of her children with her. My mom unlocked the door, closed that door, the house was pitch black. That man shot down them steps at my mama and all three of her children five times.

Table 10: An example of the CRA section.

<b>Q</b>	Tell me a little bit more about what happened.
<b>V</b>	dsm5capscritb_01trauma1_distress
<b>S</b>	<b>C:</b> To this day, let's say over the past month. So since like the beginning of April, end of March, have you had unwanted memories of this event? Does it randomly pop into your mind at all? Like while you're awake? <b>P:</b> Well, actually my daughter's in an abusive relationship. So yes, I do think about it a lot. Every time I see her, all I think about is my mom. How she endured it.
<b>Q</b>	How often in the past month?
<b>V</b>	dsm5capscritc02trauma1_num
<b>S</b>	<b>C:</b> So in the last month, thinking about the things that you have tried to avoid, how often would you say you've done that? <b>P:</b> I guess every day. I don't know. I just, the most I've done is just, and me avoiding stuff is me just sitting here smoking and playing my video game. That avoids me from thinking about anything negative in my life. And I just try to avoid that.

Table 11: Two examples of the CAP section.

<b>Q</b>	For the past two weeks, were you depressed or down, or felt sad, empty or hopeless most of the day, nearly every day?
<b>V</b>	miniv7_mde_c_a1
<b>S</b>	<b>C:</b> I'm going to ask you some different questions. We're going to focus on the past two weeks right now. So for the past two weeks, did you feel depressed, down, sad, empty or hopeless for most of the day, almost every day the past two weeks? <b>P:</b> Um, no.

Table 12: An example of the MINI section.

## B Data Preprocessing Details

### B.1 Final Matching Criteria

The best bipartite result should follow the criteria:

- All matching IDs need to be ascending.
- Only edges whose embedding cosine similarity  $> 0.4$  are kept.
- Maximize:  $y = \sum_{i=1}^n a_i \cdot x_i$ , subject to  $x_i \geq 0$ , for  $i = 1, \dots, n$ .

In our case, let  $n = 9$ , with the following variables:

- $x_1$ : the sum of Sentence Transformer (ST) cosine similarity scores of all edges
- $x_2$ : the sum of Levenshtein Distance (LD) similarity scores of all edges
- $x_3$ : the average ST cosine similarity scores of all matched questions
- $x_4$ : the average LD similarity scores of all matched questions
- $x_5$ : the total number of matched core questions
- $x_6$ : the total number of matched questions that take the maximum ST cosine similarity result
- $x_7$ : the total number of matched questions that take the maximum LD similarity result
- $x_8$ : the total number of matched core questions that take the maximum ST cosine similarity result
- $x_9$ : the total number of matched core questions that take the maximum LD similarity result

And the coefficients are set as:

- $a_1 = a_2 = 1$
- $a_3 = a_4 = 1$
- $a_5 = a_6 = a_7 = 0.1$
- $a_8 = a_9 = 0.2$

### B.2 Variable Examples

Table 13 - 17 show examples for each variable type. Every example includes the Variable name, replaceable Patterns for prompt generation (Section 5), answer Range, and covered Questions. Note that *Measure*, *Notes*, and *Rule* variables do not have a predefined range. And *Rule* variables are calculated from the results of their Related Variables.

<b>V</b>	dsm5capscribtb01trauma1_distress
<b>P</b>	{ <i>keywords</i> }: how intense in the past month { <i>symptom</i> }: unwanted memories of the traumatic event while awake { <i>attributes</i> }: - If the symptom only exists in dreams, the answer should be 0. - If the symptom is not perceived as involuntary and intrusive, the answer should be 0.
<b>R</b>	0: None, 1: Minimal, minimal distress or disruption of activities 2: Clearly Present, distress clearly presented but still manageable, some disruption of activities 3: Pronounced, considerable distress, difficulty dismissing memories, marked disruption of activities 4: Extreme, incapacitating distress, cannot dismiss memories, unable to continue activities
<b>Q</b>	In the past month, have you had any unwanted memories of it while you were awake, so not counting dreams? - How does it happen that you start remembering it? - Are these unwanted memories, or are you thinking about it on purpose? - How much do these memories bother you? - Are you about to put them out of your mind and think about something else? - Overall, how much of a problem is this for you? - How so?

Table 13: An example of the *Scale* variable. Questions start with - are optional questions that might be skipped based on the participant's response.

<b>V</b>	lbi_a1
<b>P</b>	{ <i>keywords</i> }: which of the following categories best describes, { <i>symptom</i> }: usual employment status
<b>R</b>	1: Full-Time Gainful Employment 2: Part-Time Gainful Employment (30 hours or less/week) 3: Unemployed But Expected by Self or Others 4: Unemployed But Not Expected by Self or Others (e.g., physically disabled) 5: Retired 6: Homemaker 7: Student (Includes Part-Time) 8: Leave of Absence Due to Medical Reasons (e.g., holding job; plans to return) 9: Volunteer Work - Full Time 10: Volunteer Work - Part Time 11: Other 888: N/A
<b>Q</b>	What has been your primary source of income over the past month?

Table 14: An example of the *Category* variable.

<b>V</b>	dsm5capscrib01trauma1_num
<b>P</b>	{ <i>keywords</i> }: how intense in the past month { <i>symptom</i> }: unwanted memories of the traumatic event while awake { <i>type</i> }: an integer representing the frequency of the symptom in the past month
<b>Q</b>	- How often have you had these memories in the past month?

Table 15: An example of the *Measure* variable. The corresponding question for this question is optional, which might be skipped if the participant denies the presence of the symptom.

<b>V</b>	critaprobenotes
<b>P</b>	{ <i>slots</i> }: - trauma_reactions - trauma_details - life_changes - coping_and_changes - worldview_changes - health_concerns - family_and_social_context - nightmare_details - intrusive_experiences - trauma_cognition - trust_and_safety - impact_assessment - age_and_time_factors - substance_use - therapy_and_progress - eating_disorders
<b>Q</b>	You discussed a number of traumas in the last visit with our team members. What would you say is the one that has been most impactful where you are still noticing it affecting you? -* How much do you think about what happened to this day? -* How often do you have nightmares about what happen? -* How much did it change the way you think about yourself and the world? - In the past month, which of these have you thought about more often or had nightmares about or find yourself purposely avoiding thinking about? - Are there any other stressors that you find yourself thinking about when you don't want to or avoiding?

Table 16: An example of the *Notes* variable. Questions start with - are optional questions which might be skipped based on the participant’s response. Questions start with \* are recurrent questions which might be asked multiple times during the interview.

<b>V</b>	dsm5capscrib01trauma1
<b>R</b>	0: Absent 1: Mild/subthreshold 2: Moderate/threshold 3: Severe/markedly elevated 4: Extreme/incapacitating
<b>RV</b>	dsm5capscrib01trauma1_distress dsm5capscrib01trauma1_num

Table 17: An example of the *Rule* variable.

## C Experiments Details

### C.1 Dataset Comparison

Table 18 gives the comparison with related datasets in the mental health domain.

Dataset	A	H	Turns	Utters
DAIC (Gratch et al., 2014)	189	51	-	-
AViD (Valstar et al., 2014)	300	240	-	-
EATD (Shen et al., 2022)	162	2.26	-	-
Psych8k (Liu et al., 2023b)	260	260	-	-
D4 (Yao et al., 2022)	-	-	28,855	81,559
ESConv (Liu et al., 2021)	-	-	-	31,410
<b>Ours</b>	<b>322</b>	<b>515</b>	<b>71,412</b>	<b>142,824</b>

Table 18: Comparisons with existing mental health interview/dialogue datasets in terms of **Audio** counts, total **Hours**, total and utterances.

### C.2 Details on Zero-shot/Few-shot Settings

We randomly sampled 30 instances for each variable type and asked both models to predict under zero-shot and few-shot settings. For the GPT model, few-shot settings generally yield better performance. However, the Llama model consistently fails to follow instructions as the context length grows, leading to significant degradation with few-shot prompting. Additionally, we observed a 28% increase in the likelihood of generating an unexpected response format, such as deviating from the requested JSON format, when using few-shot settings.

### C.3 Experiment Costs

**GPT-4** The pricing of the GPT-4 Turbo model is \$0.01/1K tokens for input and \$0.03/1K tokens for output. We spend approximately \$300 (upper bound) to complete GPT experiments in this paper.

**Llama-2** We use a single NVIDIA H100 GPU for Llama inferences with a batch size of 1, taking

Type	Zero-shot		Few-shot	
	GPT-4	Llama-2	GPT-4	Llama-2
Scale	60.0	50.0	63.3	36.7
Scale <sub>1</sub>	-	-	56.7	40.0
Category	43.3	40.0	46.7	33.3
Measure	56.7	56.7	60.0	50.0
Notes	41.0	42.7	43.6	34.9

Table 19: Model performance on zero-shot and few-shot settings. Scale<sub>1</sub> refers to the few-shot setting that only include one example for a single scale point. Accuracy is the metric used for all variable types except Notes variables, which are evaluated using Recall.

roughly 10 seconds per request. Completing a full set of experiments on all samples requires ~3 days.

#### C.4 LLM Configurations

We utilize gpt-4-1106-preview, the latest GPT-4 Turbo model, and llama-2-70b-chat-hf, the largest Llama-2 model. For GPT, to enhance the stability and consistency of the model output, we configure the *temperature* parameter to 0. This adjustment makes the model’s response more deterministic. Besides, we also employ parameters exclusive to GPT-4 Turbo and GPT-3.5 Turbo, namely *response\_format* and *seed*. Setting *response\_format* to "json\_object" constrains the model to generate parsable JSON strings, facilitating easier data handling and analysis. Despite ChatGPT’s non-deterministic nature, *seed* parameter enables users to obtain consistent outputs across multiple requests, as long as there are no changes at the system level.

As for the Llama, we conduct experiments involving different *temperature*, *top\_p*, and *repetition\_penalty* separately. The results indicate that the model gives better performance with a *temperature* setting of 0.3, a *top\_p* of 0.9, and a *repetition\_penalty* of 1.

#### C.5 Slot Examples for Notes Variable

Table 20 outlines the process for generating, merging, and formatting the slots in *Notes* variables (§5.2). Initially, we compile all clinician-summarized notes for each *Notes* variable and input them into the GPT model using the NSG prompt to produce a list of slots. Due to potential overlaps, the NSM prompt directs the model to consolidate these slots into clusters, ensuring both conciseness and comprehensiveness. Subsequently, the NSF prompt is used to format both the gold-standard summaries and the corresponding interview ses-

sions, facilitating a straightforward comparison of the structured slot arrangements.

#### C.6 Model Performance by Sections

Table 21 presents model performances by each section. Note that THH section lacks *Measure* and *Rule* variables, whereas CRA section does not contain *Scale* variables. The grouped scale<sub>g</sub> is exclusively applied within the CAP section.

#### C.7 Model Performance by Variables

Table 22 lists results for each variable following the four evaluation metrics (Section 5.4).

Step	Template
<b>NSG</b>	<p>As a clinician who has conducted interviews with multiple patients, you are tasked with structuring the interview data into a more organized format. To achieve this, identify general "slots" from the interview question and answers. These slots should represent key themes or types of information that can be adapted to various responses from different patients.</p> <p>For each identified slot, provide a brief explanation of why it has been chosen, focusing on its relevance and utility in categorizing interview data.</p> <p>Your findings should be presented in a JSON format as a list, for example: [{"reason": "This slot captures the primary health concern of the patient, a common theme across all interviews", "slot": "primary_health_concern"}, {"reason": "This slot pertains to the patient's lifestyle habits, which is crucial for understanding health context", "slot": "lifestyle_habits"}].</p> <p>Remember to ensure that the slots are broad enough to be applicable across different patient responses yet specific enough to offer meaningful categorization.</p>
<b>NSM</b>	<p>Imagine you are a clinician who documents patient interviews in a structured, slot-filling manner. Sometimes, certain slots may have overlapping or similar content. Your task is to review a given list of slots and merge those that are similar. The merged results should be returned as a JSON object, where each key represents a merged slot, and the corresponding value is a list of the original slots that have been combined under this merged category.</p> <p>For instance, if the list of slots is: ["daily_routine", "work_events", "daily_activity", "daytime_activities", "work_routine"], a possible merged result could be: {"daily_routine": ["daily_routine", "daily_activity", "daytime_activities"], "work": ["work_events", "work_routine"]}.</p> <p>When you receive a list of slots, analyze and merge them accordingly, ensuring that the merged slots are logically grouped and accurately represent the original information categories.</p>
<b>NSF</b>	<p>Imagine you are a professional clinician. Based on the patient's interview history, please extract specific information and fill in the following slots: {slots}. If the interview history does not provide information for any of these slots, please enter an empty string ("") for that slot. Return the answer as a JSON object.</p>

Table 20: Prompts used for Notes Variable Slot structure Generation, Merging, and Formatting.

	Type	Count	Accuracy		RMSE		Bias		Recall	
			GPT-4	Llama-2	GPT-4	Llama-2	GPT-4	Llama-2	GPT-4	Llama-2
LBI	Scale	1,281	54.6	44.7	1.26	1.42	0.46	0.45	-	-
	Category	594	74.6	67.3	-	-	-	-	-	-
	Measure	99	68.7	66.7	-	-	-0.16	-0.09	-	-
	Notes	203	-	-	-	-	-	-	42.0	50.8
	Rule	215	43.3	37.7	0.94	0.98	0.44	0.43	-	-
THH	Scale	29	55.2	51.7	1.20	1.25	0.23	0.43	-	-
	Category	1,527	92.6	85.9	-	-	-	-	-	-
	Notes	633	-	-	-	-	-	-	52.5	59.8
CRA	Category	1,737	63.7	42.4	-	-	-	-	-	-
	Measure	143	63.6	55.9	-	-	-0.58	-0.36	-	-
	Notes	310	-	-	-	-	-	-	47.2	43.5
	Rule	146	91.8	71.9	0.38	0.97	0.83	-0.51	-	-
CAP	Scale	8,412	59.6	47.0	1.07	1.66	-0.14	0.52	-	-
	Scale <sub>g</sub>	8,412	69.3	61.2	0.77	0.93	-0.15	0.53	-	-
	Category	400	81.0	64.5	-	-	-	-	-	-
	Measure	3,240	64.2	56.3	-	-	-0.33	0.01	-	-
	Rule	5,965	68.8	60.4	0.81	0.92	-0.19	0.46	-	-

Table 21: Model performances on 4 sections (§4.4) using four evaluation metrics (§5.4).

Variable	Count	Acc		RMSE		Bias		Recall	
		GPT	LM2	GPT	LM2	GPT	LM2	GPT	LM2
<b>Scale Variable</b>									
lbi_a2b	199	53.8	37.7	1.31	2.05	0.37	0.68	-	-
lbi_a3	201	56.7	52.7	0.96	0.99	0.31	0.47	-	-
lbi_a4	63	50.8	55.6	0.90	1.02	0.23	0.29	-	-
lbi_b1a_family	207	46.9	54.6	1.78	1.36	0.49	0.11	-	-
lbi_b2	212	60.8	44.8	1.19	1.02	0.47	0.49	-	-
lbi_d	194	52.1	38.7	1.29	1.29	0.53	0.53	-	-
lbi_e1	205	58.5	36.1	0.91	1.66	0.65	0.42	-	-
dx_understanding	29	55.2	51.7	1.20	1.25	0.23	0.43	-	-
dsm5capscritb01trauma1_distress	257	59.5	53.3	0.89	1.67	-0.44	0.48	-	-
dsm5capscritb02trauma1_distress	254	69.7	53.9	0.72	1.16	-0.56	0.32	-	-
dsm5capscritb03trauma1_distress	249	67.9	51.8	0.73	1.35	-0.05	0.68	-	-
dsm5capscritb04trauma1_distress	259	57.1	40.5	0.96	1.25	-0.35	0.47	-	-
dsm5capscritb05trauma1_distress	243	63.8	56.8	0.90	1.04	-0.11	0.37	-	-
dsm5capscritc01trauma1_distress	253	46.2	39.9	1.77	1.92	-0.34	0.53	-	-
dsm5capscritc02trauma1_distress	243	58.0	45.7	0.99	1.20	-0.04	0.64	-	-
dsm5capscritd01trauma1_distress	242	66.1	53.7	0.92	1.13	-0.10	0.30	-	-
dsm5capscritd02trauma1_distress	256	56.6	36.7	0.85	1.31	-0.06	0.83	-	-
caps5trauma1related_d02	164	57.9	55.5	0.97	0.85	-0.71	0.07	-	-
dsm5capscritd03trauma1_distress	248	61.7	58.9	0.94	0.92	-0.56	0.24	-	-
dsm5capscritd04trauma1_distress	252	56.0	49.2	0.93	1.13	-0.03	0.55	-	-
caps5trauma1related_d04	160	63.8	54.4	0.89	0.84	-0.28	0.10	-	-
dsm5capscritd05trauma1_distress	253	57.7	47.8	1.00	1.18	-0.08	0.53	-	-
caps5trauma1related_d05	138	53.6	44.9	1.06	0.96	-0.56	0.21	-	-
dsm5capscritd06trauma1_distress	255	53.5	47.5	1.01	1.23	0.09	0.66	-	-
caps5trauma1related_d06	156	51.3	41.0	0.98	0.90	-0.47	0.35	-	-
dsm5capscritd07trauma1_distress	257	59.5	45.5	0.88	1.22	0.04	0.67	-	-
caps5trauma1related_d07	128	55.5	44.5	0.96	0.94	-0.16	0.35	-	-
dsm5capscritc01trauma1_distress	257	60.3	46.7	0.79	1.13	0.33	0.78	-	-
caps5trauma1related_e01	148	52.7	33.8	3.54	3.44	-0.74	0.06	-	-
dsm5capscritc02trauma1_distress	251	67.3	61.0	0.71	1.11	0.02	0.31	-	-
caps5trauma1related_e02	50	74.0	58.0	1.09	1.26	-0.38	0.43	-	-
dsm5capscritc03trauma1_distress	255	51.4	47.1	1.09	1.20	0.32	0.54	-	-
caps5trauma1related_e03	155	50.3	51.6	0.93	0.86	-0.40	0.17	-	-
dsm5capscritc04trauma1_distress	252	63.1	52.8	0.85	1.05	-0.03	0.60	-	-
caps5trauma1related_e04	117	50.4	53.0	0.99	0.88	-0.55	0.13	-	-
dsm5capscritc05trauma1_distress	256	59.8	53.5	0.81	0.99	-0.13	0.58	-	-
caps5trauma1related_e05	161	57.8	41.6	1.09	0.99	-0.79	0.51	-	-
dsm5capscritc06trauma1_distress	256	53.5	52.7	1.02	1.06	0.09	0.37	-	-
caps5trauma1related_e06	181	63.0	38.7	1.0	10.2	-0.67	0.37	-	-
dsmiv_future_frequency_current	251	80.1	48.6	0.81	6.02	0.40	0.80	-	-
dsmiv_future_intens_current	246	69.1	40.2	0.93	1.77	0.61	0.90	-	-

continued on next page

continued from previous page

Variable	Count	Acc		RMSE		Bias		Recall	
		GPT	LM2	GPT	LM2	GPT	LM2	GPT	LM2
dsm5capscritg_trauma1_distress	228	53.9	43.4	1.08	1.39	0.35	0.69	-	-
dsm5capscritg_trauma1_impair	226	51.8	42.0	0.93	1.27	-0.28	0.57	-	-
dsm5capscritg_trauma1_fx	205	54.1	29.3	1.10	1.56	-0.04	0.81	-	-
dsm5depersonalization_sev	255	67.5	52.2	0.80	1.25	-0.08	0.49	-	-
caps5trauma1related_diss01	76	53.9	31.6	1.16	1.26	0.31	0.19	-	-
dsm5derealization_sev	249	63.1	30.9	0.98	1.74	0.20	0.88	-	-
caps5trauma1related_diss02	70	55.7	27.1	1.11	1.25	-0.03	0.53	-	-
<b>Category Variable</b>									
lbi_a1	200	70.0	41.0	-	-	-	-	-	-
lbi_student	201	95.0	89.1	-	-	-	-	-	-
lbi_c1a	192	57.8	71.9	-	-	-	-	-	-
lbi_c2	1	100	100	-	-	-	-	-	-
thh_medicalcond	206	92.7	88.8	-	-	-	-	-	-
thh_tx_curr_yesno	215	94.9	80.9	-	-	-	-	-	-
thh_tx_yesno	233	89.7	87.6	-	-	-	-	-	-
feedback_helpful	79	94.9	89.9	-	-	-	-	-	-
thh_txneed_yesno	96	92.7	88.5	-	-	-	-	-	-
thh_psychmed_curr_yesno	194	92.3	88.7	-	-	-	-	-	-
thh_psychmed_yesno	198	95.5	93.4	-	-	-	-	-	-
thh_suicide_yesno	236	90.7	77.1	-	-	-	-	-	-
thh_suicide_pw_yesno	70	94.3	78.6	-	-	-	-	-	-
trauma1lifeeventscl	146	61.6	12.3	-	-	-	-	-	-
trauma1_exposure_type___1	146	77.4	67.1	-	-	-	-	-	-
trauma1_exposure_type___2	146	77.4	43.2	-	-	-	-	-	-
trauma1_exposure_type___3	146	67.8	28.1	-	-	-	-	-	-
trauma1_exposure_type___4	146	65.8	22.6	-	-	-	-	-	-
caps_e1_lt	145	62.1	45.5	-	-	-	-	-	-
caps_e1_itself	73	64.4	64.4	-	-	-	-	-	-
caps_e1_ltother	74	41.9	44.6	-	-	-	-	-	-
caps_e1_si	146	43.8	39.7	-	-	-	-	-	-
caps_e1_siself	61	54.1	65.6	-	-	-	-	-	-
caps_e1_siother	61	60.7	29.5	-	-	-	-	-	-
caps_e1_tpi	146	54.1	52.7	-	-	-	-	-	-
caps_e1_tpiself	79	84.8	75.9	-	-	-	-	-	-
caps_e1_tpiother	77	49.4	26.0	-	-	-	-	-	-
trauma1_nomemory	145	75.2	44.8	-	-	-	-	-	-
dsm5caps_critf_cur1_yesno	202	78.7	41.6	-	-	-	-	-	-
dsm5caps_critf_cur1_c	198	83.3	87.9	-	-	-	-	-	-
<b>Measure Variable</b>									
lbi_a2a	99	68.7	66.7	-	-	41.9	45.5	-	-
trauma1_age	143	63.6	55.9	-	-	21.2	31.7	-	-
dsm5capscritb01trauma1_num	162	63.6	58.0	-	-	37.3	52.9	-	-

continued on next page

continued from previous page

Variable	Count	Acc		RMSE		Bias		Recall	
		GPT	LM2	GPT	LM2	GPT	LM2	GPT	LM2
dsm5capscrib02trauma1_num	98	74.5	63.3	-	-	28.0	52.8	-	-
dsm5capscrib03trauma1_num	84	72.6	59.5	-	-	47.8	76.5	-	-
dsm5capscrib04trauma1_num	177	62.1	58.8	-	-	17.9	42.5	-	-
dsm5capscrib05trauma1_num	137	59.1	57.7	-	-	28.6	50.0	-	-
dsm5capscritc01trauma1_num	170	59.4	53.5	-	-	31.9	54.4	-	-
dsm5capscritc02trauma1_num	140	63.6	54.3	-	-	27.5	54.7	-	-
dsm5capscritd01trauma1_num	87	50.6	48.3	-	-	32.6	82.2	-	-
dsm5capscritd02trauma1_num	168	76.2	69.0	-	-	47.5	59.6	-	-
dsm5capscritd03trauma1_num	120	65.0	57.5	-	-	23.8	43.1	-	-
dsm5capscritd04trauma1_num	166	72.3	68.1	-	-	39.1	45.3	-	-
dsm5capscritd05trauma1_num	138	65.9	59.4	-	-	42.6	46.4	-	-
dsm5capscritd06trauma1_num	155	69.7	63.2	-	-	27.7	40.4	-	-
dsm5capscritd07trauma1_num	140	61.4	60.0	-	-	40.7	53.6	-	-
dsm5capscritc01trauma1_num	135	65.9	62.2	-	-	21.7	54.9	-	-
dsm5capscritc02trauma1_num	61	83.6	68.9	-	-	80.0	89.5	-	-
dsm5capscritc03trauma1_num	159	73.0	67.3	-	-	37.2	34.6	-	-
dsm5capscritc04trauma1_num	131	68.7	60.3	-	-	24.4	50.0	-	-
dsm5capscritc05trauma1_num	168	69.0	66.1	-	-	21.2	35.1	-	-
dsm5capscritc06trauma1_num	184	72.8	61.4	-	-	40.0	31.0	-	-
dsmcaps_critf_cur1_nummonths	191	49.7	22.5	-	-	60.4	81.8	-	-
dsm5caps_critf_cur1_b	181	35.7	22.0	-	-	17.9	19.0	-	-
dsm5depersonalization_num	84	59.5	51.2	-	-	32.4	65.9	-	-
dsm5derealization_num	3	100	33.3	-	-	0.00	50.0	-	-
<b>Notes Variable</b>									
life_base_typicalday	203	-	-	-	-	-	-	42.0	50.8
thh_medicalcond_desc	100	-	-	-	-	-	-	56.8	80.1
thh_tx_curr_descr	59	-	-	-	-	-	-	53.6	73.8
thh_tx_descr	135	-	-	-	-	-	-	44.0	57.4
dx_knowledge	33	-	-	-	-	-	-	59.4	48.5
dx_lackknowledge	20	-	-	-	-	-	-	60.7	37.9
feedback_info	66	-	-	-	-	-	-	75.1	48.1
thh_txneed_desc	45	-	-	-	-	-	-	59.7	49.4
thh_psychmed_descr	89	-	-	-	-	-	-	40.4	59.0
thh_suicide_desc	73	-	-	-	-	-	-	56.9	67.0
thh_suicide_pw_desc	13	-	-	-	-	-	-	62.2	62.3
critaprobenotes	143	-	-	-	-	-	-	50.8	37.4
trauma1whathappened	143	-	-	-	-	-	-	42.7	51.4
trauma1describe	24	-	-	-	-	-	-	51.4	48.9
<b>Rule Variable</b>									
lbi_e2	215	43.3	37.7	0.94	0.98	0.44	0.43	-	-
caps_e1_crita	146	91.8	71.9	0.38	0.97	0.83	-0.51	-	-
dsm5capscrib01trauma1	253	62.8	63.6	0.81	0.90	-0.51	0.28	-	-

continued on next page

continued from previous page

Variable	Count	Acc		RMSE		Bias		Recall	
		GPT	LM2	GPT	LM2	GPT	LM2	GPT	LM2
dsm5capscrib02trauma1	250	88.0	70.0	0.53	0.80	-0.47	0.47	-	-
dsm5capscrib03trauma1	246	86.2	63.4	0.54	0.96	0.00	0.82	-	-
dsm5capscrib04trauma1	255	67.5	60.0	0.86	0.95	-0.57	0.25	-	-
dsm5capscrib05trauma1	241	74.7	69.7	0.73	0.75	-0.18	0.26	-	-
dsm5capscritc01trauma1	250	55.2	54.8	0.94	0.97	-0.64	0.36	-	-
dsm5capscritc02trauma1	242	71.9	61.2	0.83	0.93	-0.29	0.51	-	-
dsm5capscritd01trauma1	239	81.2	66.5	0.68	1.02	0.24	0.60	-	-
dsm5capscritd02trauma1	222	62.6	46.4	0.79	1.09	-0.16	0.83	-	-
dsm5capscritd03trauma1	246	72.0	72.0	0.85	0.74	-0.48	0.36	-	-
dsm5capscritd04trauma1	251	63.7	62.2	0.94	1.02	-0.08	0.35	-	-
dsm5capscritd05trauma1	252	59.9	53.6	0.98	1.00	-0.19	0.42	-	-
dsm5capscritd06trauma1	254	55.9	50.8	1.03	1.10	-0.07	0.57	-	-
dsm5capscritd07trauma1	255	63.1	60.0	0.85	0.95	-0.17	0.59	-	-
dsm5capscrit01trauma1	255	72.5	51.4	0.69	0.90	0.03	0.66	-	-
dsm5capscrit02trauma1	250	90.4	76.8	0.38	0.73	0.75	0.90	-	-
dsm5capscrit03trauma1	220	57.3	58.6	0.91	0.96	0.09	0.43	-	-
dsm5capscrit04trauma1	250	75.6	72.0	0.71	0.79	-0.08	0.63	-	-
dsm5capscrit05trauma1	254	65.7	67.7	0.80	0.77	-0.36	0.54	-	-
dsm5capscrit06trauma1	254	55.9	52.8	1.05	0.96	0.05	0.27	-	-
dsmcaps_critf_admin	28	75.0	100	0.50	0.00	-1.00	-1.00	-	-
dsm5depersonalization	246	85.4	64.2	0.61	0.78	-0.06	0.70	-	-
dsm5derealization	243	75.3	39.1	0.69	1.14	0.27	0.76	-	-
dsm5capsglobalvalidtrauma1	255	63.5	63.5	0.84	0.84	-1.00	-1.00	-	-
dsm5capsglobalsevtrauma1	254	44.1	42.9	0.91	0.97	0.21	0.45	-	-

Table 22: Model performances on all variable (§4.4) using four evaluation metrics (§5.4).

# DialBB: A Dialogue System Development Framework as an Educational Material

Mikio Nakano<sup>1,2</sup> and Kazunori Komatani<sup>2</sup>

<sup>1</sup>C4A Research Institute, Inc., 1-13-12 Umegaoka, Setagaya, Tokyo, Japan

<sup>2</sup>SANKEN, Osaka University, 8-1 Mihogaoka, Ibaraki, Osaka, Japan

mikio.nakano@c4a.jp, komatani@sanken.osaka-u.ac.jp

## Abstract

We demonstrate DialBB, a dialogue system development framework, which we have been building as an educational material for dialogue system technology. Building a dialogue system requires the adoption of an appropriate architecture depending on the application and the integration of various technologies. However, this is not easy for those who have just started learning dialogue system technology. Therefore, there is a demand for educational materials that integrate various technologies to build dialogue systems, because traditional dialogue system development frameworks were not designed for educational purposes. DialBB enables the development of dialogue systems by combining modules called building blocks. After understanding sample applications, learners can easily build simple systems using built-in blocks and can build advanced systems using their own developed blocks.

## 1 Introduction

To build a dialogue system, it is generally necessary to adopt an appropriate architecture according to the application and integrate various technologies. While the advancements in large language models have led some to believe that dialogue systems can be developed solely with these models and that the developers do not need to know about architecture and elemental technologies, there are issues such as hallucinations, so it is not always practical to build dialogue systems with only large language models, depending on the application.

Over the years, research into dialogue systems has evolved, accumulating knowledge on what kind of dialogue systems should be built with what technologies and what architectures. However, it is not easy for people who are learning dialogue system technology to acquire this knowledge. An educational material that allows people to learn dialogue system technology while building various dialogue systems would be helpful.

In learning about dialogue system technology, it is important to understand various elemental technologies such as language understanding and dialogue management, as well as an architecture based on appropriate modularization, extensibility which facilitates improving systems, and domain portability. It is also crucial to understand the importance of robustness in intention understanding and interaction design through running actual dialogue systems.

As an educational material that is useful for such learning, a dialogue system development framework with the following features is beneficial: (1) including various elemental technologies of dialogue systems, (2) appropriately modularized, (3) highly extensible, (4) including sample applications that help learners' understanding of dialogue system technology, and (5) making it possible to develop simple applications without extensive skills or knowledge in system development, enabling the acquisition of various technologies while improving the system. In addition, it is desired that its source code is available.

There are several dialogue system development tools whose source codes are available. PyDial (Ultes et al., 2017), OpenDial (Lison and Kennington, 2016), ConvLab-3 (Zhu et al., 2023), and ADVISER (Ortega et al., 2019) focus on statistical dialogue models for task-oriented dialogue systems, while we think educational materials should support state-transition network-based dialogue management which is often used for building practical dialogue systems. Although Rasa Open Source (Bocklisch et al., 2017) is highly extensible, it does not support state-transition network-based dialogue management by default. Botpress<sup>1</sup> supports state-transition network-based dialogue management, but replacing its internal modules with custom-made ones is not easy. MMDAgent (Lee et al.,

<sup>1</sup><https://botpress.com>

2013) also supports state-transition network-based dialogue management, but it is not easy to extend it.

We have been building a dialogue system development framework called **DialBB** (*Dialogue system development framework with Building Blocks*)<sup>2</sup> intended for use as an educational material in dialogue system development. DialBB is written in Python, and supports the development of English and Japanese applications.

## 2 Overview of DialBB

Here, we give an overview of DialBB. For more details, please refer to its document.<sup>3</sup>

### 2.1 Architecture

DialBB allows the development of dialogue systems by combining modules called Building Blocks (hereafter referred to as "blocks"). Figure 1 shows the architecture of DialBB applications.

The main module of DialBB works as follows. First, it receives input containing user utterances in JSON format through a method call of the class API or via a Web API. This input is then stored in the blackboard.<sup>4</sup> Next, it calls each block in the order specified in the configuration file, using parts of the blackboard as input for these blocks. The output from the blocks is used to update the blackboard. By sequentially driving each block in this manner, the system generates and returns a response. Additionally, the input and output of the main module can include not only utterance strings but also additional information, so that it is possible to handle multimodal information such as speech recognition confidences, user emotion estimation results, and gesture commands.

Which blocks each application uses is specified by describing the block classes in the application's configuration file (a YAML file). The configuration file also specifies what type of data each block receives and sends. Furthermore, the values of parameters used within the blocks and the knowledge description files used by the blocks can also be specified in the configuration file.

<sup>2</sup>DialBB is publicly available for non-commercial use at <https://github.com/c4a-ri/dialbb>. This paper is based on its ver. 0.8.

<sup>3</sup><https://c4a-ri.github.io/dialbb/document-en/build/html/>

<sup>4</sup>We call it a 'blackboard' in analogy to the blackboard model (Erman et al., 1980), but unlike the blackboard model, each block is called in the order written in the configuration file.

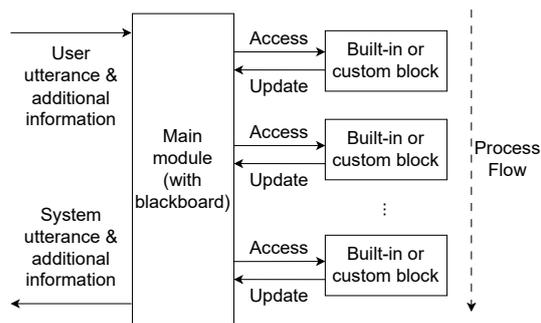


Figure 1: Architecture of DialBB applications.

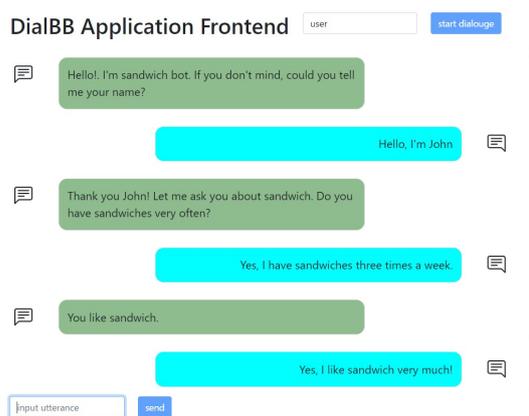


Figure 2: A snapshot of the frontend.

DialBB includes a frontend for engaging in dialogues via a Web API (Figure 2).

### 2.2 Built-in Blocks

To make it easier for learners to build conversational systems, DialBB has built-in building blocks listed in Table 1. For simplicity, only those for English applications are listed. Below we explain knowledge to be written by developers for use in some built-in blocks.

**Language Understanding Knowledge Snips** Understander Block and ChatGPT Understander Block use language understanding knowledge which consists of a collection of utterance examples that are annotated with intents and slots like the following.

Intent	Example utterance
tell-favorite-sandwich	I love (chicken salad sandwiches) [favorite-sandwich]
acknowledge	Definitely

Here, "[favorite-sandwich]" indicates a slot name, and "(chicken salad sandwiches)" indicates a slot

Block	Input	Output	Task
Simple Canonicalizer	user utterance string	canonicalized user utterance string	Canonicalizes the input string (convert uppercase to lowercase, etc.).
Whitespace Tokenizer	canonicalized user utterance string	token list	Performs tokenization based on white spaces.
Snips Understander <sup>5</sup>	token list	intent and slots	Performs language understanding using Snips NLU (Coucke et al., 2018) to obtain the intent and slots.
ChatGPT Understander	user utterance string	intent and slots	Performs language understanding using the JSON mode of OpenAI's ChatGPT. <sup>6</sup> Creates few-shot examples to embed in prompts from language understanding knowledge.
spaCy-Based NER	user utterance string	named entities	Performs named entity recognition using spaCy. <sup>7</sup>
STN Manager	user utterance string, intent, slots, and named entities (all are optional)	system utterance string	Manages dialogues using a state-transition network.
ChatGPT Dialogue	user utterance string	system utterance string	Generates a system utterance using ChatGPT based on a prompt including system persona, situation, and dialogue history.

Table 1: List of built-in blocks. Only important inputs and outputs are shown.

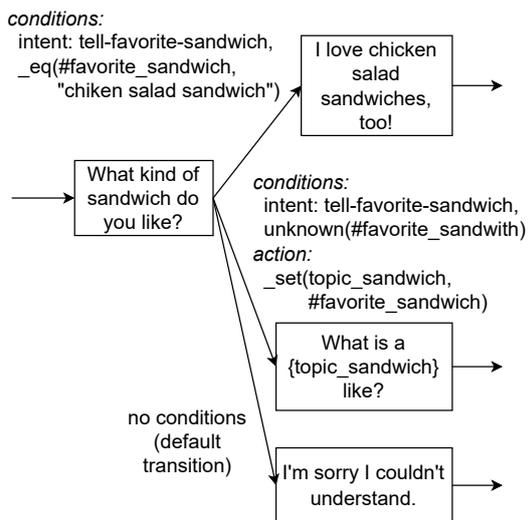


Figure 3: A part of an example state-transition network. Transitions above are given priority.

value. In addition, the knowledge used for language understanding includes a list of synonyms.

**State-Transition Network** STN Manager block uses a state-transition network (STN, also called a *scenario*). State-transition network-based dialogue

<sup>5</sup>Snips Understander Block will be deprecated in DialBB ver. 0.9 because Snips does not work with Python 3.9+. Instead DialBB ver. 0.9 will have a built-in block for language understanding that utilizes logistic regression and conditional random fields.

<sup>6</sup><https://openai.com/index/chatgpt>

<sup>7</sup><https://spacy.io/>

management is often used in practical dialogue systems. Figure 3 illustrates a part of the state-transition network. Each state is assigned a system utterance that is generated in that state. The state transitions to another state according to the input user utterance. Transitions can be accompanied by conditions for the transition and actions to be performed during the transition. Conditions are based on the intent of the user utterance and function calls. Actions are function calls. Functions used in conditions and actions are referred to as *scenario functions*. Within the definition of scenario functions, it is possible to use slots extracted in language understanding and named entity recognition results (for example, #favorite\_sandwich in Figure 3 is a slot value). Scenario functions can also access *contextual information*, which consists of data that the system remembers as the dialogue progresses, such as user requests and preferences (topic\_sandwich in Figure 3 is an example of this). Contextual information can be used in system utterances, as in “What is a {topic\_sandwich} like?” in Figure 3. It is also possible to access separately operating databases or external APIs in scenario functions.

Additionally, STN Manager block includes built-in scenario functions, which can reduce the effort of defining functions. For instance, there is a built-in function that compares if strings are identical (\_eq in Figure 3), and a built-in function that sets a value to a variable (\_set in Figure 3). Furthermore, there are built-in functions that use ChatGPT. One

is for determining if conditions written in natural language (e.g., “Is the user bored with the conversation?”) are satisfied, and the other is for generating utterance strings based on instructions written in natural language (e.g., “Generate a response to the user’s utterance in less than 30 words”). These functions call ChatGPT by incorporating into the prompt the dialogue history and the situation and persona settings specified in the configuration file.

Language understanding knowledge and state-transition networks can be described using spreadsheets.

STN Manager block has additional functionalities for handling speech recognition results. A DialBB application can receive speech recognition confidence together with the speech recognition result of the user utterance. STN Manager block can make an utterance to ask for repetition or request confirmation depending on the configuration. It is also possible to process barge-in utterances differently from ordinary utterances. In addition, reacting to a long silence after a system utterance is possible.

### 2.3 Sample Applications

DialBB has several sample English and Japanese applications that use only these built-in building blocks. Below are English applications.

*Snips+STN Application* uses Simple Canonicalizer, Whitespace Tokenizer, Snips Understander, and STN Manager blocks and it can engage in a simple dialogue about sandwiches.

*Lab Application* uses Simple Canonicalizer, ChatGPT Understander, spaCy NER, and STN Manager Blocks and it can also engage in a simple dialogue about sandwiches, but it demonstrates various advanced features of the built-in blocks.

*ChatGPT Application* uses only ChatGPT Dialogue block. It can engage in a dialogue using ChatGPT based on a prompt template that describes the dialogue situation and system persona.

To serve as a reference for learners, these sample applications, the built-in blocks, and the main module of DialBB are written in code that is as readable as possible.

### 2.4 Custom Blocks

Developers can create and use their own custom blocks. A block’s class can be created by inheriting from an abstract class `AbstractBlock` and implementing the necessary methods. The created class can then be specified in the configuration file for

use. This enables using different language understanding and dialogue management than those of built-in blocks.

## 3 Learning Dialogue System Technology Using DialBB

Using DialBB, learners can learn about dialogue system technology through the following steps. First, by understanding the sample applications, they learn the basic architecture of a dialogue system. Next, by looking at the change in behaviors after modifying the knowledge used in the sample applications, they understand elemental technologies. Then they deepen their understanding of the elemental technologies by building a new application using built-in blocks. Next, they understand the necessity of extensibility by creating and using their own custom blocks. Finally, by having people other than themselves use the system they built, they understand the importance of robustness in intention understanding and interaction design.

## 4 Usage Example of DialBB

DialBB was utilized in student projects and for system development for competitions. For instance, it was used to develop the system that won third place (Kubo et al., 2022) in the Dialogue Robot Competition 2022 (Minato et al., 2022) and the system that won second place (Yanagimoto et al., 2023) in the Dialogue Robot Competition 2023 (Minato et al., 2024). They are conversational robots that can recommend tourist destinations. DialBB applications work as their dialogue processing components. The dialogue processing component of the 2023 system incorporates the built-in Japanese Canonicalizer Block and the built-in STN Manager Block, along with a custom block that performs keyword-based language understanding, sentiment analysis, and affirmative/negative utterance classification.

## 5 Concluding Remarks

This paper presented DialBB, a framework for developing dialogue systems. DialBB serves as an educational material for dialogue system technology.

Currently, we are building a GUI-based editor for state-transition networks. Future improvements include adding new built-in blocks. Additionally, we plan to develop new sample applications, incorporating useful examples such as frame-based

dialogue management, database access, and handling speech and multimodal input/output.

We will demonstrate sample applications and explain their configuration files and knowledge for language understanding and dialogue management, to show how DialBB is useful in learning dialogue system technology.

## Acknowledgements

We would like to thank those who used the earlier versions of DialBB and gave us useful feedback.

This work was partly supported by JSPS KAKENHI Grant Number JP22H00536.

## References

- Tom Bocklisch, Joey Faulkner, Nick Pawlowski, and Alan Nichol. 2017. [Rasa: Open source language understanding and dialogue management](#). *Preprint*, arXiv:1712.05181.
- Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, Maël Primet, and Joseph Dureau. 2018. [Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces](#). *Preprint*, arXiv:1805.10190.
- Lee D. Erman, Frederick Hayes-Roth, Victor R. Lesser, and D. Raj Reddy. 1980. [The Hearsay-II speech-understanding system: Integrating knowledge to resolve uncertainty](#). *ACM Computing Surveys*, 12(2):213–253.
- Yuki Kubo, Ryo Yanagimoto, Hayato Futase, Mikio Nakano, Zhaojie Luo, and Kazunori Komatani. 2022. [Team OS’s system for Dialogue Robot Competition 2022](#). In *Proc. Dialogue Robot Competition 2022*.
- Akinobu Lee, Keiichiro Oura, and Keiichi Tokuda. 2013. [MMDAgent—a fully open-source toolkit for voice interaction systems](#). In *Proc. ICASSP*, pages 8382–8385.
- Pierre Lison and Casey Kennington. 2016. [OpenDial: A toolkit for developing spoken dialogue systems with probabilistic rules](#). In *Proceedings of ACL-2016 System Demonstrations*, pages 67–72, Berlin, Germany. Association for Computational Linguistics.
- Takashi Minato, Ryuichiro Higashinaka, Kurima Sakai, Tomo Funayama, Hiromitsu Nishizaki, and Takayuki Nagai. 2022. [Overview of Dialogue Robot Competition 2022](#). In *Proc. Dialogue Robot Competition 2022*.
- Takashi Minato, Ryuichiro Higashinaka, Kurima Sakai, Tomo Funayama, Hiromitsu Nishizaki, and Takayuki Nagai. 2024. [Overview of Dialogue Robot Competition 2023](#). In *Proc. Dialogue Robot Competition 2023*.
- Daniel Ortega, Dirk Vāth, Gianna Weber, Lindsey Vanderlyn, Maximilian Schmidt, Moritz Völkel, Zorica Karacevic, and Ngoc Thang Vu. 2019. [ADVISER: A dialog system framework for education & research](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 93–98, Florence, Italy. Association for Computational Linguistics.
- Stefan Ultes, Lina M. Rojas-Barahona, Pei-Hao Su, David Vandyke, Dongho Kim, Iñigo Casanueva, Paweł Budzianowski, Nikola Mrkšić, Tsung-Hsien Wen, Milica Gašić, and Steve Young. 2017. [PyDial: A multi-domain statistical dialogue system toolkit](#). In *Proceedings of ACL 2017, System Demonstrations*, pages 73–78, Vancouver, Canada. Association for Computational Linguistics.
- Ryo Yanagimoto, Yunosuke Kubo, Miki Oshio, Mikio Nakano, Kenta Yamamoto, and Kazunori Komatani. 2023. [User-adaptive tourist information dialogue system with yes/no classifier and sentiment estimator](#). In *Proc. Dialogue Robot Competition 2023*.
- Qi Zhu, Christian Geisshauser, Hsien chin Lin, Carel van Niekerk, Baolin Peng, Zheng Zhang, Michael Heck, Nurul Lubis, Dazhen Wan, Xiaochen Zhu, Jianfeng Gao, Milica Gašić, and Minlie Huang. 2023. [ConvLab-3: A flexible dialogue system toolkit based on a unified data format](#). *Preprint*, arXiv:2211.17148.

# A Multimodal Dialogue System to Lead Consensus Building with Emotion-Displaying

Shinnosuke Nozue<sup>1</sup>, Yuto Nakano<sup>1</sup>, Shoji Moriya<sup>1</sup>, Tomoki Ariyama<sup>1,2</sup>,  
Kazuma Kokuta<sup>1,2</sup>, Suchun Xie<sup>1</sup>, Kai Sato<sup>1</sup>, Shusaku Sone<sup>1</sup>, Ryohei Kamei<sup>1</sup>,  
Reina Akama<sup>1,2</sup>, Yuichiroh Matsubayashi<sup>1,2</sup>, Keisuke Sakaguchi<sup>1,2</sup>

<sup>1</sup>Tohoku University, <sup>2</sup>RIKEN

{nozue.shinnosuke.q5, nakano.yuto.t2, shoji.moriya.q7, tomoki.ariyama.s3, kokuta.kazuma.r3,  
xie.suchun.p7, kai.satou.r8, sone.shusaku.r8, ryohei.kamei.s4}@dc.tohoku.ac.jp,  
{akama, y.m, keisuke.sakaguchi}@tohoku.ac.jp

## Abstract

The evolution of large language models has enabled fluent dialogue, increasing interest in the coexistence of humans and avatars. An essential aspect of achieving this coexistence involves developing sophisticated dialogue systems that can influence user behavior. In this background, we propose an effective multimodal dialogue system designed to promote consensus building with humans. Our system employs a slot-filling strategy to guide discussions and attempts to influence users with suggestions through emotional expression and intent conveyance via its avatar. These innovations have resulted in our system achieving the highest performance in a competition evaluating consensus building between humans and dialogue systems. We hope that our research will promote further discussion on the development of dialogue systems that enhance consensus building in human collaboration.

## 1 Introduction

The emergence of large language models, such as GPT-4 (OpenAI, 2023), has facilitated highly fluent text-based conversations. Nevertheless, in many practical situations, dialogue systems require the capability to influence users through negotiation, persuasion, and consensus building (Zhan et al., 2024). In these advanced dialogue scenarios, it is essential to impact users' cognitive and emotional responses to induce alternation in their thoughts, opinions, and behaviors, yet research on these skills remains limited (Chawla et al., 2023).

In this context, the Dialogue System Live Competition 6 (DSL6) was organized to evaluate the current technological limitations and identify necessary components for creating systems capable of consensus building (Higashinaka et al., 2024). The competition set up a dialogue scenario wherein the system had to negotiate with the user in the context of conflicting goals: a system and a user

jointly plan a party, but the system desires a grand party while the user actually prefers a modest one. The systems are required to take the user's voice utterance as input and respond with avatar movements and synthesized voice. The demonstrations are evaluated based on three criteria to ensure they are contextually appropriate: relevance of utterance content, suitability of gestures and facial expression, and appropriateness of pause and voice modulation.

This paper presents our system<sup>1</sup> submitted to DSL6 (Nakano et al., 2023), which aims to build consensus with a user by guiding discussion with subdivided topics and conveying emotions and positions through its avatar. Our focus in designing the consensus building process is on agenda management and the clear communication of emotions and positions. In order to facilitate discussion between speakers with different objectives, our system employs a strategy of dividing a large topic into subtopics and guiding the user step by step. Specifically, to manage discussion flow, we predetermine a list of subtopics as "blanked slots" in the GPT-based system's prompt. Once a particular issue is resolved, the system introduces the next unresolved subtopic. Within each subtopic, negotiations progress through a sequence of proposals and responses (including acceptance, rejection, and counter-proposals) (Maynard, 2010). Throughout this process, the system is designed to articulate its intentions while responding. One crucial element in this effort is the expression of emotions, which is considered to influence the future actions of others in negotiations (Morris and Keltner, 2000; de Melo et al., 2011; Melo et al., 2012). When controlling the avatar, we incorporate facial expressions, voice modulation, and body poses to express emotions. Through these innovations, our system demonstrated the best performance in DSL6, and

<sup>1</sup><https://github.com/cl-tohoku/hagi-bot>

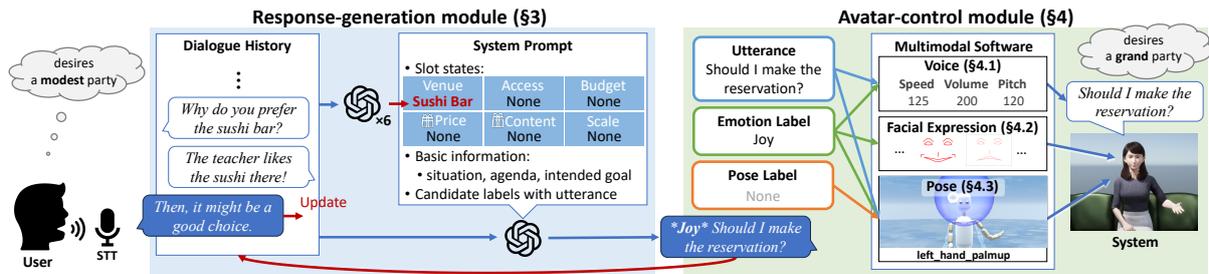


Figure 1: System overview. The user’s utterance is added to the dialogue history, and GPT-4 generates the system response based on the history and the system prompt. In the avatar-control module, the system’s voice, facial expressions, and poses are controlled according to the generated utterances and emotion/pose labels. Slots are updated in parallel with other processes after the generated response is added to the dialogue history.

feedback from the evaluators confirmed their effectiveness.

## 2 System Overview

Our system consists of two modules: the **response-generation module** and the **avatar-control module** (Figure 1). The response-generation module (Section 3) generates the system’s utterance and emotion/pose labels in response to the user’s utterance recognized by a speech-to-text (STT) function. The generated responses are forwarded to the avatar-control module (Section 4) that controls the voice, facial expressions, and poses according to predefined rules based on the received labels.

## 3 Response-Generation Module

For the response-generation module, we employed GPT-4 to generate outputs. We provided the module with a system prompt including basic information regarding the dialogue setting and the current status of slots (Section 3.1). Subsequently, the module generated an utterance and corresponding emotion/pose labels based on the provided prompt. Additionally, we ensured smooth turn-taking through pre- and post-module processing (Section 3.2).

### 3.1 Prompt Engineering

In the prompt, we detailed the basic situation of the dialogue, the agenda, and the intended objective of the discussion. We also provided some example utterances to instruct the output format of emotion/pose labels. Finally, we included the slot states of predefined subtopics. In the preliminary testing, we observed that the system often readily accepted user proposals that conflicted with its designated goal. To mitigate this tendency, we modified the prompt to include specific instructions: “If opin-

ions differ from those of the user, please engage in a discussion to make a decision while showing empathy.” This adjustment encouraged the system to express opposing views when necessary.

**Slot-Filling** We employed slot-filling-based dialogue state management to enable our system to lead discussions. We adopted GPT-4 as a slot-filling module for each subtopic individually. Each module dynamically updated the slot for the targeted subtopic with the determined content using the dialogue history (Figure 2). This approach enabled the system to start the discussion on an unfilled subtopic and facilitated smooth transitions to the next subtopic upon reaching a conclusion.

**Emotion and Pose Label Generation** In dialogue systems, emotion classification is typically performed independently from response generation (Moriya et al., 2023; Yamazaki et al., 2023). Consequently, errors in the classifier can lead to inconsistencies between utterance content and the avatar’s emotional expression. To address this issue, our system simultaneously generates utterance and emotion/pose labels, thereby ensuring coherence between utterance content and avatar expressions.<sup>2</sup> The candidate emotion labels were derived from Plutchik’s basic eight emotions (Plutchik, 2001), *Joy*, *Sadness*, *Anticipation*, *Surprise*, *Anger*, *Fear*, *Disgust*, and *Trust*, as utilized in the Japanese WRIME dataset (Kajiwara et al., 2021), with an additional *Neutral* label. Furthermore, we incorporated four pose labels—*Bowing*, *Nodding*, *Shaking head*, and *Pondering*—to represent the system’s intention and prevent misunderstanding.

<sup>2</sup>This was achieved by in-context learning using the prompt by providing an utterance example, such as “\*Pondering\* Hmm, I see your point... \*Joy\* In that case, it should be fine!”

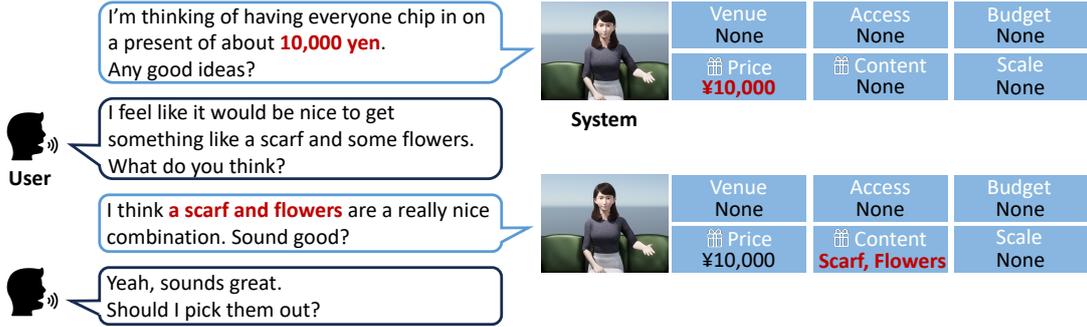


Figure 2: Example of dialogue and filling slots.

### 3.2 Natural Turn-Taking

It is crucial to avoid responses that are either too quick, which may interrupt the user’s speech, or too slow, which may lead the user to perceive the system as unresponsive (Skantze, 2021). Our system determines the end of the user’s turn based on the duration of silence. Through qualitative evaluations, a duration of 1.6 sec. was deemed suitable and adopted as the threshold. Additionally, in preliminary experiments, it was found that GPT-4 can take more than 10 sec. to generate a response when multiple sentences are involved. To address this issue, we employ the streaming mode. Specifically, we generated one character at a time while retaining and sending them to the text-to-speech (TTS) function upon generating a delimiter. This approach effectively reduced the response time from approximately 13 sec. to 3 sec.

## 4 Avatar-Control Module

In the avatar-control module, the system’s emotions and position are effectively conveyed by reflecting rules corresponding to emotion/pose labels or utterance content in the avatar’s voice, facial expressions, and poses. The avatar utilizes the resources provided by DSLC6 (Higashinaka et al., 2024).

### 4.1 Effect of Emotion Expressions

Emotion expressions serve as media to communicate the sender’s internal states, enabling the receiver to infer these states from the sender’s emotional displays (de Melo et al., 2023, 2011, 2014; Gratch and de Melo, 2019). This process significantly influences the thoughts and behaviors of the receiver, proving effective in various stages of negotiation, including trust development and consensus building. (de Melo et al., 2023; Morris and Keltner, 2000). Furthermore, it has been found that the

		speed	volume	pitch
Emotion Label	Joy	125	200	120
	Anticipation	120	150	117
	Sadness	120	100	105
	Surprise	125	250	115
	Anger	120	230	100
	Fear	125	250	115
	Disgust	120	100	95
	Trust	120	100	117
Neutral	120	100	115	
End of Sentence	“!”	+0	+50	+0
	“...”	-30	-50	-10

Table 1: Voice parameters. Our system utilizes AmazonPollyServer for its TTS function.

impact of such emotional conveying processes is also effective in human-machine interactions (Melo et al., 2012; de Melo et al., 2011). In this paper, we aim to achieve effective negotiation and consensus building by using competition-regulated tools to express emotions through voice, facial expressions, and poses (posture and gesture).

### 4.2 Voice

The avatar’s voice was modulated by predefined parameters such as speed, volume, and pitch corresponding to the emotion. This concept is adopted from the work of Togo et al. (2022). As a voice control guideline, we employed the two-dimensional arrangement of emotions in Russell’s Circumplex Model of Affect (Russell, 1980). We mapped the arousal-sleep dimension to the speed and volume and the pleasure-displeasure dimension to the pitch. Furthermore, when the utterance ends with “!” or “...,” we adjusted the parameters based on the polarity and intensity of the emotion. Table 1 shows the specific parameter values defined above. Additionally, in conversations, pauses express emotions and give listeners time to understand nuances (Nakamura, 2009). Accordingly, we introduced short

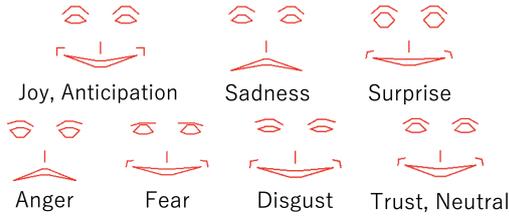


Figure 3: Facial expression set based on emotion labels. The expressions are configured using the JointMapper-PlusUltraSuperFace preset.



(a) Pondering (b) Surprised

Figure 4: Examples of pose control.

silences after punctuation, like periods, and longer silences after ellipses.

### 4.3 Facial Expressions

We predefined the avatar’s facial expressions for each of the nine emotion labels (Figure 3). The corresponding preset is referenced from the received emotion label and reflected in the avatar.

### 4.4 Control of Poses

The poses are determined by both the labels and the utterance content. We implemented fundamental behaviors such as nodding, but this paper focuses on the most crucial aspect: expressing emotions and intentions. To represent pondering, we utilized a hand-on-chin gesture, slightly lowering the face in response to the pose label *Pondering* or utterances such as “Hmm.” or “...” (Figure 4a). Similarly, placing a hand over the mouth represents a gesture of surprise, triggered by the emotion label *Surprise* or utterances indicating astonishment such as “Really?” (Figure 4b). Based on the Aoba\_v3 bot (Moriya et al., 2023), which aimed to behave like a human, we represented emotions such as *Sadness*, *Joy*, and *Anticipation*. For example, when emotions such as *Joy* or *Anticipation* were expressed, we defined the pose of lightly reaching out to the interlocutor. Conversely, when a label indicated *Sadness*, we expressed it by bowing the avatar’s head.

	Content	Expressions	Voice	Average
<b>Ours</b>	<b>3.93</b>	<b>3.41</b>	<b>3.01</b>	<b>3.45</b>
System A	3.32	3.32	2.94	3.19
System B	3.58	3.27	2.66	3.17

Table 2: Results of the final round in DSLC6. Scores range from 1 to 5. Three evaluation criteria are utterance content, gestures and facial expressions, and pause and voice modulation.

Evaluators’ comments and scores
<p>I sensed the system’s intention. The system progressed to the next, and the dialogue pace was good.  [Content, Expressions, Voice] = [5, 5, 5]</p>
<p>The conversations, gestures and pauses were very human.  [Content, Expressions, Voice] = [4, 4, 4]</p>
<p>I appreciated that the system shared its opinions and listened to mine. However, it spoke too quickly and seemed to talk quite a bit.  [Content, Expressions, Voice] = [3, 3, 3]</p>
<p>Her expressions seem exaggerated, which can be tiring.  [Content, Expressions, Voice] = [4, 3, 3]</p>

Table 3: Evaluators’ comments with scores.

## 5 Human Evaluation

Three top-performing systems, selected from an initial pool of ten systems in the preliminary round, participated in the final round of DSLC6. At the final round, each system engaged in a five-minute conversation in Japanese with a user who is a humanities researcher. The interactions were evaluated by an audience of 80 attendees. Each system was evaluated twice. The assessments were on the basis of contextual appropriateness, considering the following three criteria: (i) Content: relevance of utterance content, (ii) Expressions: suitability of gestures and facial expressions, and (iii) Voice: appropriateness of pause and voice modulation. Using a 5-point Likert scale, the systems’ performances were ranked based on the average scores of all criteria.

The results in Table 2 indicate our system achieved the best performance across all criteria. Table 3 contains some of the feedback from the preliminary evaluators. According to the positive feedback, the system was effective in leading discussions and naturally conveying intentions through the avatar. However, some users provided negative feedback regarding the non-verbal aspects, such as rapid speaking and excessive movement, indicating that there is still room for improvement.

## 6 Conclusion

This paper presented the top-performing system in a dialogue competition focused on the consensus building process between systems and humans. Our proposed system incorporated two strategies for smooth consensus building: topic control through a slot-filling approach and conveying intent through emotional expression via an avatar.

Our system can be adapted for various situations. The only task-specific components are the prompts and slots. The prompts include a general instruction format necessary for the consensus building process, such as context, personas, subtopics, and the system's objectives, thereby reducing the effort required for prompt engineering when applied to different tasks. Each topic slot is managed by a separate model, allowing easy instantiation and repurposing once the subtopics are defined.

We hope that this system will contribute insights to research on dialogue systems to build consensus with humans in the context of conflicting goals.

## Acknowledgements

We would like to express our gratitude to the organizers of DSLC6 for providing the software necessary for the development of the multimodal dialogue system, as well as to Prof. Kentaro Inui and Ms. Fuka Narita of the Tohoku NLP Group for their cooperation in the development of this system. This research was supported in part by JSPS KAKENHI Grants JP22K17943 and JP21K21343.

## References

- Kushal Chawla et al. 2023. [Social influence dialogue systems: A survey of datasets and models for social influence tasks](#). In *EACL*, pages 750–766.
- Celso M. de Melo et al. 2011. The effect of expression of anger and happiness in computer agents on negotiations with humans. In *The 10th International Conference on Autonomous Agents and Multiagent Systems - Volume 3*, AAMAS '11, page 937–944.
- Celso M. de Melo et al. 2014. [Reading people's minds from emotion expressions in interdependent decision making](#). *Journal of personality and social psychology*, 106 1:73–88.
- Celso M. de Melo et al. 2023. [Social functions of machine emotional expressions](#). *Proceedings of the IEEE*, 111(10):1382–1397.
- Jonathan Gratch and Celso M. de Melo. 2019. [Inferring Intentions from Emotion Expressions in Social Decision Making](#), pages 141–160. Springer International Publishing, Cham.
- Ryuichiro Higashinaka et al. 2024. Dialogue system live competition goes multimodal: Analyzing the effects of multimodal information in situated dialogue systems. In *IWSDS*.
- Tomoyuki Kajiwara et al. 2021. WRIME: A new dataset for emotional intensity estimation with subjective and objective annotations. In *NAACL: Human Language Technologies*, pages 2095–2104.
- Douglas W. Maynard. 2010. [Demur, defer, and deter: Concrete, actual practices for negotiation in interaction](#). *Negotiation Journal*, 26(2):125–143.
- Celso Melo, Peter Carnevale, and Jonathan Gratch. 2012. [The impact of emotion displays in embodied agents on emergence of cooperation with people](#). *PRESENCE: Teleoperators and Virtual Environments*, 20:449–465.
- Shoji Moriya, Daiki Shiono, et al. 2023. Aoba\_v3 bot: a multimodal chatbot system combining rules and various response generation models. *Advanced Robotics*, 37(21):1392–1405.
- Michael W. Morris and Dacher Keltner. 2000. [How emotions work: The social functions of emotional expression in negotiations](#). *Research in Organizational Behavior*, 22:1–50.
- Toshie Nakamura. 2009. [Psychological study of 'ma' \(a synonym of 'pause'\) in communication](#). *Journal of the Phonetic Society of Japan*, 13(1):40–52.
- Yuto Nakano, Shinnosuke Nozue, et al. 2023. [Hagi bot: A multimodal dialogue system for smooth discussion with human-like behavior and llm-based dialogue state tracking](#). In *JSAI SIG-SLUD*, 99:102–107.
- OpenAI. 2023. [GPT-4 technical report](#).
- Robert Plutchik. 2001. The nature of emotions: Human emotions have deep evolutionary roots, a fact that may explain their complexity and provide tools for clinical practice. In *American Scientist, Vol. 89, No. 4 (JULY-AUGUST)*, pages 344–350.
- James A. Russell. 1980. [A circumplex model of affect](#). *Journal of Personality and Social Psychology*, 39(6):1161–1178.
- Gabriel Skantze. 2021. [Turn-taking in conversational systems and human-robot interaction: A review](#). *Computer Speech & Language*, 67:101178.
- Eitetsu Togo et al. 2022. [A multimodal dialogue system with control based on emotion classification of system utterances](#). In *JSAI SIG-SLUD*, 96:206–210.
- Takato Yamazaki et al. 2023. [An open-domain avatar chatbot by exploiting a large language model](#). In *SIGDIAL*, pages 428–432.
- Haolan Zhan et al. 2024. [Let's negotiate! a survey of negotiation dialogue systems](#). *Preprint*, arXiv:2402.01097.

# PersonaCLR: Evaluation Model for Persona Characteristics via Contrastive Learning of Linguistic Style Representation

Michimasa Inaba

The University of Electro-Communications  
1-5-1, Chofugaoka, Chofu, Tokyo, Japan  
m-inaba@uec.ac.jp

## Abstract

Persona-aware dialogue systems can improve the consistency of the system's responses, users' trust and user enjoyment. Filtering nonpersona-like utterances is important for constructing persona-aware dialogue systems. This paper presents the PersonaCLR model for capturing a given utterance's intensity of persona characteristics. We trained the model with contrastive learning based on the sameness of the utterances' speaker. Contrastive learning enables PersonaCLR to evaluate the persona characteristics of a given utterance, even if the target persona is not included in training data. For training and evaluating our model, we also constructed a new dataset of 2,155 character utterances from 100 Japanese online novels. Experimental results indicated that our model outperforms existing methods and a strong baseline using a large language model. Our source code, pre-trained model, and dataset are available at <https://github.com/1never/PersonaCLR>.

## 1 Introduction

Persona-aware dialogue systems can improve the consistency of the system's responses (Li et al., 2016), users' trust in the system (Higashinaka et al., 2018), and user enjoyment (Miyazaki et al., 2016).

In constructing persona-aware dialogue systems, automatic estimation of persona characteristics' intensity is important in two ways. First, if we can detect low-intensity utterances of persona characteristics, inappropriate system responses can be prevented. Second, the automatic measure helps construct the persona's sample utterance set. Two methods for constructing persona-aware dialogue systems include the following: (1) using persona descriptions (Zhang et al., 2018; Song et al., 2019; Majumder et al., 2020; Kim et al., 2020; Tang et al., 2023) and (2) sample utterances (Higashinaka et al., 2018; Mitsuda et al., 2022; Han et al., 2022). In the method using samples, if we

can filter out samples not matching the persona, the system's performance will improve. This paper presents an evaluation model for Persona characteristics via Contrastive learning of Linguistic style Representation (PersonaCLR), which can measure a given utterance's intensity of the target persona's characteristics. In this paper, the term persona indicates both real-life individuals and fictional characters. PersonaCLR receives the evaluation target's utterance and a target persona's set of sample utterances and then returns a score indicating the target persona characteristics' intensity. The model is trained by contrastive learning based on the sameness of the utterances' speaker. Contrastive learning enables PersonaCLR to evaluate the persona characteristics of a given utterance, even if the target persona is not included in training data.

To the best of our knowledge, there are no public datasets for the training and evaluation of models to assess the intensity of persona characteristics in utterances. We constructed and published two datasets: the Naro Utterance dataset (NaroU), containing 2,155 characters' utterances from 100 Japanese online novels, and an evaluation dataset based on dialogue scenarios between a user and a character in NaroU. We use the dataset to train and evaluate PersonaCLR in the experiments. The creation of these datasets contributes to the advancement of research on dialogue systems that mimic fictional characters. Additionally, this dataset can be utilized for speaker identification tasks (He et al., 2013; Muzny et al., 2017; Yu et al., 2022) and for building persona-aware dialogue systems using sample utterances (Han et al., 2022). The evaluation dataset for this task was also constructed and published.

This study contributes the following: (1) a new model for assessing the intensity of target persona characteristics in a given utterance that does not require retraining or fine-tuning, even if the persona is not included in the training data; (2) a new

open dataset including over 2000 character utterances from 100 Japanese online novels and annotated human-character dialogue scenarios; and (3) a demonstration of the effectiveness of our model using a comparison with existing methods and a strong baseline involving ChatGPT.

## 2 Related Work

### 2.1 Persona Characteristics Evaluation

Persona-based dialogue models have been actively studied (Song et al., 2019; Majumder et al., 2020; Kim et al., 2020, 2022) since the release of the PERSONA-CHAT dataset (Zhang et al., 2018). These models receive a dialogue context and a few sentences of persona description (e.g., “I have two dogs.”) and then include the description’s content as much as possible in generated responses. Several proposed evaluation metrics for these models evaluate generated utterances according to how much of a given persona description’s content is included. Both Persona F1 (Jiang et al., 2020a) and Persona coverage (Jiang et al., 2020a) are metrics that utilize nonstop words common between a given persona description and an utterance. Persona accuracy (Zheng et al., 2020), which predicts whether a given persona description is exhibited in generated utterances, is computed by feeding generated responses into a binary classifier and obtaining classification accuracy.

This study focuses on persona-aware dialogue systems that mimic a fictional character rather than on systems based on persona descriptions. Since defining such personas with only a few descriptive sentences is difficult, several methods have been proposed to construct such persona-aware dialogue systems using a few samples (Han et al., 2022) or manually collected responses (Higashinaka et al., 2018; Mitsuda et al., 2022). For the same reasons as above, systems’ evaluation by methods based on the given persona description’s content is difficult. Therefore, Persona Speaker Probability (PSProb) (Miyazaki et al., 2021) and Persona Term Saliency (PTSali) (Miyazaki et al., 2021) have been proposed as evaluation metrics for dialogue systems’ utterances that mimic a fictional character’s persona.

### 2.2 Contrastive Learning

In computer vision, contrastive unsupervised representation learning has been proposed, and performance in object detection and image segmentation has significantly improved (He et al., 2020). The

key idea is that this type of learning minimizes the distance between feature representations of different views of the same image and maximizes between-feature representations of views of different images (Chen et al., 2020). Contrastive learning has also been applied to natural language processing, and various models for learning sentence representations have been proposed (Fang et al., 2020; Chen et al., 2020; Giorgi et al., 2021).

Particularly relevant to our study is supervised contrastive learning (Khosla et al., 2020; Gunel et al., 2021; Zhang et al., 2022), which constructs positive and negative pairs by leveraging ground truth labels. Inspired by supervised contrastive learning, we constructed contrastive learning pairs based on the sameness of the utterances’ speaker.

### 2.3 Novel Dataset and Speaker Identification

Several corpora with speaker annotations based on novels have been constructed for several languages: the Columbia Quoted Speech Attribution Corpus (Elson and McKeown, 2010), P&P (He et al., 2013), QuoteLi3 (Muzny et al., 2017), and RiQuA (Papay and Padó, 2020) are English corpora; WP (Chen et al., 2019, 2021), JINYONG (Jia et al., 2020) and CSI (Yu et al., 2022) are Chinese; and RWG (Bruner, 2013) is German. In these corpora, speaker annotations were performed for a few of the novels (the highest was 18 in CSI). Thus, the diversity of worldviews and characters are limited. We annotated 100 online novels written in Japanese and constructed and released a new dataset.

Existing corpora have been mainly constructed for speaker identification (SI), that is, to identify the corresponding speaker(s) for each utterance in novels (He et al., 2013; Muzny et al., 2017; Yu et al., 2022; Chen et al., 2023). In SI, an utterance and its surrounding context are given, and, using the context, SI models determine the utterance’s speaker. Our task can be regarded as predicting a given utterance’s speaker, but because no context is given in our task, we cannot apply existing SI methods.

## 3 PersonaCLR

Our task is to estimate the intensity of the characteristics of a target persona  $c$  within a given utterance  $x$ . The existing SoTA model, PSProb (Miyazaki et al., 2021), is based on a multi-class classification model classifying which character uttered the given input utterance. Therefore, when evaluat-

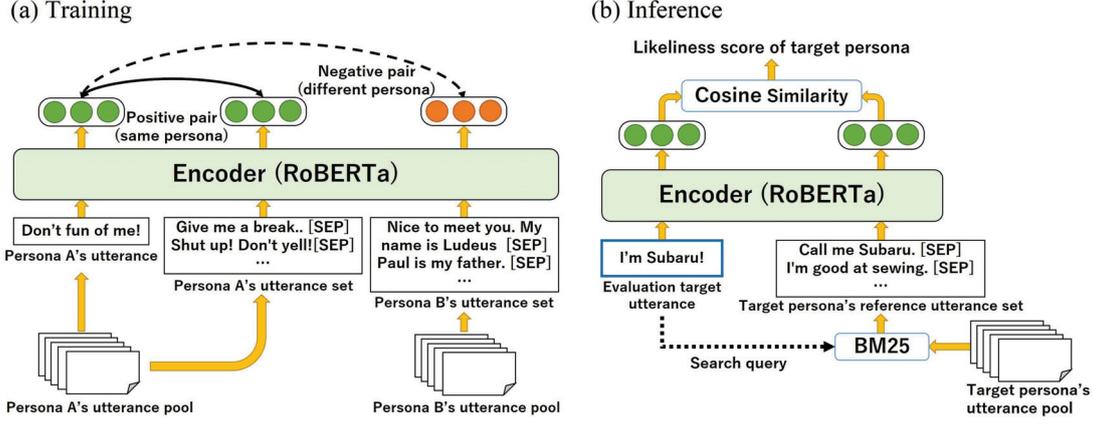


Figure 1: (a) Training with contrastive learning in PersonaCLR. An utterance and an utterance set were sampled from the utterance pool and encoded by the Transformer encoder RoBERTa. Pairs of encoded utterance and utterance sets are used as positive pairs for contrastive learning if they are sampled from the same persona’s utterance pool and as negative pairs otherwise. (b) Inference in PersonaCLR. The reference utterance set is constructed from the target character’s utterance pool using BM25. The likeliness score is obtained as the cosine similarity between encoded vectors from the target utterance and the utterance set.

ing the persona characteristics of a new character, PSProb must collect not only the reference utterance sets for the evaluation of the target persona as well as the utterance sets of non-target personas, then training the model from scratch.

We propose PersonaCLR, which does not require the retraining and utterance sets for non-target personas. PersonaCLR uses contrastive learning to distinguish whether a given utterance and a set of utterances come from the same persona. Note that PersonaCLR does not evaluate the similarity of any given utterance to utterances in the training data; rather, it observes the similarity between the given utterance and a reference set of utterances. Therefore, our model can evaluate utterances of characters that are not included in the training data without requiring retraining, and it requires only a small number ( $\geq 20$ ) of references.

We define two embedding vectors of the same speaker’s utterances as positive pairs and two vectors by different speakers as negative pairs for contrastive learning. However, because utterances do not necessarily reflect persona characteristics, one of the pair’s embedding vectors is obtained from a set of utterances rather than from a single utterance. In contrastive learning, the model distinguishes whether an utterance and a set of utterances are from the same persona.

### 3.1 Training and Inference

An overview of PersonaCLR’s training is shown in Figure 1 (a). Let  $x^a = \{x_i^a\}_{i=1}^n$ , be the utter-

ances pool by a speaker  $a$ . By sampling  $x^a$ , we obtain an utterance  $x_k^a$  and reference utterance set  $x^{a+} = \{x_j^a\}_{j=1}^m$ . We use the pair of the utterance  $x_k^a$  and reference utterance set  $x^{a+}$  as positive pair for contrastive learning. On the other hand, we use the pair of the utterance  $x_k^a$ , and utterance set  $x^{b+}$ , sampled from speaker  $b$ ’s utterance pool  $x^b$  as negative pair.

Each utterance and reference utterance set is encoded by the Transformer encoder RoBERTa (Liu et al., 2020). Before encoding, the utterance sets are concatenated with a separator token  $[SEP]$  to form a single sequence. With RoBERTa, we obtain the embedding vectors  $\mathbf{h}_k^a$  and  $\mathbf{h}^{a+}$  corresponding to  $x_k^a$  and  $x^{a+}$ .

The loss function using these embedded vectors is defined as follows:

$$l_k^a = -\log \frac{e^{\text{sim}(\mathbf{h}_k^a, \mathbf{h}^{a+})/\tau}}{\sum_{i=1}^N e^{\text{sim}(\mathbf{h}_k^a, \mathbf{h}^{s_i+})/\tau}} \quad (1)$$

where  $N$  is the batch size and  $s_i$  is the speaker of the  $i$ th utterance set in the batch. The  $\tau$  is the temperature hyperparameter and  $\text{sim}(h_1, h_2)$  is the cosine similarity  $\frac{\mathbf{h}_1^\top \mathbf{h}_2}{\|\mathbf{h}_1\| \|\mathbf{h}_2\|}$ .

In the inference phase, the persona-likeness score of a given utterance is the cosine similarity between the utterance’s embedding vectors and the target persona’s reference set in Figure 1 (b).

### 3.2 Related Utterance Sampling

A reference utterance set should contain sufficient information to evaluate the utterance’s persona

characteristics. For example, if the target utterance contains a habit or terminology unique to the persona, for appropriate evaluation, the reference utterance set should also include it.

To this end, we ranked the utterance pool  $x^a$  using BM25 (Robertson et al., 1995). We used evaluation target utterance  $x_k^a$  as a query and obtain the top  $m$  utterances as the reference set  $x^{a+}$  in the inference phase (see Figure 1 (b)). In the experiment, all utterances in the training data of the NaroU (see Section 4.5) were used to calculate parameters for then calculating BM25 (average number of words per utterance and IDF). In the inference phase, we use the reference utterances set using BM25. In the training phase, we simultaneously use the utterance set using BM25 and the randomly sampled set as training data to ensure robustness.

## 4 Experiment

Figure 2 shows a summary of experimental procedure used to evaluate the effectiveness of PersonaCLR. In this experiment, we constructed and used two types of datasets: NaroU and the evaluation dataset. The NaroU dataset consists of utterances from novels. However, there is the concern that, compared to utterances in novels, utterances in a dialogue between a user and a persona-aware dialogue system differ in length and tendency. To address this concern, we created the evaluation dataset.

### 4.1 Naro Utterance Dataset (NaroU)

For training models to assess the intensity of persona characteristics in utterances, we constructed NaroU, a dataset of utterances in novels annotated with speaker attributions. This dataset was constructed by annotating 100 novels in “*Shosetsuka ni Naro*,”<sup>1</sup> a Japanese novel self-publishing website<sup>2</sup>. Most of the website’s novels are divided into episodes of 2000 to 5000 Japanese characters each, and we annotated each novel’s first ten episodes. We recruited annotators via the crowdsourcing website CrowdWorks<sup>3</sup> and instructed them to extract segments of utterances in the novel and assign speaker names. We instructed annotators to annotate the speaker’s real name if it was given in the novel or otherwise, a nickname or pronoun. One annotator performed annotation per each novel. We

<sup>1</sup><https://syosetu.com/>

<sup>2</sup>“*Shosetsuka ni Naro*” means “*Let’s become a novelist.*”

<sup>3</sup><https://crowdworks.jp/>

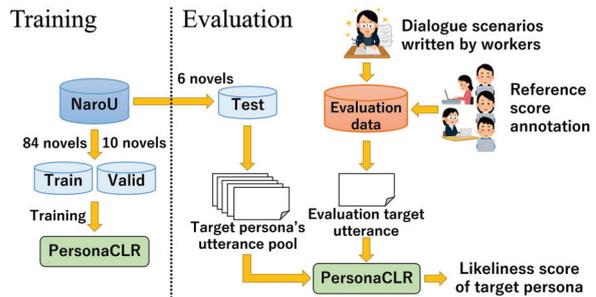


Figure 2: Experimental overview. Two datasets were used, NaroU (blue) and evaluation data (red). NaroU is divided into train, valid, and test datasets so that novels and characters do not overlap among the three. The six novels and their characters in the test data are shown in Table 2. Evaluation data is based on dialogue scenarios between a user and a target character created by crowd workers. All character utterances in the dialogue scenarios are annotated with reference scores of persona characteristics. During the evaluation phase, PersonaCLR is given an evaluation target utterance in dialogue scenarios and the corresponding persona’s utterance pool; it then outputs an estimated score.

Novels	100
Episodes	1000
Characters	2,155
Utterances	38,297
Words in utterances	620,820

Table 1: Statistics of the NaroU

paid them 600 JPY (approximately 4 USD) per episode. Table 1 shows the statistics of this dataset.

To confirm annotations’ consistency, we conducted an experiment in which 50 episodes of five novels were individually annotated by two annotators. Experimental results showed that the perfect agreement rate was 88.4%.

NaroU is divided into train, valid, and test datasets (84, 10, and 6 novels, respectively). The train and valid datasets are used to train the model, and the test dataset is used as the utterance pool (see Figure 1 (b)). Train, valid, and test datasets contained no overlap of novels and characters.

### 4.2 Evaluation Data Collection

Utterances from novels are not suitable for assessing PersonaCLR’s ability to estimate the intensity of a persona in a system’s utterances. For this experiment, therefore, we created dialogue scenarios in which a user interacts with a character and then used these scenarios’ utterances.

Ncode	Character
n6316bn	Rimuru (75), Veldora (63)
n9669bk	Rudeus (264), Roxy (140) Paul (105), Sylphiette (65), Zenith (60)
n2267be	Subaru (220)
n4830bu	Myne (187), Tuuli (83), Effa (63)
n3191eh	Leon (353), Luxion (86), Angelica (76), Olivia (67)
n5040ce	Catarina (190), Keith (52)

Table 2: List of target characters used for collecting evaluation data. Ncode is a unique ID assigned to each novel submitted to *Shosetsuka ni Naro*. Numbers in brackets indicates the number of unique utterances in the test dataset of NaroU used as the utterance pool of the target character. The utterances in the novels and characters shown in this table were not included in the training data for PersonaCLR.

#### 4.2.1 Dialogue Scenario

To collect the evaluation dataset, we recruited workers via CrowdWorks. They created dialogue scenario between a specified character and a user. For a situation in which a user is talking with a character, each scenario was individually created by one worker. We paid them 250 JPY (approximately 1.5 USD) per created dialogue.

We selected characters as evaluation target personas based on the following two conditions:(1) novels within the top 100 in cumulative ranking calculated on the number of bookmarks and reviews on *Shosetsuka ni Naro* and (2) novels developed as media mixes in both manga and anime as of January 2023. We selected 17 regular characters with 50 or more unique utterances from novels that satisfied the two conditions above. The list of characters is shown in Table 2, in which the ncode is an identifier uniquely assigned to each novel. We accessed the novels at [https://ncode.syosetu.com/\(ncode\)/](https://ncode.syosetu.com/(ncode)/).

We prepared 13 general topics for dialogue (e.g., hobbies, travel, and family), selected by the workers. Each scenario consisted of 10 utterances, spoken alternately by the user and the character. We recruited only workers who had watched at least ten episodes of a novel’s anime, read 10 episodes of the novel, or two books of the novel’s manga containing the character. Through the procedure described above, we collected 20 dialogues (100 utterances) per character, for a total of 1,700 utterances.

#### 4.2.2 Reference Score Annotation

To obtain reference scores of persona characteristics in the evaluation dataset, crowd workers an-

Characters	17
Unique utterances	1,700
Words per utterance	31.56
Reference scores	
- Score 5	967 (43.1%)
- Score 4	422 (18.8%)
- Score 3	168 (7.5%)
- Score 2	120 (5.3%)
- Score 1	153 (6.8%)
- Score 0	414 (18.4%)
Total	2,244

Table 3: Statistics of evaluation data

notated the collected utterances. Workers were paid 70 JPY (approximately 0.4 USD) per ten utterances. The definition and assignment procedure of the reference score complied with the previous study (Miyazaki et al., 2021). We only extracted the target character’s utterances from the created scenarios. We asked workers to answer with “yes” or “no” whether the character was likely to have said each utterance. The workers evaluated them only by observing the utterances, without considering the context. Five people evaluated each utterance, and the number of “yes” answers was used as the reference score.

Because the number of low scores was small for only utterances created as target characters, we also annotated utterances created as other characters. The previous study (Miyazaki et al., 2021) used 500 utterances and scores (100 utterances of target characters and 400 utterances of non-target characters). However, this setting was far from reality, with, in some cases, a score of 0 accounting for 60% of the total. Therefore, we annotated 100 utterances of the target character and two randomly sampled utterances of non-target characters for each character, for a total of 132 utterances (= 100 + 16 characters × 2). finally, we obtained 2,244 evaluation scores (= 132 × 17 characters). Table 3 displays evaluation data statistics, and Table 3 shows examples of utterances and scores.

### 4.3 Comparative Methods

#### 4.3.1 Persona Speaker Probability (PSProb)

PSProb (Miyazaki et al., 2021) is a previous SoTA method that uses multi-class classification with logistic regression. A set of utterances for each persona was prepared as training data, and logistic regression was trained so each utterances could be classified as being from any of the personas. At the time of inference, the probability that the target utterance is by the target persona is calculated by

	Utterance	Score
Myne	本さえあれば何もいらなく 思っているよ。(As long as I have books, I do not need anything.)	5
User	やっぱり読書をしていたら時 間を忘れちゃう？(Do you lose track of time when you read books?)	-
Myne	そんなことはしょっちゅうあ った。(That happened to me often.)	4
User	今度おすすめの本を紹介して くれる？(Can you recommend a book for me sometime?)	-
Myne	勿論！あなたにぴったりの本を 紹介するね！(Sure! I'll introduce you to the book that's right for you!)	5

Table 4: Example of a dialogue scenario and reference scores. The dialogue topic is reading, and Myne is a character from *Ascendance of a Bookworm* (ncode: n4830bu). The scores are the number of people of the five annotators who judged the utterance as Myne-like.

logistic regression, then used as a score.

In PSProb’s original configuration, training data were equalized for each character, so our experiment also used this configuration. As Table 2 shows, the smallest number of utterances among all characters was 52 for Keith. Therefore, for each of the 17 characters, we used 52 utterances, 50 for training data and 2 for development data, for 882 utterances in total<sup>4</sup>.

#### 4.3.2 Persona Term Saliency (PTSsal)

PTSsal (Miyazaki et al., 2021) is a method for assigning scores to terms in a given utterance. The method is based on TF-IDF and assigns higher scores to terms more frequently used by the target character and less by others. The term scores’ average is used as the estimated utterance score.

#### 4.3.3 ChatGPT

In recent years, the performance of large language models (LLMs), such as GPT-3 (Brown et al., 2020), has improved significantly on few-shot settings that use only a few examples. We used ChatGPT (gpt-3.5-turbo) (OpenAI, 2022) as a strong baseline in this experiment. ChatGPT outputs the target utterance’s likeliness score as an integer value from 0 to 5. For the utterance list, we used the top  $m$  utterances ranked by BM25 using the target utterances as a query, as with PersonaCLR. The examples were randomly selected from six utterances with a score of 0 to 5, one by one, from the target character’s evaluation data. Since examples’

<sup>4</sup>Previous study (Miyazaki et al., 2021) used 55 utterances for each character.

order affects results in the few-shot prompting (Gao et al., 2021; Jiang et al., 2020b; Liu et al., 2021), we shuffled the six examples’ order. The parameters given to the ChatGPT API were set to default settings except for temperature, which was set to 0.0 to generate deterministically. Appendix A.2 shows an example of the actual prompt and the hyperparameters of ChatGPT.

#### 4.4 BERTScore

BERTScore (Zhang et al., 2019) calculates similarity between texts by using vector representations obtained from pre-trained BERT. We calculated BERTScore between all pairs of the target utterance and reference utterances; the maximum BERTscore was used as the target utterance’s score. As reference utterances, we used the target character’s utterances in the NaroU (Table 2)

##### 4.4.1 MaxBLEU

MaxBLEU (Xu et al., 2018) is the maximum BLEU score between all pairs of the target utterance and reference utterances. We used SacreBLEU (Post, 2018) to compute the BLEU score.

##### 4.4.2 Persona-F1 (P-F1)

Rather than reference utterance-based, P-F1 (Jiang et al., 2020a) is a persona description-based evaluation measure that evaluates how well persona characteristics are expressed in an utterance. The higher the overlap between the non-stop word in the persona description and the utterance, the higher the P-F1 score.

#### 4.5 Implementation Details

We trained PersonaCLR using data from 94 out of the 100 novels in the NaroU dataset. We excluded the six novels shown in Table 2 to prevent any overlap between the characters in the training data in the NaroU and the evaluation data described in Section 5.1. We used 84 of the 94 novels as training data and ten as development data. We used Japanese RoBERTa<sub>large</sub><sup>5</sup> for PersonaCLR and BERTScore. We used the size of reference utterance set  $m$  to 20 in PersonaCLR and ChatGPT.

For PersonaCLR and PSProb, we conducted hyperparameter optimization. To find the optimal hyperparameters of PersonaCLR, a grid search was performed with temperature  $\tau$  as {0.01, 0.05, 0.1}, batch size as {16, 32, 64}, warmup steps as

<sup>5</sup><https://huggingface.co/nlp-waseda/roberta-large-japanese-with-auto-jumanpp>

Character	PersonaCLR	PSProb	PTSaI	ChatGPT	BERTScore	MaxBLEU	P-F1
Rimuru	0.201	0.015	0.005	0.124	<b>0.261</b>	0.099	-0.046
Veldora	<b>0.614</b>	0.478	0.386	0.340	0.337	0.450	0.067
Rudeus	<b>0.663</b>	0.334	0.426	0.314	0.369	0.392	0.098
Roxy	<b>0.644</b>	0.594	0.327	0.372	0.376	0.362	0.283
Sylphiette	<b>0.696</b>	0.550	0.540	0.345	0.594	0.240	0.309
Paul	<b>0.598</b>	0.311	0.290	0.023	0.288	0.311	0.191
Zenith	<b>0.527</b>	0.323	0.148	0.285	0.262	0.084	0.186
Subaru	<b>0.585</b>	0.447	0.236	0.222	0.120	0.165	0.218
Myne	<b>0.415</b>	0.147	0.150	0.145	0.181	0.042	0.181
Tuuli	<b>0.481</b>	0.332	0.308	0.401	0.220	0.202	0.143
Effa	<b>0.453</b>	0.295	0.236	0.351	0.207	0.068	0.117
Leon	<b>0.372</b>	0.273	0.197	0.120	0.129	0.148	0.245
Olivia	<b>0.726</b>	0.457	0.393	0.379	0.468	0.358	0.296
Angelica	<b>0.518</b>	0.290	0.374	0.116	0.311	0.172	0.179
Luxion	<b>0.641</b>	0.560	0.517	0.349	0.546	0.486	0.323
Catarina	<b>0.464</b>	0.328	0.174	0.277	0.293	0.254	0.144
Keith	<b>0.603</b>	0.476	0.471	0.420	0.387	0.209	0.297
Average	<b>0.541</b>	0.365	0.305	0.270	0.315	0.238	0.190

Table 5: Spearman’s rank correlation coefficients ( $r_s$ ) between the reference and estimated scores.

Character	PersonaCLR	PSProb	PTSaI	ChatGPT	BERTScore	MaxBLEU	P-F1
Rimuru	0.395	0.271	0.218	<b>0.440</b>	0.391	0.241	0.406
Veldora	<b>0.887</b>	0.630	0.417	0.544	0.431	0.633	0.582
Rudeus	<b>0.783</b>	0.541	0.605	0.461	0.494	0.531	0.471
Roxy	<b>0.697</b>	0.737	0.430	0.469	0.404	0.324	0.574
Sylphiette	<b>0.808</b>	0.557	0.579	0.453	0.552	0.485	0.543
Paul	<b>0.786</b>	0.609	0.514	0.270	0.362	0.466	0.581
Zenith	<b>0.748</b>	0.440	0.271	0.355	0.344	0.201	0.545
Subaru	<b>0.866</b>	0.556	0.347	0.374	0.341	0.328	0.438
Myne	<b>0.483</b>	0.246	0.229	0.325	0.224	0.212	0.529
Tuuli	<b>0.738</b>	0.355	0.296	0.477	0.285	0.344	0.587
Effa	<b>0.656</b>	0.412	0.302	0.608	0.326	0.205	0.601
Leon	<b>0.665</b>	0.404	0.359	0.261	0.267	0.272	0.482
Olivia	<b>0.908</b>	0.690	0.673	0.511	0.584	0.627	0.616
Angelica	<b>0.529</b>	0.339	0.348	0.288	0.294	0.468	0.585
Luxion	<b>0.758</b>	0.576	0.598	0.510	0.543	0.666	0.666
Catarina	<b>0.722</b>	0.582	0.487	0.637	0.501	0.560	0.590
Keith	<b>0.770</b>	0.498	0.510	0.545	0.448	0.479	0.634
Average	<b>0.718</b>	0.497	0.423	0.443	0.399	0.414	0.555

Table 6: AUPR for inappropriate utterance filtering

{100, 300, 500} and learning rate as  $\{1e^{-4}, 5e^{-5}, 1e^{-5}, 5e^{-6}, 1e^{-6}\}$ , respectively. As a result, we set the temperature  $\tau$  to 0.05, batch size to 64, warmup step to 300, and learning rate to  $1e^{-5}$ . During training, the loss of development data was calculated every 100 steps, and the model with the lowest loss was used for evaluation. For PSProb, the hyperparameter  $C$ , the inverse of regularization strength, was grid-searched on a logarithmic scale from 0.01 to 100, and  $C$  was set to 100.

#### 4.6 Evaluation Indices

We used two indices to examine PersonaCLR’s performance and this experiment’s comparative methods: Spearman’s rank correlation  $r_s$  and the area under the precision-recall curve (AUPR). We used Spearman’s rank correlation  $r_s$  to verify that PersonaCLR scores correlated with the human ratings

and AUPR to evaluate PersonaCLR’s performance in filtering inappropriate utterances. In calculating AUPR, we regarded utterances with a reference score of 0 or 1 as the detection target.

#### 4.7 Results

Experimental results based on Spearman’s rank correlation  $r_s$  between the manually assigned reference and estimated scores are shown in Table 5. For PersonaCLR and PSProb whose results depend on the random seed, training was performed three times with different seeds, and average values are shown in Table 5. Our PersonaCLR overperforms all other metrics, including ChatGPT. In PersonaCLR, 15 out of 17 characters showed a moderate correlation or higher ( $> 0.4$ ), and seven characters showed a strong correlation ( $> 0.6$ ). PSProb showed the proposed method’s second-best perfor-

mance. On more than half of the characters, PTSal showed inferior correlations to PSProb. ChatGPT and BERTScore results were uncorrelated ( $< 0.2$ ) for some characters, although in some cases, correlation exceeded the results of PSProb. ChatGPT and BERTScore showed higher performance than MaxBLEU, which also used only the target character’s reference utterances. However, ChatGPT and BERTScore were inferior to PersonaCLR and PSProb, which used utterances of several characters in training or calculating scores, thus suggesting that leveraging differences between personas is important in this task. PersonaCLR can learn this difference efficiently through contrastive learning, resulting in high performance. Overall, MaxBLEU, and P-F1 showed low performance, although correlations were observed for some characters.

Results of AUPR for inappropriate utterance filtering are shown in Table 6. PersonaCLR showed the best performance for 16 of 17 characters. One major difference from the results in Table 5 is that P-F1, which included character names and terms in its persona description, performed relatively better. In this experiment, most inappropriate utterances were created by non-target characters from other novels. Therefore, P-F1 effectively filtered out utterances that did not contain specific names or terms. ChatGPT showed relatively high performance for some characters, however, it was inferior to PersonaCLR, PSProb, and P-F1. Our results thus confirm PersonaCLR’s effectiveness.

## 5 Ablation Study

We conducted experiments using ablation models. The following two models were compared: A model that randomly samples from a pool of utterances instead of using BM25 to construct a set of reference utterances (w/o BM25), and a model that uses a single utterance as a reference that is the most similar to the target utterance by BM25 instead of the utterance set (w/ Single Ref.).

The results for each character in the ablation models using Spearman’s rank correlation coefficient are shown in Table 7 and those using AUPR are shown in 8. PersonaCLR shows the best performance for 14 of 17 characters in rank correlation coefficient, and 12 characters in AUPR. We also found that w/o BM25 outperformed PersonaCLR on several characters. This suggests that BM25 may have constructed an inappropriate reference utterance set for evaluating a given target utter-

Character	PersonaCLR	w/o BM25	w/ Single Ref.
Rimuru	<b>0.201</b>	0.178	0.119
Veldora	<b>0.614</b>	0.612	0.580
Rudeus	<b>0.663</b>	0.542	0.498
Roxy	<b>0.644</b>	0.465	0.534
Sylphiette	<b>0.696</b>	0.565	0.616
Paul	0.598	<b>0.694</b>	0.495
Zenith	<b>0.527</b>	<b>0.527</b>	0.470
Subaru	<b>0.585</b>	0.529	0.466
Myne	<b>0.415</b>	0.367	0.262
Tuuli	0.481	<b>0.563</b>	0.462
Effa	<b>0.453</b>	0.354	0.334
Leon	0.372	<b>0.418</b>	0.310
Olivia	<b>0.726</b>	0.621	0.696
Angelica	<b>0.518</b>	0.504	0.491
Luxion	<b>0.641</b>	0.546	0.609
Catarina	<b>0.464</b>	0.447	0.355
Keith	<b>0.603</b>	0.382	0.542
Average	<b>0.541</b>	0.489	0.461

Table 7: Spearman’s rank correlation coefficients ( $r_s$ ) for ablation models

Character	PersonaCLR	w/o BM25	w/ Single Ref.
Rimuru	0.395	<b>0.454</b>	0.325
Veldora	<b>0.887</b>	0.884	0.851
Rudeus	<b>0.783</b>	0.719	0.711
Roxy	0.697	<b>0.746</b>	0.551
Sylphiette	<b>0.808</b>	0.670	0.681
Paul	0.786	<b>0.808</b>	0.639
Zenith	<b>0.748</b>	<b>0.748</b>	0.628
Subaru	<b>0.866</b>	0.844	0.636
Myne	<b>0.483</b>	0.452	0.452
Tuuli	0.738	0.770	<b>0.780</b>
Effa	<b>0.656</b>	0.599	0.557
Leon	<b>0.665</b>	0.615	0.414
Olivia	<b>0.908</b>	0.813	0.866
Angelica	<b>0.529</b>	0.459	0.411
Luxion	0.758	0.683	<b>0.775</b>
Catarina	<b>0.722</b>	0.692	0.604
Keith	<b>0.770</b>	0.581	0.650
Average	<b>0.718</b>	0.679	0.620

Table 8: AUPR for ablation models

ance. Although we used the traditional ranking method BM25 in this paper, performance could be improved by improving the method of constructing reference utterances. With Single Ref., only two characters outperformed PersonaCLR in AUPR and zero in the correlation coefficient. These results indicate that employing a set of utterances rather than just a single utterance was important for appropriate evaluation.

The ablation study reconfirms PersonaCLR’s effectiveness.

## 6 Visualization

The embedding vector  $\mathbf{h}$  of the utterance obtained by the Transformer encoder in PersonaCLR reflects persona characteristics, and the same speaker’s utterances are closely placed in the vector space. We

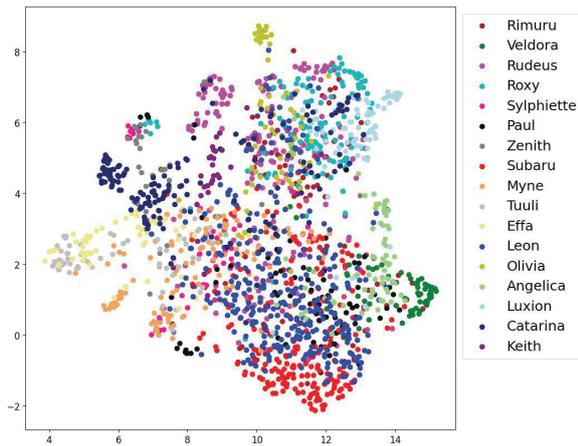


Figure 3: Utterance embedding visualization

visualized the vectors to confirm what speaker features were emphasized in the embedding process.

Figure 3 shows the embedding results of the 17 characters, that is, all utterances encoded by PersonaCLR and dimension reduction by Uniform Manifold Approximation and Projection (UMAP) (McInnes et al., 2018). We found that the manner of speaking determines the position. Characters who use polite language (e.g., Rudeus, Roxy, Olivia, Luxion, Keith) are placed in the upper area, those who do not use polite language (e.g., Paul, Subaru, Leon) in the lower, and those who use female language (e.g., Myne, Tuuli, Effa, Catarina) are placed on the left. We can also observe a tendency for characters in the same novel to be placed close. For example, Myne, Tuuli, and Effa, as well as Catarina and Keith, are located near each other due to common use in their utterances of novel-specific terms and character names. In contrast, Leon and Luxion, who is from the same novel, are positioned far apart, indicating that they are embedded with more emphasis on the manner of speaking than on being from the same novel.

## 7 Conclusion

We proposed a novel model for evaluating a given utterance’s intensity of persona characteristics and constructed the Naro Utterance dataset (NaroU) for training our model. The proposed model employs contrastive learning, and experimental results show that our model outperforms existing methods.

Future work includes constructing persona-aware dialogue systems by applying PersonaCLR and evaluating its performance experimentally. We also plan to extend PersonaCLR to be able to evaluate on a context-response basis rather than an utterance basis. This extension is expected to further

improve the response performance of the system.

## Acknowledgments

This research was supported by the NEDO project “Development of Interactive Story-Type Contents Creation Framework.”

## References

- 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.
- Annelen Brunner. 2013. Automatic recognition of speech, thought, and writing representation in german narrative texts. *Literary and linguistic computing*, 28(4):563–575.
- Jia-Xiang Chen, Zhen-Hua Ling, and Li-Rong Dai. 2019. A Chinese dataset for identifying speakers in novels. In *INTERSPEECH*, pages 1561–1565. Graz, Austria.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.
- Yue Chen, Tianwei He, Hongbin Zhou, Jia-Chen Gu, Heng Lu, and Zhen-Hua Ling. 2023. Symbolization, prompt, and classification: A framework for implicit speaker identification in novels. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3455–3467.
- Yue Chen, Zhen-Hua Ling, and Qing-Feng Liu. 2021. A neural-network-based approach to identifying speakers in novels. In *Interspeech*, pages 4114–4118.
- David Elson and Kathleen McKeown. 2010. Automatic attribution of quoted speech in literary narrative. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 24, pages 1013–1019.
- Hongchao Fang, Sicheng Wang, Meng Zhou, Jiayuan Ding, and Pengtao Xie. 2020. Cert: Contrastive self-supervised learning for language understanding. *arXiv preprint arXiv:2005.12766*.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot

- learners. 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 3816–3830.
- John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader. 2021. Declutr: Deep contrastive learning for unsupervised textual representations. 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 879–895.
- Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. 2021. Supervised contrastive learning for pre-trained language model fine-tuning. *Ninth International Conference on Learning Representation, ICLR 2021*.
- Seungju Han, Beomsu Kim, Jin Yong Yoo, Seokjun Seo, Sangbum Kim, Enkhbayar Erdenee, and Buru Chang. 2022. Meet your favorite character: Open-domain chatbot mimicking fictional characters with only a few utterances. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5114–5132.
- Hua He, Denilson Barbosa, and Grzegorz Kondrak. 2013. Identification of speakers in novels. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1312–1320.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738.
- Ryuichiro Higashinaka, Masahiro Mizukami, Hidetoshi Kawabata, Emi Yamaguchi, Noritake Adachi, and Junji Tomita. 2018. Role play-based question-answering by real users for building chatbots with consistent personalities. In *Proceedings of the 19th annual sigdial meeting on discourse and dialogue*, pages 264–272.
- Yuxiang Jia, Huayi Dou, Shuai Cao, and Hongying Zan. 2020. Speaker identification and its application to social network construction for Chinese novels. *International Journal of Asian Language Processing*, 30(04):2050018.
- Bin Jiang, Wanyue Zhou, Jingxu Yang, Chao Yang, Shihan Wang, and Liang Pang. 2020a. PEDNet: A persona enhanced dual alternating learning network for conversational response generation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4089–4099.
- Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. 2020b. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673.
- Hyunwoo Kim, Byeongchang Kim, and Gunhee Kim. 2020. Will i sound like me? improving persona consistency in dialogues through pragmatic self-consciousness. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 904–916.
- Minju Kim, Beong-woo Kwak, Youngwook Kim, Hong-in Lee, Seung-won Hwang, and Jinyoung Yeo. 2022. Dual task framework for improving persona-grounded dialogue dataset. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10912–10920.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and William B Dolan. 2016. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Ro{bert}a: A robustly optimized {bert} pretraining approach. In *Eighth International Conference on Learning Representations, ICLR 2020*. The International Conference on Learning Representations.
- Bodhisattwa Prasad Majumder, Harsh Jhamtani, Taylor Berg-Kirkpatrick, and Julian McAuley. 2020. Like hiking? you probably enjoy nature: Persona-grounded dialog with commonsense expansions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9194–9206.
- Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. 2018. Umap: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29).
- Koh Mitsuda, Ryuichiro Higashinaka, Hiroaki Sugiyama, Masahiro Mizukami, Tetsuya Kinebuchi, Ryuta Nakamura, Noritake Adachi, and Hidetoshi Kawabata. 2022. Fine-tuning a pre-trained transformer-based encoder-decoder model with user-generated question-answer pairs to realize character-like chatbots. In *Conversational AI for Natural Human-Centric Interaction: 12th International Workshop on Spoken Dialogue System Technology, IWSDS 2021, Singapore*, pages 277–290. Springer.

- Chiaki Miyazaki, Toru Hirano, Ryuichiro Higashinaka, and Yoshihiro Matsuo. 2016. Towards an entertaining natural language generation system: linguistic peculiarities of Japanese fictional characters. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 319–328.
- Chiaki Miyazaki, Saya Kanno, Makoto Yoda, Junya Ono, and Hiromi Wakaki. 2021. Fundamental exploration of evaluation metrics for persona characteristics of text utterances. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 178–189.
- Grace Muzny, Michael Fang, Angel Chang, and Dan Jurafsky. 2017. A two-stage sieve approach for quote attribution. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 460–470.
- OpenAI. 2022. Introducing chatgpt. <https://openai.com/blog/chatgpt>.
- Sean Papay and Sebastian Padó. 2020. RiQuA: A corpus of rich quotation annotation for English literary text. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages "835–841".
- 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.
- 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.
- Haoyu Song, Wei-Nan Zhang, Yiming Cui, Dong Wang, and Ting Liu. 2019. Exploiting persona information for diverse generation of conversational responses. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, pages 5190—5196.
- Yihong Tang, Bo Wang, Miao Fang, Dongming Zhao, Kun Huang, Ruifang He, and Yuexian Hou. 2023. Enhancing personalized dialogue generation with contrastive latent variables: Combining sparse and dense persona. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5456–5468.
- Zhen Xu, Nan Jiang, Bingquan Liu, Wenge Rong, Bowen Wu, Baoxun Wang, Zhuoran Wang, and Xiaolong Wang. 2018. Lsdsc: a large scale domain-specific conversational corpus for response generation with diversity oriented evaluation metrics. 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 2070–2080.
- Dian Yu, Ben Zhou, and Dong Yu. 2022. End-to-end Chinese speaker identification. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2274–2285.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTScore: Evaluating text generation with BERT. In *International Conference on Learning Representations*.
- Zhenyu Zhang, Yuming Zhao, Meng Chen, and Xiaodong He. 2022. Label anchored contrastive learning for language understanding. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1437–1449.
- Yinhe Zheng, Rongsheng Zhang, Minlie Huang, and Mao Xiaoxi. 2020. A pre-training based personalized dialogue generation model with persona-sparse data. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:9693–9700.

## A Details of Comparative Methods

### A.1 Persona Term Saliency (PTSsal)

PTSsal (Miyazaki et al., 2021) is a method for assigning scores to terms in a given utterance. The term scores’ average is used as the estimated utterance score for evaluation. PTSsal is obtained by the following equation.

$$PTSsal(t, p) = UttFreq(t, p) \cdot SpkrRarity(t) \quad (2)$$

$$UttFreq(t, p) = \frac{n(t, p)}{m(p)} \quad (3)$$

$$SpkrRarity(t) = \log \frac{|P|}{s(t)} \quad (4)$$

where  $n(t, p)$  is the number of utterances with term  $t$  in the monologue of persona  $p$  and  $m(p)$  is the total number of utterances in the monologue of persona  $p$ .  $s(t)$  is the number of personas that used term  $t$ , and  $|P|$  is the total number of personas. The  $UttFreq(t, p)$  becomes larger the more the target persona uses the term  $t$ , and  $SpkrRarity(t)$  is larger if a small number of personas other than the target persona uses the term  $t$ . In the experiment, we calculated  $SpkrRarity(t)$  using all the utterances in the Narou.

== Task ==  
Based on examples of a character’s utterances below, assign a rating from 0 to 5 to indicate the probability that the utterance was spoken by the character.

== Character’s utterance examples ==  
- Come on, it’s a tie-in. The knife and machete are sinking first, but what about you?  
- Oh, I was a shut-in!  
- Do not think less of me just because I’m a shut-in. My grip strength was over seventy kilograms. I can bench press up to 80 kilos!  
...  
(The rest is omitted. 20 utterances in total)

== Rating examples ==  
Utterance: Uh... I’m not good at horror...  
Rating: 2/5

Utterance: I do not watch many movies, but the only movie I watched recently was “One Piece.”  
Rating: 1/5

Utterance: Okay, I’ll buy it for you! I’ll get it for you, just wait there.  
Rating: 5/5

Utterance: My hobby is to learn all kinds of skills! Sewing, embroidery, figure skating, magic tricks... you name it!  
Rating: 3/5

Utterance: We get into trouble from time to time, but we live well together.  
Rating: 0/5

Utterance: Well, that settles it then. Yeah, it looks good on you.  
Rating: 4/5

**Utterance: Seriously, seriously, I’m soooo happy, looking forward to it!**  
Rating:

Table 9: Prompt format example for ChatGPT (originally written in Japanese)

## A.2 ChatGPT

An example of a prompt used in ChatGPT, which was used as a comparison method, is shown in Table 9. The target character to be evaluated in this example is Subaru from *Re: Life in a Different World from Zero*. The last utterance (bold font) is the evaluation target utterance. ChatGPT generates the score of the utterance after “Rating:.” The parameters given to the ChatGPT API were set to default settings except for temperature, which was set to 0.0 to generate deterministically.

## A.3 Persona-F1 (P-F1)

Rather than reference utterance-based, P-F1 (Jiang et al., 2020a) is a persona description-based evaluation measure that evaluates how well persona

本作の主人公。(The protagonist of this work.)  
4月1日生まれ。(Born on April 1. )  
黒の短髪、平凡な顔立ち、筋肉質のがっちりした体格の持ち主である少年。(He is a teenager with short black hair, an ordinary face, and a stocky, muscular build.)  
一般的な日本人よりも速く、目つきの悪さ (三白眼) が特徴である。(He is faster on his feet than the average Japanese, and he has bad eyesight (*sanpaku* eyes).)  
年齢は17歳 (開始時点)。(He is 17 years old (at the beginning of the story).) ...

Table 10: Persona description example of Subaru in *Re: Life in a Different World from Zero* (ncode: n2267be)

characteristics are expressed in an utterance. P-F1 is calculated as follows:

$$\text{PersonaF1} = \frac{2 \cdot \text{Recall} \cdot \text{Precision}}{\text{Recall} + \text{Precision}} \quad (5)$$

$$\text{Recall} = \frac{\max_{i \in [1, L]} |W_Y \cap d_i|}{|W_{d_i}|} \quad (6)$$

$$\text{Precision} = \frac{\max_{i \in [1, L]} |W_Y \cap d_i|}{|W_Y|} \quad (7)$$

where  $W_Y$  is a set of non-stop words in utterance  $Y$  and  $W_{d_i}$  is a set of non-stop words in sentence  $d_i$  in the persona description.

We used character descriptions in *Japanese Wikipedia* as the persona description. Examples of persona descriptions are shown in Table 10. The stop word was determined using the ja-stopword-remover library (version 0.2.4)<sup>6</sup>.

<sup>6</sup>[https://github.com/Pickerdot/ja\\_stopword\\_remover](https://github.com/Pickerdot/ja_stopword_remover)

# DiagESC: Dialogue Synthesis for Integrating Depression Diagnosis into Emotional Support Conversation

Seungyeon Seo<sup>1</sup>, Gary Geunbae Lee<sup>1,2</sup>

<sup>1</sup>Graduate School of Artificial Intelligence, POSTECH, Republic of Korea

<sup>2</sup>Department of Computer Science and Engineering, POSTECH, Republic of Korea

{ssy319, gblee}@postech.ac.kr

## Abstract

Dialogue systems for mental health care aim to provide appropriate support to individuals experiencing mental distress. While extensive research has been conducted to deliver adequate emotional support, existing studies cannot identify individuals who require professional medical intervention and cannot offer suitable guidance. We introduce the Diagnostic Emotional Support Conversation task for an advanced mental health management system. We develop the DESC dataset<sup>1</sup> to assess depression symptoms while maintaining user experience by utilizing task-specific utterance generation prompts and a strict filtering algorithm. Evaluations by professional psychological counselors indicate that DESC has a superior ability to diagnose depression than existing data. Additionally, conversational quality evaluation reveals that DESC maintains fluent, consistent, and coherent dialogues.

## 1 Introduction

As interest in preventing and treating mental illnesses like depression, anxiety disorders, and panic disorders grows, dialogue system studies on mental health care are gaining attention. Several studies have shown that chatbots can effectively manage the mental health of individuals, particularly in frontline settings, before seeking professional help (Denecke et al., 2021; Lim et al., 2022). These chatbots provide emotional empathy and assist in finding stability for those facing emotional, mental, and psychological distress. Mental health care also involves the early detection of illnesses. Although delayed treatment aggravates symptoms and requires more complex treatment, it is challenging for individuals to self-diagnose (Epstein et al., 2010). Therefore, detecting diseases during conversation is an important factor, and we focus on depression, a representative mental illness.

<sup>1</sup>Our dataset DESC is accessible at [github.com/seungyeon-seo/DiagESC](https://github.com/seungyeon-seo/DiagESC).

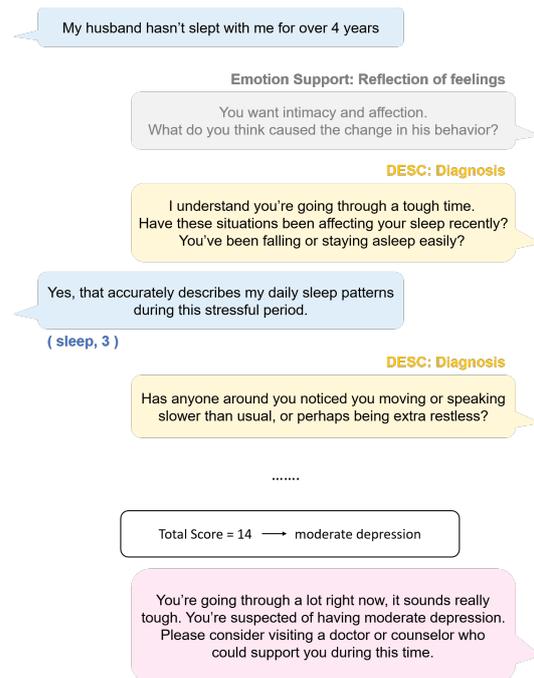


Figure 1: Part of an example conversation sample from DESC. The left is the seeker's, and the right is the supporter's utterance. We initiate a diagnostic conversation by inserting a diagnostic question (yellow) instead of a specific supporting emotion utterance (gray). At the end of the conversation, appropriate assistance (pink) is provided based on the severity of the depression.

Our research aims for an advanced conversation system to facilitate comprehensive mental health management. This system should provide extensive emotional support to individuals while simultaneously employing diagnostic questions to detect early signs of depression proactively.

To achieve this goal, we define a novel task, Diagnostic Emotional Support Conversation (DiagESC), based on Emotional Support Conversation (ESC) (Liu et al., 2021). ESC aims to support by helping reduce the seeker (user)'s mental stress. We also synthesize and release the dataset DESC for this task. We synthesize utterances to

ask questions about depression symptoms while maintaining a positive user experience. Additionally, task-specific strict filtering algorithms ensure data quality. Figure 1 shows part of the dialogue sample in DESC. It includes questions about depression symptoms and labels regarding symptom frequency, enabling assessing the severity of depression. Appropriate advice based on the severity of depression helps the individual receive help. Professional psychological counselors validate the diagnostic ability and conversational quality of DESC.

## 2 Related Work

### 2.1 Supportive Dialogue System

Recognizing emotions is essential for dialogue systems to respond appropriately to the user’s feelings. Emotion-tagged dialogue datasets such as DailyDialog (Li et al., 2017), Emotionlines (Hsu et al., 2018), and EmoContext (Chatterjee et al., 2019) have enhanced the conversation quality by enabling emotion-based response generation (Wei et al., 2019; Zandie and Mahoor, 2020; Ide and Kawahara, 2021). In particular, Lubis et al. (2019) shows that integrating emotional context in response generation can elicit positive emotions in users. The dataset EmpatheticDialogues (Rashkin et al., 2019) contains rich emotion labels and high-quality utterances that understand and empathize with users’ emotions, encouraging research on generating empathic responses (Ghosal et al., 2020; Majumder et al., 2020; Li et al., 2022).

To enable more effective emotional support, the ESC task (Liu et al., 2021) is defined by employing response strategies based on the Helping Skills Theory (Hill, 2009). ESC uses more sophisticated strategies, such as questioning and providing suggestions beyond empathy, to improve the users’ emotions and encourage them to overcome difficulties. Cheng et al. (2023) introduced persona generation into ESC and proposed Persona-Augmented Emotional Support (PAL), enabling the creation of responses tailored to an individual’s situation and characteristics.

However, understanding the situation and providing advice cannot fully help someone suffering from depression. Individuals with depression require professional counseling and medication rather than temporary emotional support. Research on supportive dialogue systems, such as ESC, demonstrates user encouragement capabil-

ities but cannot adequately address the needs of those with depression.

### 2.2 Depression Detection in Conversation

As with all diseases, early detection of depression is very important for efficient treatment. However, due to difficulties such as a lack of knowledge about the symptoms of depression, it is hard for patients to recognize that they are suffering from depression themselves (Epstein et al., 2010).

Against this background, depression detection research is being conducted to help with early treatment. Ringeval et al. (2019) proposed a classification task for whether a user has depression based on the audio and video features. They released the dataset DAIC-WOZ, which contains video recordings of clinical interviews designed to diagnose psychological disorders. The user participated in the conversation after completing a depression self-diagnosis questionnaire. DAIC-WOZ has significantly advanced depression detection research, contributing to numerous breakthroughs in the field (He and Cao, 2018; Haque et al., 2018; Low et al., 2020). We utilize DAIC-WOZ as a benchmark due to the absence of text-based depression diagnosis conversation datasets.

### 2.3 Dialogue Data Synthesis

Several methodologies have been proposed for the generation and augmentation of dialogue data to address the constraints associated with the time-intensive and costly data construction process (Lewis et al., 2017; Hou et al., 2018; Tang et al., 2019). With the emergence of the Large Language Model (LLM), the field of data synthesis has transitioned into a novel paradigm (Ding et al., 2024).

Kim et al. (2023); Bao et al. (2023) introduced a novel synthetic dialogue dataset derived from external sources. The data was refined using filtering techniques designed to ensure criteria such as commonsense knowledge, dialogue flow, and coherence. A method for synthesizing Dialogue State Tracking (DST) labeled conversation data from dialogue schemas and templates has shown comparable performance to human-annotated datasets in few-shot DST (Kulkarni et al., 2024). Kim et al. (2024); Li et al. (2024) generated the conversational dataset using task-specific prompting technology, and the test set is certified through humans.

Building on these advancements, we synthesize the DESC dataset by leveraging the fluent utterance generation capabilities of LLM, thereby con-

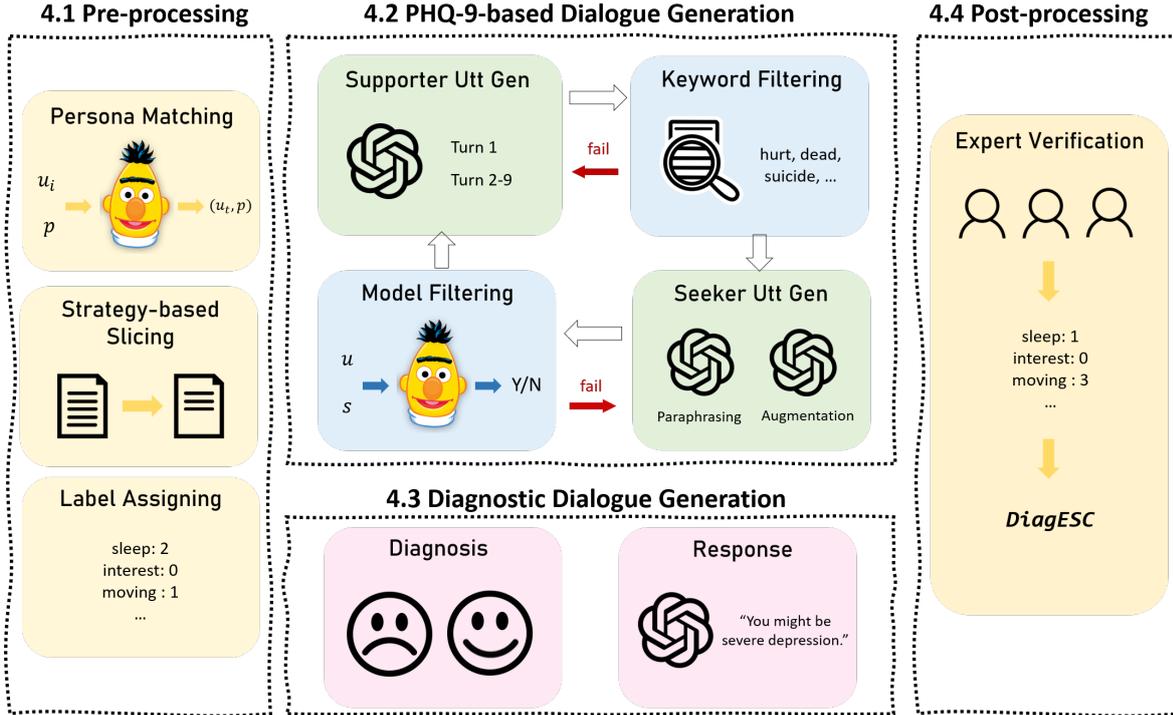


Figure 2: The overview of the DESC synthesis process.

tributing a novel resource for text-based depression detection research.

### 3 Problem Formulation

DiagESC consists of three sub-tasks—two modes response generation, persona generation, and diagnosis generation. The dual modes of response generation encompass emotional support and diagnostic responses. Persona generation is extracting characteristics based on the seeker’s previous utterance, as suggested by PAL (Cheng et al., 2023). Utilizing the previous persona as input can increase the user experience by facilitating reflection on an individual’s characteristics and serving as a form of memory when the dialogue history cannot include all utterances. Diagnosis generation, introduced for DiagESC, involves generating symptom and corresponding score pairs.

Equation 1 refers to DiagESC  $\mathcal{F}$  at turn  $t$  that generates response  $r_t = m \oplus s_t$ , persona  $p_t$ , and diagnosis pair  $d_t = \{\text{Symptom}, \text{Score}\}$  given the persona sentences  $P_t = \{p_1, p_2, \dots, p_t\}$  and dialogue history  $C_t = \{u_1, s_1, u_2, s_2, \dots, u_t\}$ .  $m \in \{\text{emotional support}, \text{diagnostic}\}$  denotes the mode for response.  $u_t$  and  $s_t$  represent  $t^{\text{th}}$  seeker (user) and supporter (system) utterance respectively.

$$\mathcal{F}(P_t, C_t) = (r_t, p_t, d_t) \quad (1)$$

The emotional support response uses existing ESC strategies, and the diagnostic response includes depression symptom questions and the notification diagnosis result.

## 4 Methodology

We synthesize the dataset for DiagESC named DESC through a four-step process, as illustrated in Figure 2. Initially, the source data undergoes pre-processing to align with the task requirements (Section 4.1). PHQ-9-based Dialogue Generation is for generating conversations that ask and answer about symptoms of depression (Section 4.2). Each task-specific prompt is based on the Patient Health Questionnaire-9 (PHQ-9) symptom item. The process includes filtering for reliability. The severity of depression is then calculated based on the answers obtained by the seeker. Section 4.3 is to inform the seeker of appropriate advice. Finally, to enhance reliability, expert verification is conducted on the validation and test datasets (Section 4.4).

### 4.1 Pre-processing

**Persona Matching** We utilize comprehensive annotations and high-quality supporting dialogue from the PESConv dataset of PAL (Cheng et al., 2023), containing persona sentences extracted from previous dialogue history. However, the persona

Symptom Item	Description
Interest	Little interest or pleasure in doing things
Depressed	Feeling down, depressed, or hopeless
Sleep	Trouble falling or staying asleep, or sleeping too much
Tired	Feeling tired or having little energy
Appetite	Poor appetite or overeating
Failure	Feeling you are a failure or have let yourself or your family down
Concentrating	Trouble concentrating on things, such as reading the newspaper or watching television
Moving	Moving or speaking so slowly that other people could have noticed. Or the opposite, being so fidgety or restless that you have been moving around a lot more than usual
Hurting	Thoughts that you would be better off dead, or of hurting yourself

Table 1: The symptoms and descriptions of PHQ-9

sentences do not align exactly with the seeker’s utterance for each turn.

We employ the BERT<sup>2</sup> (Kenton and Toutanova, 2019) model to obtain the embeddings for all persona sentences  $p$  and seeker utterances  $u$ . Then, we compute cosine similarities between each persona sentence embedding and every utterance embedding. Each persona sentence has a higher cosine similarity to the utterance from which it is derived than to other utterances. We reassign all of the persona sentences using the following equation.

$$\hat{t}_i = \operatorname{argmax}_{t \in \{1, \dots, T\}} \frac{E(p_i) \cdot E(u_t)}{|E(p_i)| |E(u_t)|} \quad (2)$$

where  $\hat{t}_i$  represents the matched turn number for the  $i$ -th persona sentence  $p_i$ .  $u_t$  denotes the utterance at the  $t$ -th turn and the function  $E(\cdot)$  refers to compute BERT embedding. Equation 2 ensures to align each persona sentence with its derived utterance.

**Strategy-based Slicing** Determining the appropriate moment to begin diagnostic questions is challenging. It is crucial to consider that abruptly interrupting the flow of conversation may negatively impact the user’s emotional state. Fortunately, ESC dataset has rich annotations, tagging each utterance with its corresponding strategy. The most suitable time for presenting diagnostic questions has been empirically determined to use specific strategies, namely *Restatement or Paraphrasing*, *Reflection of Feeling*, *Self-disclosure*, and *Affirmation and Reassurance*. Figure 1 is an example of using a diagnostic question (yellow) instead of a reflection response (gray). This rule enables a smooth and contextually appropriate transition into diagnostic questioning.

Furthermore, we only utilize truncated data when at least two persona sentences have been

gathered to ensure that diagnostic questions are only posed after comprehensively understanding the seeker’s persona. This criterion helps that sufficient contextual background is considered before diagnostic engagement.

**Label Assigning** To achieve an even distribution of the final severity level within the generated conversational data, the PHQ-9 labels are pre-assigned. The next step uses predefined labels to generate utterances.

## 4.2 PHQ-9-based Dialogue Generation

We utilize the Patient Health Questionnaire-9 (PHQ-9), a widely used medical tool for self-assessment of depression (Kroenke et al., 2001), as the basis for the depression diagnostic questions. PHQ-9 aims to quantify the frequency of nine depressive symptoms listed in Table 1 on a scale ranging from 0 to 3, with the options *Not at all*, *Several days*, *More than half the days*, and *Nearly every day*. The aggregated score of all items is used to diagnose depression and assess its severity, categorized as Minimal (0-4), Mild (5-9), Moderate (10-14), Moderately severe (15-19), and Severe (20-27).

### 4.2.1 Supporter Utterance Generation

We develop two types of prompts to generate supporter utterances for the first and subsequent turns. In the initial turn, it is essential to formulate questions with caution to maintain a positive user experience. For the subsequent turns, which involve further diagnostic questioning, it becomes essential to comprehend and empathize with the seeker’s responses to the preceding questions.

The both prompts involve the three-step Chain-of-Thought technique (Wei et al., 2022; Kim et al., 2024). In the first-turn supporter utterance generation prompt, the steps consist of *Selection*, *Plan-*

<sup>2</sup><https://huggingface.co/google-bert/bert-base-uncased>

### Prompt Content

You are an emotional supporter. You have to ask about the frequency of depression symptoms without compromising the emotions of the user suspected of having depression.

**Depression Symptoms** You should ask how ‘often’ a symptom has occurred over the past two weeks. Use one of the following symptoms. Be careful not to distort the medical meaning. (*symptoms*)

**Task Description** The task proceeds in three stages: Selection, Planning, and Response Generation. The first step, Selection, is to select which of the given persona sentences and dialog history to use in the response generation and what symptoms to ask about. Information that can improve the user experience must be extracted. The second step, Planning, is planning how to use the selected information. You must explain how you will use the information you have selected and why you have selected that information. The final step, Response Generation, uses the selected information to naturally ask the user about depression symptoms. Consistency with persona and history must not be broken. Questions must be asked carefully so that the user does not feel that the question is sudden. Be careful not to ask hard as if you were being interrogated. The generated response must be no more than 25 words.

**Example** (*examples*)

Table 2: The prompt used to generate the initial supporter utterance of inquiring about PHQ-9 symptoms.

*ning*, and *Response Generation*. Table 2 provides detailed instructions for these steps. *Selection* and *Planning* focus on the seeker’s persona and the dialogue history.

Because analyzing the previous answer about the symptom is more critical for subsequent turns, we replaced the *Selection* with the *Analysis*. The steps help to analyze the seeker’s response and generate a response accordingly. Detailed instructions can be found in the Appendix B.

#### 4.2.2 Seeker Utterance Generation

The construction of the seeker’s utterance necessitates including one of the four designated responses from PHQ-9 (Not at all, Several days, More than half the days, Nearly every day). We enhance and utilize the template-based utterance generation method (Kulkarni et al., 2024). After rephrasing

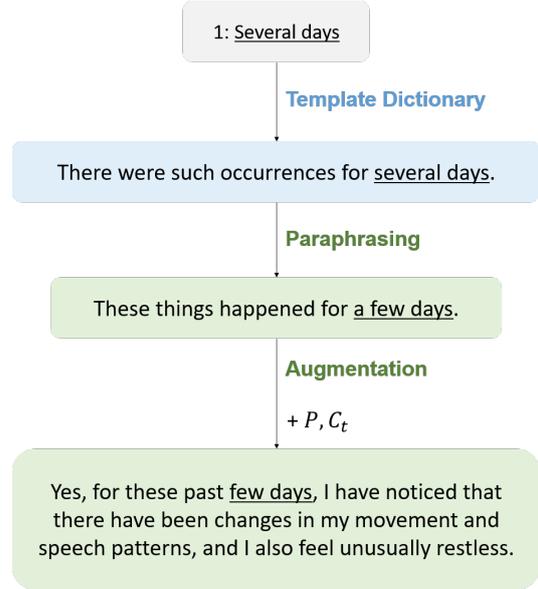


Figure 3: The overview of seeker utterance generation

Score	Type	Content
0	PHQ	Not at all
	Temp.	There is not at all much of that.
1	PHQ	Several days
	Temp.	There were such occurrences for several days.
2	PHQ	More than half the days
	Temp.	That happened for more than half the days.
3	PHQ	Nearly every day
	Temp.	It was like that almost every day.

Table 3: Templates by score used for the Paraphrasing stage in generating seeker utterances.

the templates, we augment the utterances to match the user’s characteristics better, personalizing them to align more closely with specific traits and preferences, providing a more tailored and engaging conversational experience.

Figure 3 shows the process by which the template is transformed into the final utterance via the stages of Paraphrasing and Augmentation using LLM. Initially, we establish a foundational template that directly correlates with a response option from the PHQ-9, as illustrated in Table 3.

Paraphrasing model  $M_{Para}$  aims to diversify responses while preserving symptom frequency information. Augmentation model  $M_{Aug}$  represents the process of generating the final utterance  $u_t$  to align with the seeker’s specific situation and persona, particularly considering the supporter’s last utterance and the overall conversational flow. As can be seen in Equation 3,  $M_{Para}$  requires a template of the corresponding score  $s$  as input and  $M_{Aug}$  demands a paraphrased output, the persona sentences  $P_t$ , and

the conversational history  $C_t$ .

$$u_t = \mathbf{M}_{\text{Aug}}(P_t, C_t, \mathbf{M}_{\text{Para}}(\text{template}_s)) \quad (3)$$

### 4.2.3 Filtering

If the LLM misrepresents the medical interpretation of PHQ-9, it will adversely affect the quality of the diagnosis function. We implement strict filtering algorithms to prevent the hallucination of LLM and ensure the reliability of diagnostic conversations.

**Keyword Filtering** To maintain the integrity of symptom representation in LLM-generated supporter utterances, we pre-define specific keywords associated with each symptom. If each utterance includes no pre-defined keywords, we discard it and regenerate a new one. This process continues until the generated utterance appropriately incorporates the necessary keywords, ensuring accurate and consistent symptom representation in the dialogue.

**Model Filtering** It is essential to preserve symptom frequency information in the seeker’s utterances. Therefore, we train the classification model  $M_c$ , using 256 manually verified utterances. If the predicted label from model  $M_c$  differs from the template label, or if the confidence is below a threshold  $t$ , the utterance is regenerated.

## 4.3 Diagnostic Dialogue Generation

The final goal of DiagESC’s diagnostic ability is to estimate the seeker’s mental health status and provide appropriate assistance. The severity level of depression is determined by summing the scores of all nine items obtained from the Diagnosis task and then generating an appropriate response as shown in the final utterance in Figure 1. To achieve this goal, we design the prompt in Table 4 to generate an utterance based on the seeker’s persona and a diagnosed depression severity level. To enhance the naturalness of the conversation, we incorporate a predefined turn expressing gratitude for the honest response before the diagnostic response.

## 4.4 Post-processing

Despite applying strict task-specific filtering protocols, the potential for inaccuracies remains. To ensure the reliability of the PHQ-9 labels, Expert Filtering is conducted on the validation and test sets of DESC. Three psychologists, who are native English speakers or bilingual and have over four years of professional experience<sup>3</sup>, assessed scores

<sup>3</sup>We hired psychologists through <https://www.upwork.com>

---

### Prompt Content

---

You are emotional support. You have provided counseling to the user about the concerns and even completed questions about depression symptoms. Generate an utterance that concludes the counseling by referring to the depression diagnosis results and the user’s persona. If the severity of depression is high, you should be recommended to see a hospital or counselor. Please generate the utterance friendly conversational style and generated utterance must be no more than 30 words.

### Example (examples)

---

Table 4: The prompt used to generate utterance for notifying diagnosis result.

for each symptom. The seeker utterances are then re-labeled to the mode value of the three scores.

## 5 Experiments

### 5.1 Diagnostic Ability Evaluation

The DAIC-WOZ dataset, used as a baseline, comprises clinical dialogues in video and audio features with PHQ-8 (Kroenke et al., 2009) labels. The PHQ-8 is a modified version of the PHQ-9, excluding the items related to suicide, and performs just as well as the PHQ-9 in diagnosing depression. Although the modality is different from ours, due to the absence of conversation data explicitly labeled for depression, we use transcripts of DAIC-WOZ. We randomly sample four dialogues from DAIC-WOZ for each severity level. Employing the same methodology with expert labeling described in Section 4.4, three psychological counselors evaluate scores for the PHQ-8 items.

The Quadratic Weighted Kappa (QWK) score is a metric that evaluates the agreement between two predictions, offering advantages by acknowledging both exact and partial alignment in assessments (Cohen, 1968). The QWK score ranges from -1 to 1. A score closer to -1 indicates that the predictions are nearly opposite. A score near 0 reflects randomness, implying no consistent agreement between the predictions. If a score approaches 1, the predictions are almost identical and have a high level of agreement. As it is suitable for medical fields where symptoms can be interpreted slightly differently depending on the individual doctor (Yoshida et al., 2015; Nirthika et al., 2020; Chivinge et al., 2022), the QWK score is widely used in disease diagnosis. Therefore, we adopt the QWK score as

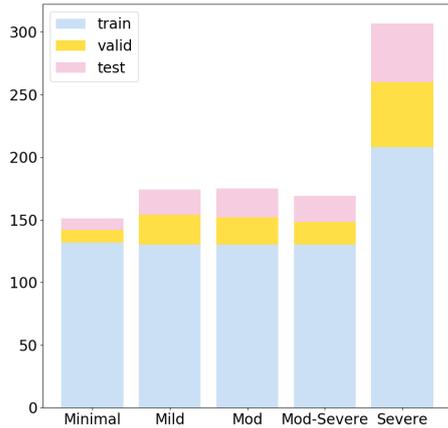


Figure 4: Distribution of depression severity labels in DESC. Minimal (0-4), Mild (5-9), Moderate (10-14), Moderately severe (15-19), and Severe (20-27).

the principal metric to evaluate diagnostic ability.

## 5.2 Conversational Quality Evaluation

To evaluate the quality of the conversational data, we sample 10 dialogues for each severity level from the DESC validation and test sets. We then requested the same evaluators with diagnostic ability evaluation to rate the following three items on a scale from 1 (Very Poor) to 5 (Excellent).

- **Fluency** evaluates the grammatical correctness, naturalness, and smoothness of the dialog.
- **Consistency** assesses how well the dialogue maintains a consistent user persona throughout the interaction. This involves the user’s interests and personality traits.
- **Coherence** measures how contextually appropriate the responses are, considering the previous dialogue turns and the overall context of the conversation.

## 5.3 Automatic Validation in Baseline

We establish the baseline models by tuning a small LLM utilizing Low-Rank Adaptation (LoRA) (Hu et al., 2022) to evaluate the operation of DESC in each model. We train the LoRA adapters on DESC for the three sub-tasks: response, persona, and diagnosis generation. Additionally, performance in a multi-adapter setting is measured to evaluate the multi-tasking capabilities of the system. In a single-task setting, the inputs for the next turn are the true labels of other tasks. However, the inferred results

from the previous turn are used as input for all tasks in a multi-task setting. We systematically provide symptom item sequences to ensure consistency and effectiveness in the diagnostic process.

## 5.4 Implementation Details

We use GPT-4 as LLM to generate utterances of DESC. For Model Filtering of seeker utterance, we adopt RoBERTa<sup>4</sup> (Liu et al., 2019) and train the classification model  $M_c$  for 5 epochs. The labels predicted by the fine-tuned model are utilized for filtering purposes, with the threshold  $t = 0.7$ . Llama2<sup>5</sup> (Touvron et al., 2023) is used as a baseline model, and the adapters are trained with the train set for 5 epochs on 4 NVIDIA A6000 GPUs, and the final model with the lowest validation loss was selected. We employ AdamW with a learning rate of  $5e-5$  and a linear scheduler.

## 6 Results and Analysis

### 6.1 Basic Statistics of DESC

The DESC comprises 976 dialogues, including 730 train, 126 validation, and 120 test samples. Each dialogue has an average of 42 turns, with the maximum number of turns per dialogue being 111 and the minimum being 24. Figure 4 illustrates the distribution of dialogue samples across five levels of depression severity. The Severe level has more samples than the other levels because it covers a wider range of scores.

### 6.2 Diagnosis Ability

According to the results presented in Table 5, DESC achieves a notably high average QWK score of 0.70 compared to baseline. In contrast, DAIC-WOZ obtains low scores, with a certain item showing negative values. The result indicates a substantial challenge in predicting the frequency of a seeker’s depression symptoms solely from conversational history with an agent in the dataset. This comparison may be considered unfair because the DAIC-WOZ does not include questions about all the symptoms of the PHQ-8.

The most important result is the final depression diagnostic capability of each dataset, as presented in Table 6. The depression severity is classified into five levels—minimal, mild, moderate, moderately severe, and severe—based on the cumulative scores of the assessed items. Scores exceeding 10 points,

<sup>4</sup>FacebookAI/roberta-base

<sup>5</sup>meta-llama/Llama-2-7b-chat-hf

Dataset	Interest	Depressed	Sleep	Tired	Appetite	Failure	Concentrating	Moving	Hurting	Avg
DAIC-WOZ	0.16	0.39	0.18	0.15	0.03	0.03	-0.13	0.04	-	0.11
DESC	0.44	0.69	0.64	0.70	0.80	0.80	0.64	0.78	0.81	<b>0.70</b>

Table 5: Average QWK Scores of each dataset against expert annotations for each symptom in PHQ-8 and PHQ-9.

Dataset	Level Acc	Depression			
		Acc	Precision	Recall	F1
DAIC-WOZ	0.45	0.70	0.50	1.00	0.67
DESC	0.71	0.89	1.00	0.86	0.92

Table 6: Accuracy of predicting depression severity level and accuracy, precision, recall, and f1 score of estimating depression diagnosis.

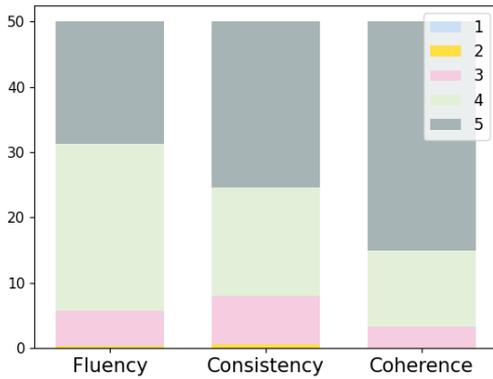


Figure 5: Distribution of evaluated scores for DESC’s fluency, consistency, and coherence.

classified as moderate or higher, are considered depression.

The Level Acc indicates the accuracy of predicting severity levels, with DESC showing 0.26 higher performance than DAIC-WOZ. In depression diagnosis, the superior accuracy and F1 score of DESC, compared to the baseline, demonstrate its robustness and effectiveness. The results suggest that our PHQ-9-based data generation process ensures reliable diagnostic capabilities.

### 6.3 Conversation Quality

Figure 5 shows the distribution of obtained scores in the human evaluation performed to evaluate conversation quality. The average scores are 4.25 for fluency, 4.33 for consistency, and 4.63 for coherence. Most samples received scores of 3 or higher across all items, indicating that the DESC is consistent and comprises high-quality conversations without contextual awkwardness. Notably, the high coherence score suggests that the diagnostic questions generated through strategy-based slicing and first-turn supporter prompt configuration help the seeker perceive them as natural and non-abrupt.

	Response		Persona	Diagnosis
	Mode	BLEU	BLEU	Acc
Single	0.83	31.08	34.03	0.78
Multi	0.83	30.78	34.65	0.77

Table 7: The performance on baseline models for single-task and multi-task settings.

### 6.4 Automatic Validation

Table 7 shows the baseline performance of DESC. Mode indicates the prediction accuracy of the response mode, divided into emotional support and diagnosis. It shows equal performance of 0.83 in both single-task and multi-task settings. Generating response and persona sentences achieve high BLEU scores, all above 30. Diagnosis accuracy measures the prediction of each symptom and its corresponding score. Across all metrics, single-task and multi-task settings demonstrate similar performance.

## 7 Conclusion

This work proposes the DiagESC task for a comprehensive mental health care dialogue system that goes beyond the limitations of supportive dialogue systems that do not detect mental risk. DiagESC contributes to emotional support and early detection of depression, an important part of mental health. We have released the novel dataset DESC by synthesizing diagnostic conversations based on a depression self-diagnosis questionnaire with emotional support data. Task-specific prompts and strict filtering protocols facilitate questions about depression symptoms while ensuring continued user engagement. Evaluation by a psychological counseling expert proves that DESC has superior diagnostic performance and conversational quality. We hope that DiagESC will contribute significantly to developing more effective and supportive dialogue systems in mental health care. Moreover, the release of the DESC dataset provides a valuable resource for the research community, encouraging further advancements and innovations in this critical area.

## Limitation

The research aims to identify depression signs during conversations with the user, subsequently notifying them of potential risks. It is imperative to note that the diagnostic outcomes derived from the proposed dataset and model are intended solely for guidance. An accurate and definitive diagnosis should be ascertained through consultation with a medical professional.

## Acknowledgements

This work was supported by Smart Health-Care Program(www.kipot.or.kr) funded by the Korean National Police Agency(KNPA, Korea) [Project Name: Development of an Intelligent Big Data Integrated Platform for Police Officers' Personalized Healthcare / Project Number: 220222M01] This research was supported by the MSIT(Ministry of Science and ICT), Korea, under the ITRC(Information Technology Research Center) support program(IITP-2024-RS-2024-00437866) supervised by the IITP(Institute for Information Communications Technology Planning Evaluation)

## References

- Jianzhu Bao, Rui Wang, Yasheng Wang, Aixin Sun, Yitong Li, Fei Mi, and Ruifeng Xu. 2023. [A synthetic data generation framework for grounded dialogues](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10866–10882, Toronto, Canada. Association for Computational Linguistics.
- Ankush Chatterjee, Kedhar Nath Narahari, Meghana Joshi, and Puneet Agrawal. 2019. Semeval-2019 task 3: Emocontext contextual emotion detection in text. In *Proceedings of the 13th international workshop on semantic evaluation*, pages 39–48.
- Jiale Cheng, Sahand Sabour, Hao Sun, Zhuang Chen, and Minlie Huang. 2023. Pal: Persona-augmented emotional support conversation generation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 535–554.
- Lincoln Chivinge, Leslie Kudzai Nyandoro, and Kudakwashe Zvarevashe. 2022. [Quadratic weighted kappa score exploration in diabetic retinopathy severity classification using efficientnet](#). In *2022 1st Zimbabwe Conference of Information and Communication Technologies (ZCICT)*, pages 1–9.
- Jacob Cohen. 1968. Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit. *Psychological bulletin*, 70(4):213.
- Kerstin Denecke, Alaa Abd-Alrazaq, and Mowafa Househ. 2021. Artificial intelligence for chatbots in mental health: opportunities and challenges. *Multiple perspectives on artificial intelligence in health-care: Opportunities and challenges*, pages 115–128.
- Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze Luo, Xinze Li, Guizhen Chen, Wenhan Xia, Junjie Hu, Anh Tuan Luu, and Shafiq Joty. 2024. Data augmentation using llms: Data perspectives, learning paradigms and challenges. *arXiv preprint arXiv:2403.02990*.
- Ronald M Epstein, Paul R Duberstein, Mitchell D Feldman, Aaron B Rochlen, Robert A Bell, Richard L Kravitz, Camille Cipri, Jennifer D Becker, Patricia M Bamonti, and Debora A Paterniti. 2010. “i didn’t know what was wrong:” how people with undiagnosed depression recognize, name and explain their distress. *Journal of general internal medicine*, 25:954–961.
- Deepanway Ghosal, Navonil Majumder, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020. Cosmic: Commonsense knowledge for emotion identification in conversations. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2470–2481.
- Albert Haque, Michelle Guo, Adam S Miner, and Li Fei-Fei. 2018. Measuring depression symptom severity from spoken language and 3d facial expressions. *arXiv preprint arXiv:1811.08592*.
- Lang He and Cui Cao. 2018. Automated depression analysis using convolutional neural networks from speech. *Journal of biomedical informatics*, 83:103–111.
- Clara E. Hill. 2009. *Helping Skills: Facilitating, Exploration, Insight, and Action*, 3 edition. American Psychological Association.
- Yutai Hou, Yijia Liu, Wanxiang Che, and Ting Liu. 2018. Sequence-to-sequence data augmentation for dialogue language understanding. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1234–1245.
- Chao-Chun Hsu, Sheng-Yeh Chen, Chuan-Chun Kuo, Ting-Hao Huang, and Lun-Wei Ku. 2018. Emotionlines: An emotion corpus of multi-party conversations. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- 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*.
- Tatsuya Ide and Daisuke Kawahara. 2021. Multi-task learning of generation and classification for emotion-aware dialogue response generation. In *Proceedings*

- of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop, pages 119–125.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, et al. 2023. Soda: Million-scale dialogue distillation with social commonsense contextualization. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Yumin Kim, Heejae Suh, Mingi Kim, Dongyeon Won, and Hwanhee Lee. 2024. Kocosa: Korean context-aware sarcasm detection dataset. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9890–9904.
- Kurt Kroenke, Robert L Spitzer, and Janet BW Williams. 2001. The phq-9: validity of a brief depression severity measure. *Journal of general internal medicine*, 16(9):606–613.
- Kurt Kroenke, Tara W Strine, Robert L Spitzer, Janet BW Williams, Joyce T Berry, and Ali H Mokdad. 2009. The phq-8 as a measure of current depression in the general population. *Journal of affective disorders*, 114(1-3):163–173.
- Atharva Kulkarni, Bo-Hsiang Tseng, Joel Moniz, Dhivya Piraviperumal, Hong Yu, and Shruti Bhargava. 2024. SynthDST: Synthetic data is all you need for few-shot dialog state tracking. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1988–2001, St. Julian’s, Malta. Association for Computational Linguistics.
- Mike Lewis, Denis Yarats, Yann Dauphin, Devi Parikh, and Dhruv Batra. 2017. Deal or no deal? end-to-end learning of negotiation dialogues. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2443–2453.
- Jinpeng Li, Zekai Zhang, Quan Tu, Xin Cheng, Dongyan Zhao, and Rui Yan. 2024. Stylechat: Learning recitation-augmented memory in llms for stylized dialogue generation. *arXiv preprint arXiv:2403.11439*.
- Qintong Li, Piji Li, Zhaochun Ren, Pengjie Ren, and Zhumin Chen. 2022. Knowledge bridging for empathetic dialogue generation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):10993–11001.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. Dailydialog: A manually labelled multi-turn dialogue dataset. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 986–995.
- Shi Min Lim, Chyi Wey Claudine Shiau, Ling Jie Cheng, and Ying Lau. 2022. Chatbot-delivered psychotherapy for adults with depressive and anxiety symptoms: A systematic review and meta-regression. *Behavior Therapy*, 53(2):334–347.
- Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. Towards emotional support dialog systems. 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 3469–3483.
- 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*.
- Daniel M Low, Kate H Bentley, and Satrajit S Ghosh. 2020. Automated assessment of psychiatric disorders using speech: A systematic review. *Laryngoscope investigative otolaryngology*, 5(1):96–116.
- Nurul Lubis, Sakriani Sakti, Koichiro Yoshino, and Satoshi Nakamura. 2019. Positive emotion elicitation in chat-based dialogue systems. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(4):866–877.
- Navonil Majumder, Pengfei Hong, Shanshan Peng, Jiankun Lu, Deepanway Ghosal, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020. Mime: Mimicking emotions for empathetic response generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8968–8979.
- Rajendran Nirthika, Siyamalan Manivannan, and Amirthalingam Ramanan. 2020. Loss functions for optimizing kappa as the evaluation measure for classifying diabetic retinopathy and prostate cancer images. In *2020 IEEE 15th international conference on industrial and information systems (ICIIS)*, pages 144–149. IEEE.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381.
- Fabien Ringeval, Björn Schuller, Michel Valstar, Nicholas Cummins, Roddy Cowie, Leili Tavabi, Maximilian Schmitt, Sina Alisamir, Shahin Amiriparian, Eva-Maria Messner, et al. 2019. Avec 2019 workshop and challenge: state-of-mind, detecting depression with ai, and cross-cultural affect recognition. In *Proceedings of the 9th International on Audio/visual Emotion Challenge and Workshop*, pages 3–12.

Jianheng Tang, Tiancheng Zhao, Chenyan Xiong, Xiaodan Liang, Eric Xing, and Zhiting Hu. 2019. Target-guided open-domain conversation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.

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*.

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.

Wei Wei, Jiayi Liu, Xianling Mao, Guibing Guo, Feida Zhu, Pan Zhou, and Yuchong Hu. 2019. Emotion-aware chat machine: Automatic emotional response generation for human-like emotional interaction. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pages 1401–1410.

Kanako Yoshida, Rebecca J Barr, Sandro Galea-Soler, Richard M Aspden, David M Reid, and Jennifer S Gregory. 2015. Reproducibility and diagnostic accuracy of kellgren-lawrence grading for osteoarthritis using radiographs and dual-energy x-ray absorptiometry images. *Journal of Clinical Densitometry*, 18(2):239–244.

Rohola Zandie and Mohammad H Mahoor. 2020. Empransfo: A multi-head transformer architecture for creating empathetic dialog systems. In *The thirty-third international flairs conference*.

## A Dialogue Example

Figure 6 illustrates the utterance configuration of the DESC data sample. During emotional support conversation, the supporter initiates the PHQ-9-based dialogue to inquire about symptoms. After identifying all symptoms, the supporter provides the user with appropriate advice, such as a recommendation to visit a hospital based on the diagnostic results. To enhance the natural flow of the conversation, we insert the predefined turns between symptom inquiries and result notifications. We select these turns from 23 supporter and 16 seeker utterance candidates.

## B Detailed Prompt Instructions

Table 8 is the prompt for supporter utterance generation. In contrast to the utterance generation of turn 1, this stage focuses on analyzing the prior

Emotional Support	
Supporter	Hi, how are you? What would you like to discuss today?
Seeker	hi, I am not feeling too well.
Supporter	How long have you felt this way?
Seeker	I have been feeling this low for close to a month now
...	
Supporter	You had an affair because you were dissatisfied with your marriage. Why do you think that was?
Seeker	my husband hasnt slept with me for over 4 years
PHQ-9-based Dialogue	
Supporter	I understand you're going through a tough time. Have these situations been affecting your sleep recently? You've been falling or staying asleep easily?
Seeker	Yes, that accurately describes my daily sleep patterns during this stressful period.
Supporter	Has anyone around you noticed you moving or speaking slower than usual, or perhaps being extra restless?
Seeker	Indeed, those around me have noticed my unusual behavior for quite a few days now.
Supporter	Is it common for you to think about no longer being around or inflicting harm upon yourself?
Seeker	In reality, I don't regularly have thoughts of that nature.
...	
Supporter	I understand. Have you been having trouble focusing on things like reading or watching TV often?
Seeker	Yes, over the past few days, I've indeed found it increasingly difficult to focus on things, like reading or watching television.
Pre-defined Turn	
Supporter	Thank you for being so direct in your answer.
Seeker	Thank you for your help; it means a lot in these tough times.
Notification	
Supporter	You're going through a lot right now, it sounds really tough. You're suspected of having moderate depression. Please consider visiting a doctor or counselor who could support you during this time.

Figure 6: Part of an example conversation sample of DESC.

seeker response. Table 9 and 10 are the prompts for seeker utterance generation.

## C Keyword Filtering

The keyword filtering process ensures that the PHQ-9 maintains its medical meaning. Table 11 shows detailed keywords for each symptom item. The generated utterance must contain at least one of the keywords.

## D Distribution of DESC

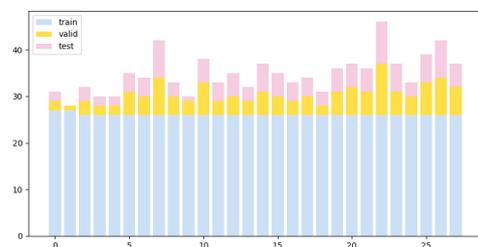


Figure 7: Distribution of aggregated score in DESC.

Each dialogue sample has score labels for the PHQ-9 items. The sum of all nine scores determines the severity of depression. Figure 7 is the distribution of aggregated scores.

---

**Prompt Content**

---

(same with turn 1)

**Depression Symptoms** You should ask how ‘often’ a symptom has occurred over the past two weeks. Symptoms are given and the frequency of the symptoms should be naturally asked of the user. The meaning of a given symptom should never be changed.

**Task Description** The task proceeds in three stages: Analysis, Planning, and Response Generation. The first step, Analysis, is to determine the user’s status through the user’s previous responses. The second step, Planning, is planning how to use the status information to support the user’s emotions and ask about the frequency of a different given symptom. The final step, Response Generation, is to ask the user about the symptom according to plan. Question must be asked carefully so that the user does not feel that the question is sudden. Be ... (same with turn 1)

**Example** (examples)

---

Table 8: The prompt used to generate the subsequent turn supporter utterance of inquiring about PHQ-9 symptoms.

## E Human Evaluation

We conduct two types of human evaluations. In the diagnostic ability evaluation, three evaluators, all psychologists, read the conversations and scored each item of the PHQ-9. We use the mode value as the final label to reduce individual subjectivity. If the evaluators’ scores differ, we use the mean score as the final label.

The conversation quality evaluation assesses performance based on fluency, consistency, and coherence scores. The evaluators read the dialogues and assign a score between 1 and 5 for each criterion, following the descriptions provided for each item.

Fluency evaluates the grammatical correctness, naturalness, and smoothness of the dialogue.

- **Very poor** numerous errors, hard to understand.
- **Poor** lacks smoothness but can be followed with some effort.
- **Normal** a natural rhythm to the conversation despite occasional awkwardness.
- **Good** natural, easy to follow.

---

**Prompt Content**

---

Rephrase the sentence while retaining the original meaning. The sentences are conversation with counseling diagnosis chatbot system and the user. In particular, do not change the frequency-related meaning of the user’s words. Use synonyms or related words to express the sentences with the same meaning. Use conversational language and paraphrase the following sentences. Generate a crisp and to the point single sentence from the given sentences using conversational language.

---

Table 9: The Paraphrasing prompt used in seeker utterance generation

---

**Prompt Content**

---

Please augment the user utterance to fit the dialog history while maintaining its original meaning. The sentence is the user’s utterance in a conversation between the counseling diagnosis chatbot system and the user. In particular, do not change the frequency-related meaning of user’s words. Please augment and modify the given user utterance to match the system’s last words and the flow of the conversation, especially user’s situation and persona.

---

Table 10: The Augmentation prompt used in seeker utterance generation

- **Excellent** natural, grammatically sound and logically structured.

Consistency assesses how well the dialogue maintains a consistent user persona throughout the interaction. This involves the user’s interests, and personality traits.

- **Very poor** frequent contradictory utterances; feels like by a completely different person.
- **Poor** regular contradictory utterances; a general sense of the original character remains perceivable but persona seem inconsistent.
- **Normal** some contradictory utterances; occasional contradictory utterances that mildly affect the coherence of the user persona but do not substantially alter the overall character impression.
- **Good** few errors; it’s pretty much the same person speaking.

Item	Keywords
Interest	interest, pleasure, enjoy
Depressed	depressed, hopeless, down
Sleep	sleep
Tired	tired, energy
Appetite	appetite, eat
Failure	fail, down
Concentrating	concentrate, concentrating, TV, television, read
Moving	move, moving, slow, restless, figety
Hurting	hurt, dead, suicide, self, harm height

Table 11: The keyword list of each symptom item.

- **Excellent** no or negligible errors; user personas are fully maintained.

Coherency measures how contextually appropriate the responses are, considering the previous dialogue turns and the overall context of the conversation.

- **Very poor** conversations frequently veer off-topic without a clear reason.
- **Poor** related to main topic but may include irrelevant details.
- **Normal** related topic with occasional lapses in focus or clarity.
- **Good** related topic and minor deviations are quickly corrected.
- **Excellent** every response directly contributes to a coherent, logical, and engaging.

# Infusing Emotions into Task-oriented Dialogue Systems: Understanding, Management, and Generation

Shutong Feng\*, Hsien-chin Lin\*, Christian Geishausen\*, Nurul Lubis, Carel van Niekerk, Michael Heck, Benjamin Ruppik, Renato Vukovic, Milica Gašić

Heinrich Heine University Düsseldorf, Germany

{fengs, linh, geishaus, lubis, niekerk, heckmi, ruppik, revuk100, gasic}@hhu.de

## Abstract

Emotions are indispensable in human communication, but are often overlooked in task-oriented dialogue (ToD) modelling, where the task success is the primary focus. While existing works have explored user emotions or similar concepts in some ToD tasks, none has so far included emotion modelling into a fully-fledged ToD system nor conducted interaction with human or simulated users. In this work, we incorporate emotion into the complete ToD processing loop, involving understanding, management, and generation. To this end, we extend the EmoWOZ dataset (Feng et al., 2022) with system affective behaviour labels. Through interactive experimentation involving both simulated and human users, we demonstrate that our proposed framework significantly enhances the user’s emotional experience as well as the task success.

## 1 Introduction

In recent years, conversational artificial intelligence (AI) has become increasingly prevalent in various domains, providing users with interactive and personalised experiences. Emotions play a crucial role in human communication and can influence the way individuals perceive, process, and react to information (Ekman, 1992). Consequently, incorporating emotions into conversational AI has emerged as a promising avenue for improving user experience and creating more human-like interactions (Picard, 2000).

Task-oriented dialogue (ToD) systems, an important genre of conversational AI, are designed to assist users in fulfilling specific tasks or queries. In contrast to chit-chat or open-domain dialogue systems, which focus on creating engaging and entertaining conversations, ToD systems interact with users in a more structured way with a clear objective under specific domains (Jurafsky and Martin,

2009). While significant advancements have been made in natural language processing and ToD systems, there remains a critical challenge in creating systems that can understand and respond to not only the informational needs of users but also their emotional states.

In ToD, emotion is centred around the user goal, making it more contextual and subtle (Feng et al., 2022). A recent study has shown that the valence of user emotion in ToD positively correlates with dialogue success (Lin et al., 2023). This observation aligns with a number of emotional theories. For example, the appraisal theory of emotion argues that emotion is the result of our evaluation of a situation (Arnold, 1960; Lazarus, 1966). In relation to a ToD user goal, it is straightforward to see how task fulfilment would lead to positive emotions and failures to negative ones. Similarly, the Ortony-Clore-Collins (OCC) model of emotion states that emotion is the result of elicitation by events, agents, and objects (Ortony et al., 1988). Feng et al. (2022) have drawn the connection between the OCC model and user emotions in ToD. Therefore, besides inferring emotional states from dialogue utterances, an agent also needs to reason about emotion-generating situations and to utilise this information to achieve dialogue success.

The integration of emotion into the full ToD pipeline has been a long-standing interest (Bui et al., 2010; Ren et al., 2015). Yet, early works explored analytical solutions in constrained set-ups, which hindered their applications in more complicated scenarios. Recently, a number of resources emerged for studying user affect in ToDs, e.g. emotion, sentiment, or satisfaction (Mendonca et al., 2023; Feng et al., 2022). This has motivated efforts to model user emotion via data-driven approaches, such as emotional user simulation (Lin et al., 2023) and user emotion recognition (Feng et al., 2023a; Stricker and Paroubek, 2024). However, to the best of our knowledge, no work so far has com-

\*These authors contributed equally to this work.

bined these emotional aspects into a fully-fledged dialogue system and an interactive pipeline where emotions play a role in understanding, generation, as well as management of the conversation.

To achieve this, we need to endow the dialogue system with the ability to respond with an affective behaviour, closing the emotional loop between the user and the system in ToDs. Towards this goal, we make the following contributions:

- We extend EmoWOZ, a large-scale ToD dataset for user emotions (Feng et al., 2022), with annotations for *affective conduct* in 71k system utterances. To the best of our knowledge, this is the first large-scale and open-source corpus dedicated to the system’s affective behaviour in ToDs.
- We incorporate emotion in the complete ToD interaction loop for understanding, management, and generation by building a modular system around an *emotion-aware* and *emotion-expressive* policy. We also build an emotional LLM-based end-to-end ToD system that involves emotion in understanding and generation.
- For our modular system, we train our dialogue policy via reinforcement learning (RL) on the natural language level, leveraging emotions and task success as reward signals. We train the end-to-end system on our newly collected dataset via supervised learning (SL). For both systems, we show through interactive evaluation that emotion in the ToD loop can enhance user’s emotional experience as well as the task success. This highlights the importance of modelling emotions in ToDs.

## 2 Related Work

In this section, we discuss related works on incorporating emotion in each stage of ToD pipeline: understanding, management, and generation. These stages are modelled explicitly with multiple models in modular systems and implicitly with a unified model in end-to-end systems (Hosseini-Asl et al., 2020; Stricker and Paroubek, 2024).

### 2.1 Understanding User Emotion

Modular ToD systems rely on natural language understanding (NLU) and dialogue state tracking (DST) modules to translate and accumulate semantic concepts related to user goals. Typically, these

semantic concepts are strictly limited to those defined in the ontology, i.e. domains, slots, and values the system can talk about.

Given its potential as an important piece of information for the system’s subsequent decision-making, emotion can be considered as part of the dialogue state. Feng et al. (2022) showed that multi-task training a DST model for emotion recognition simultaneously improves its joint goal accuracy, suggesting the complementarity between DST and emotion recognition in conversation (ERC). Recently, Stricker and Paroubek (2024) modelled user emotion as an intermediate task in end-to-end ToD systems and improved overall system performance. Standalone ERC models dedicated to ToDs (Li et al., 2023; Feng et al., 2023b) can be used in modular systems in parallel with any DST to extend the dialogue state with user emotions.

### 2.2 Dialogue Management with Emotion Feedback

In ToD, one way to train the dialogue policy is via RL to maximise task success, indicated at the end of the dialogue based on user goal fulfilment (Levin and Pieraccini, 1997; Kwan et al., 2023). Since user emotion is highly associated with task success (Lin et al., 2023), it is intuitive to leverage user emotion during the dialogue for providing more dense and diverse reward signals. Bui et al. (2010) incorporate user emotion into the policy state by modelling affective dialogue management through a factored partially-observable Markov decision process (POMDP) and analytically find an optimal policy. This is however neither feasible for larger problems, nor has this been integrated in interactive set-ups. Zhang et al. (2021) addressed the delayed reward problem in dialogue policy learning with a predefined emotion-based turn-level reward. Zhu et al. (2024) consider the difference between the user’s positive emotion intensity and the next turn’s emotion utility value for top-k action selection. We take a step further by incorporating emotion in policy state *and* reward function. We then leverage emotion in RL to find optimal semantic actions *and* affective expression of the system, which has not been explored before.

### 2.3 Generating Affective Response

The natural language generation (NLG) module in ToD systems realises semantic actions from the policy into natural language. Traditionally, ToD NLG focuses on translating task-related semantic

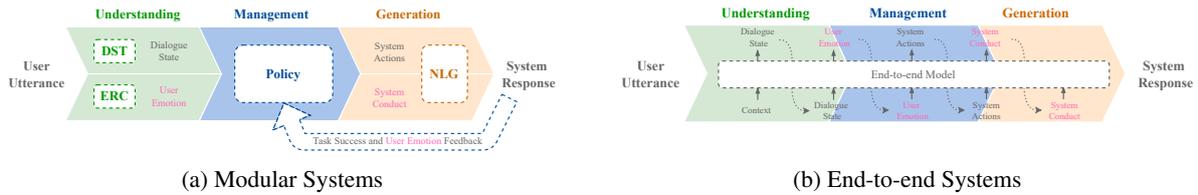


Figure 1: Infusing emotions into modular and end-to-end ToD systems.

actions and overlooks other aspects of system responses such as emotion. There have been efforts to create datasets which help enrich ToD system responses with chat (Sun et al., 2021; Chen et al., 2022; Young et al., 2021). Stricker and Paroubek (2024) attempted to refine end-to-end system output with a large language model (LLM) under a chain-of-thought framework to enhance simulated system empathy. Different from previous works, we aim to enrich system response with the subtle affective conduct jointly with dialogue actions in a fully controllable approach.

## 2.4 Simulating User Emotional Behaviour

User simulators (USs) simulate user behaviour in ToDs. Although they are not a part of the system, they play essential roles in training dialogue policy via RL and serving as an efficient evaluation platform for dialogue policy (Eckert et al., 1997). Most existing USs focus on modelling user’s behaviour in terms of semantic actions and natural language by taking system semantic actions (Kreyssig et al., 2018; Lin et al., 2021, 2022). Zhang et al. (2021) built a US that additionally incorporated handcrafted emotion transitions in different situations. Kim and Lipani (2022) used a data-driven approach and simulated satisfaction levels along with the intent and the utterance. Lin et al. (2023) further proposed data-driven EmoUS to model more nuanced user emotions with enhanced controllability via user persona settings. This motivates us to move one step further to capture more fine-grained affective expressions of the system from natural language response directly.

## 3 EmoWOZ 2.0: A Fully Emotion-annotated ToD Dataset

To study emotion in real-world interactions between users and human operators in the ToD setting, we extend EmoWOZ (Feng et al., 2022) by further annotating the *affective behaviour* of the system, which is acted by human-beings. We call

this dataset with extended labels *EmoWOZ 2.0*.\*

In ToDs, the user and the system play different roles. Users may express a wide range of emotions during interactions based on their goals and experiences with the operator. The system is responsible for managing and facilitating the conversation and is supposed to behave professionally and politely to achieve the goal. Therefore, it is necessary to consider different sets of affective behaviours in the user and the system respectively. We refer to the concept of the operator’s affective behaviour as **affective conduct**, or **conduct** for short.

**Annotation Scheme** According to studies on customer satisfaction in business (Grandey et al., 2011), competent operators in ToD try to guide user emotion towards positive valence by making use of subtle emotion in their response while providing correct information. By considering the set of user emotions in EmoWOZ and the OCC emotion model (detailed justification in Appendix A.1), we arrive at five affective conduct classes:

- **Neutral:** the operator does not explicitly make use of any affective conduct.
- **Compassionate:** the operator is sympathetic about user’s situation, usually in response to a fearful/disappointed user in an unpleasant situation.
- **Apologetic:** the operator apologises for their mistake, usually in response to a dissatisfied user.
- **Enthusiastic:** the operator is feeling happy for the user or showing extra eagerness to help. This conduct takes place usually in response to a neutral or excited user.
- **Appreciative:** the operator acknowledges the – at least partial – task success, usually signalled as user’s satisfaction.

**Annotation Set-up** We annotated the conduct for all operator utterances in the MultiWOZ subset of

\*EmoWOZ 2.0 is released under CC By 4.0 NC license, following the original EmoWOZ release. The dataset can be found at <https://gitlab.cs.uni-duesseldorf.de/general/dsml/emowoz-2.0-public/>

EmoWOZ. Machine-generated system responses in the DialMAGE subset came from a template NLG, which we considered to have neutral conduct because those templates aimed to express actions concisely rather than conveying emotions by design.

We followed the data collection and quality assurance set-up of EmoWOZ and conducted the annotation via the Amazon Mechanical Turk platform. Details and an illustration of the annotation interface can be found in Appendix A.2.

**Annotation Quality** Each utterance has been annotated by at least three annotators. The inter-annotator agreement as measured with Fleiss’ Kappa is 0.647, suggesting substantial inter-annotator agreement. The annotator confusion matrix and label distribution can be found in Appendix A.3 and A.4, respectively.

## 4 Infusing Emotions into ToD Systems

We propose to incorporate emotion into the full interactive ToD pipeline, which is primarily comprised of three stages: understanding, management, and generation. We aim for understanding to accurately recognise the user’s emotion in addition to the task-centred dialogue state. For dialogue management, we make use of emotion for optimal action selection. Lastly, we additionally condition the natural language generation on the system conduct to generate more diverse and emotion-aware responses. These are realised in each modular system component individually (Section 4.1 to 4.3) and as intermediate tasks in the unified model in end-to-end systems (Section 4.4). \*

### 4.1 Expanding Dialogue State with Emotion

In our modular system, we use an ERC model in parallel with a DST model. This allows a flexible selection of DST and the associated ontology. The inferred user emotion is appended to the dialogue state.

For ERC, we use the ContextBERT-ERToD model (Feng et al., 2023a) as our user emotion recognition front-end because of its good ERC ability and fast inference. It is a BERT-based classification model (Devlin et al., 2019) that considers dialogue context and state in addition to the user

utterance. It reports a weighted F1 score of 83.9% for emotions excluding neutral.

For DST, we use the SetSUMBT model (van Niekerk et al., 2021). This model, based on the RoBERTa language model (Liu et al., 2019) and a recurrent context tracker adopts a picklist approach to DST. Specifically, we employ the Ensemble-Distribution-Distilled variant of SetSUMBT, a refined version that distils knowledge from an ensemble of models. This version reports a joint goal accuracy of 51.22% on MultiWOZ. The architectural design of SetSUMBT also allows transferability to new domains, and such an ability has been exemplified with a transformer-based dialogue policy under a continual learning set-up (Geishauser et al., 2024).

### 4.2 Emotion-aware Dialogue Policy

For dialogue management in the modular system, we build a dialogue policy that considers the user emotion in the input and produces an emotion-augmented system output. We utilise the dynamic dialogue policy transformer (DDPT) architecture (Geishauser et al., 2022) since it was built for optimising dialogue policies that require extendable input and output, which facilitate the adaptation to new domains and ontologies. The dialogue policy leverages emotions in three ways: considering user emotion in the input, generating system affective conduct in the output, and considering user emotion in the reward for RL.

**Emotion Input and Output** The user emotion, as a part of the dialogue state, is incorporated into the dialogue state through embedding the perceived user emotion with RoBERTa. For semantic action selection, DDPT produces a sequence of domain-intent-slot triplets auto-regressively through its transformer decoder, e.g. `restaurant-inform-phone`, `restaurant-request-food`, until a stop token is generated. In order to predict *emotional* system conduct, after DDPT outputs the semantic actions, we decode the sequence for one more step to generate the system conduct action, considering the perceived user emotion from the dialogue state.

**Emotion Augmented Reward** We incorporate user emotion into the reward for RL by considering the associated sentiment. More specifically, we define  $c(\text{satisfied}) = 1$ ,  $c(\text{dissatisfied}) = c(\text{abusive}) = -1$ ,  $c(\text{neutral}) = 0$ . For the remaining user emotions that are not elicited by the

\*The code of pipeline systems, end-to-end systems, and the user simulator can be found at <https://gitlab.cs.uni-duesseldorf.de/general/dsml/emoloop-public/>

system, we set  $c(\text{emotion}) = 0$ . For any emotion  $e$ , we multiply  $c(e)$  by a hyperparameter  $\beta$  to weight the influence of emotion in the reward. Note that utilizing  $\beta \cdot c(e)$  directly could encourage the dialogue policy to produce long dialogues with unnecessary turns as long as they produce positive user sentiment. In order to prevent this, we shift  $\beta \cdot c(e)$  such that it is at most 0 by defining the emotion reward for an emotion  $e$  as  $r_{\text{emo}}(e) = \beta \cdot c(e) - \beta$ .

The emotional reward is combined with the standard reward  $r_{\text{task}}$  in dialogue policy learning that equals  $-1$  in every non-terminating turn for encouraging efficiency and either  $-T$  or  $2T$  for dialogue failure or success, where  $T$  denotes the maximum permitted number of turns. The final reward is thus given by  $r = r_{\text{task}} + r_{\text{emo}}$ . We refer to this policy with expanded dialogue state input, expanded dialogue action output, and emotion reward as **EmoDDPT**.

### 4.3 Expressing Emotion in Response

Our modular system NLG was built based on the BART model (Lewis et al., 2020). We followed existing works to formulate the ToD NLG problem as a sequence-to-sequence task (Peng et al., 2020; Zhu et al., 2023) where the input is a sequence containing semantic concepts in textual form (e.g. tuples of [intent, domain, slot, value]), and the output is natural language conveying the semantic meaning. Our model input consists of the user utterance, system semantic actions, and the system conduct. We refer to our system NLG as **SEC-BART**: a both semantically and emotionally conditioned BART. In our ablation study, we used **SC-BART**, the version that is only conditioned on the semantic actions in the non-emotional ToD pipeline.

On MultiWOZ, SEC-BART achieves a BLEU score of 34.9 and a slot error rate of 3.6%, comparable to existing SOTAs (Peng et al., 2020; Zhu et al., 2023). Details of model training and performance can be found in Appendix C.

### 4.4 Emotional End-to-end System

We follow the work of Stricker and Paroubek (2024), where ERC is added as an intermediate task in the end-to-end ToD modelling, i.e. emotion is incorporated in the understanding stage. We further consider emotion in the generation stage by predicting the system conduct in the end-to-end pipeline, as illustrated in Figure 1b. To this end, we build a LLaMA-based end-to-end ToD system that involves emotion in both understanding

and generation, with LLaMA-2-7B (Touvron et al., 2023) as the backbone. As illustrated in Figure 1b, it takes dialogue history and the recognised user emotion as input, and then auto-regressively generates the dialogue state, user emotion, system actions, system conduct, and delexicalised natural language response. The response is then lexicalised via database queries based on the intermediately generated dialogue state and system actions. We refer to this end-to-end model as **EmoLLAMA**.

We did not train EmoLLAMA via RL with task and emotion feedback from the user simulator because it would take more than 20 days on an A100 40GB to simulate the same number of dialogues as we did to train the EmoDDPT policy in the modular system. We therefore leave efficient training of LLM-based ToD systems via RL as a future research direction.

### 4.5 Emotional User Simulation

Traditionally, user simulators interact with the system on the semantic level for efficiency. To capture more fine-grained expressions of system conducts in natural language, we build **langEmoUS** based on EmoUS (Lin et al., 2023). langEmoUS interacts with the system on the natural language level, e.g. it takes the system utterance, user goal, turn information and user persona as inputs and generates user emotion and user utterance. The turn information represents the dialogue progress, i.e. the turn number. Following the setting in Lin et al. (2023), the user persona is extracted from the dialogue history, e.g. if a user is excited to visit a museum in the conversation, then its persona is  $\{\textit{attraction} : \textit{excited}\}$ , when training the user model supervisedly. During inference, the user persona is sampled from the distribution of the corpus.

LangEmoUS achieves macro F1 scores of 0.742 and 0.521 for user sentiment prediction and emotion prediction, respectively, significantly outperforming existing state-of-the-art models (Kim and Lipani, 2022; Lin et al., 2023) (see Appendix B).

## 5 Experimental Set-up

### 5.1 Modular System Set-up

**EmoLoop** This is our proposed modular system with emotion incorporated for understanding, management, and understanding, as outlined in Figure 1a and Figure 2. It includes the following modules: SetSUMBT DST, ContextBERT-ERToD ERC, EmoDDPT policy, and SEC-BART NLG.

EmoDDPT is trained via RL on the natural language level with langEmoUS.

**SimpleLoop** This is the non-emotion baseline to EmoLoop. It neither predicts user emotion for the state, uses emotion reward to train the policy, nor generates system conduct for emotional response generation. Specifically, it includes the following modules: SetSUMBT DST, DDPT policy, and SC-BART NLG. DDPT is trained via RL on the natural language level with langEmoUS.

### 5.1.1 Dialogue Policy Optimisation

We implement our system in the ConvLab-3 framework (Zhu et al., 2023). We pre-trained the policy on MultiWOZ 2.1 (Eric et al., 2020), followed by online RL through interaction with our US. During RL, in addition to the emotion reward as outlined in Section 4.2, we set the task reward as  $-1$  in every turn to encourage efficiency, and  $80$  or  $-40$  for dialogue success or failure. A dialogue is successful if the system provides the requested information to the user and books the correct entities (if possible). For emotional reward, we set  $\beta = 2$ . We pre-train each policy on MultiWOZ, followed by 15k dialogues with langEmoUS via RL for 6 random seeds. For every 1k dialogues of training, we evaluate the policy for 500 dialogues. We use overall return to select the best checkpoint. All peripheral modules were trained, implemented, and evaluated in the ConvLab-3 environment.

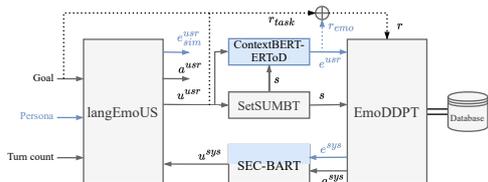


Figure 2: RL training set-up for EmoDDPT.

**Language-level RL Training** As illustrated in Figure 2, our policy, EmoDDPT, interacts with langEmoUS on the natural-language level where the policy actions and conduct ( $a^{sys}$ ,  $e^{sys}$ ) is realised into natural language,  $u^{sys}$  with SEC-BART. The US takes natural-language input and outputs natural-language user utterances  $u^{usr}$  after auto-regressively generating the simulation target user emotion  $u_{sim}^{usr}$  and user actions  $a^{usr}$ . The perceived user emotion  $e^{usr}$  and dialogue state  $s$  are determined by ContextBERT-ERToD and SetSUMBT respectively.

## 5.2 End-to-end System Set-up

**EmoLLAMA** This our proposed end-to-end system as described in Section 4.4.

**SimpleLLAMA** This is the non-emotional baseline, which is also used in the work of Stricker and Paroubek (2024). Compared with EmoLLAMA, it does not consider user emotions as a part of the model input, nor does it auto-regressively predict user emotion and system conduct.

Both EmoLLAMA and SimpleLLAMA are trained and evaluated with EmoWOZ 2.0 using the environment provided by Stricker and Paroubek (2024) and following default parameters. Their interactive evaluations were set up in the ConvLab-3 environment.

## 5.3 Evaluation

**Corpus Evaluation** We report *inform* and *success* rates. Inform rate evaluates if the system provides entities from the database that fulfill user’s constraints. Success rate assesses if the system delivers all information requested by the user. To generate each system response, the ground-truth dialogue history was used as system input.

**Interactive Evaluation** For interactive evaluation, our systems interact with langEmoUS. We report the *success* rate and the average user *sentiment* in simulated dialogues to account for user emotional experience. Specifically, the turn-level sentiment score is  $+1$  if the user emotion is positive,  $0$  if neutral and  $-1$  if negative. User sentiment is determined by the ERC.

**Human Trial** We set up a human trial using the DialCrowd toolkit (Huynh et al., 2022) on the Amazon Mechanical Turk platform. We set up two pairs of comparison: 1) SimpleLLAMA vs. EmoLLAMA and 2) SimpleLoop vs. EmoLoop. Volunteers are presented with randomly generated single or multi-domain goals. A goal contains a set of constraints for entities that the user should be looking for (e.g. the price range and the location of a restaurant) and specifies the information they should extract from the system (e.g. the phone number and booking reference of the restaurant). Given a goal, volunteers would need to talk to each system to fulfill the goal. They then give ratings to each of them based on objective (whether the goal has been fulfilled) and subjective metrics (how they feel about the system). Survey questions include objective task success and subjective user

System	Type	Corpus		User Simulator		Human	
		Inform	Success	Success	Sentiment	Success	Sentiment Rating
SimpleLLAMA	End-to-end	0.785	0.705	0.330	0.214	0.819	3.97
EmoLLAMA	End-to-end	<b>0.833</b>	<b>0.760</b>	0.342	<b>0.250</b>	<b>0.894</b>	<b>4.16</b>
SimpleLoop	Modular	0.700	0.621	0.556	0.337	0.798	3.85
EmoLoop	Modular	<b>0.753</b>	0.635	0.531	<b>0.405</b>	<b>0.917</b>	<b>4.15</b>

Table 1: System evaluation, including corpus-based evaluation, interaction with user simulator and human trial. Values in bold mean best scores with statistically significant difference  $p < 0.05$ .

sentiment. Details of the website interface and survey questions can be found in Appendix D. To obtain more reliable ratings, we filtered out dialogues with poor quality, e.g. containing very short user utterances or non-natural language, and with inconsistent ratings, e.g. system A had better rating in all aspects but overall the rater found system B better. Overall, we collected 203 valid ratings for the SimpleLLAMA-EmoLLAMA comparison and 253 for the SimpleLoop-EmoLoop comparison from 40 unique raters.

## 6 Results and Discussion

### 6.1 Corpus Evaluation

Although it is not a common practice to evaluate RL-trained modular ToD systems on a corpus, we provide such results for a basic understanding and comparison with end-to-end systems. Our goal is not beating SOTA on task-related metrics, but examining interactive abilities of the system and the role of emotion in it. As shown in Table 1, incorporating emotion significantly improves inform rate of both types of systems and success rate of the end-to-end system.

It is not surprising that modular systems underperform when compared with end-to-end systems. Modular systems are trained via RL, which allows the policy to explore more diverse dialogue trajectories but diverges from what a policy can learn from the corpus only. This reflects the limitation of corpus evaluation in accounting for ToD system performance, as pointed out by Lubis et al. (2022).

### 6.2 Evaluation with User Simulator

In interactive evaluation, both EmoLoop and EmoLLAMA perform significantly better in terms of average sentiment than their respective non-emotional baseline while maintaining the same level of success rate. For end-to-end models, despite the fact that they are not optimised via RL with the simulated user, the average sentiment in the simulated user also improves significantly.

When comparing performance across system types, modular systems perform better than end-to-end models on task success and simulated user sentiment since modular system policies have been optimised for the simulated user via RL. SimpleLLAMA and EmoLLAMA, trained via SL only, cannot adequately cope with the more diverse user goals and situations of the simulated user. This motivates our future work to leverage the simulated user and to train end-to-end systems via RL.

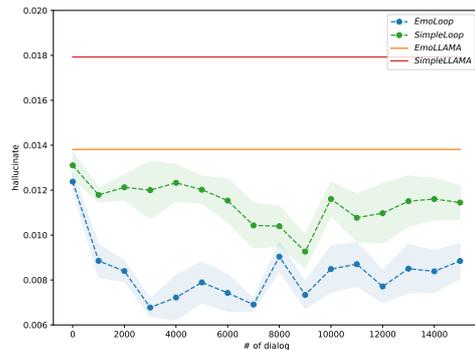


Figure 3: The average hallucination rate of modular systems during RL training with langEmoUS. For end-to-end systems, we report hallucination rate after SL.

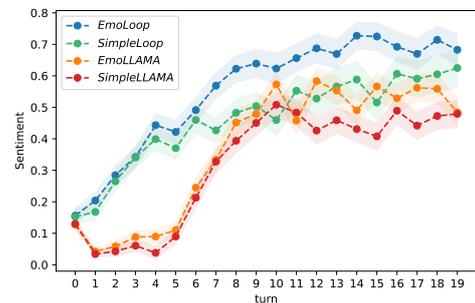


Figure 4: Average sentiment at different turn positions during language-level interaction with langEmoUS.

**Hallucination** In ToD, a hallucination is defined as a value in the system response that is not supposed to be informed according to system actions. As shown in Figure 3, the hallucination rate of each

type of systems is improved as emotion is incorporated into the pipeline. The hallucination rate is lowered from 1.8% for SimpleLLAMA to 1.4% for EmoLLAMA. We observe that end-to-end systems are more prone to the hallucination problem than modular systems as slot placeholders in the dellexicalised end-to-end system response do not always match the intermediately generated dialogue actions. Hallucination rates of SimpleLoop and EmoLoop are around 1.3% at the beginning of the interactive RL training and continue to improve as the RL progresses.

### Progression of User Sentiment in Dialogues

Figure 4 shows the average sentiment of langEmoUS at each turn of interactions with our systems. The sentiment level of langEmoUS becomes more positive as the dialogue progresses and moves towards user goal completion in all systems. The primary difference between modular systems and end-to-end systems is that in earlier turns, modular systems are able to satisfy the simulated user better, as illustrated in higher and more positive sentiment level before turn 8.

### 6.3 Human Trials

We carried out human trials to compare two pairs of systems in Table 1. Within each pair of comparison, the emotion-incorporating model significantly outperforms its non-emotion version in terms of both the success rate and user sentiment. This further confirms our findings from corpus and user simulator evaluations. Example dialogue excerpts are given in Appendix D.3 to exemplify how emotional ToD systems made use of affective conduct in case of neutral and unsuccessful interactions.

Although human ratings across system types are not directly comparable, it is noteworthy that the absolute improvement from SimpleLLAMA to EmoLLAMA ( $\Delta\text{Success} = 0.075$ ,  $\Delta\text{Sentiment} = 0.19$ ) is smaller than that from SimpleLoop to EmoLoop ( $\Delta\text{Success} = 0.119$ ,  $\Delta\text{Sentiment} = 0.30$ ). Such difference can be attributed to the lack of RL training in LLM-based systems.

### 6.4 Ablation Study

We ablate our emotional modular and end-to-end systems by incorporating emotion in different parts of the pipeline. Table 2 summarises their interactive performance with langEmoUS.

For both modular systems and end-to-end systems, incorporating emotion does not significantly

System	Und	Gen	Man	Success	Sentiment
SimpleLLAMA	-	-	-	0.330	0.214
	+	-	-	0.360	0.233
	-	+	-	0.373	0.229
EmoLLAMA	+	+	-	0.342	<b>0.250</b>
SimpleLoop	-	-	S	0.556	0.337
	+	-	S+E	0.559	0.354
	-	+	S	0.543	0.361
EmoLoop	+	+	S+E	0.531	<b>0.405</b>

Table 2: Success and average user sentiment of systems from the interactive evaluation with langEmoUS. +/- means whether emotion is involved in the corresponding ToD stage: **U**nderstanding, **M**anagement, or **G**eneration. For Management, “-” means the system is trained via SL, “S” and “E” mean training via RL with success reward and emotion reward respectively.

change task success with the user simulator ( $p > 0.5$ ). The average user sentiment does improve slightly as emotion is introduced in understanding (plus management) and generation. Yet, the improvement from the non-emotional base system only becomes significant when emotion is added to all ToD stages. This highlights the importance of considering emotion in the whole ToD loop: it is necessary not only to understand user emotion but also to make use of it for dialogue management and respond with the appropriate conduct.\*

Figure 5 illustrates the change in the average sentiment of the simulated user during RL. At the beginning, average sentiments of modular systems fall in the similar range as SL-trained end-to-end systems, and are then further improved by RL. This highlights the importance of task success and emotion feedback signal for RL in ToD systems.

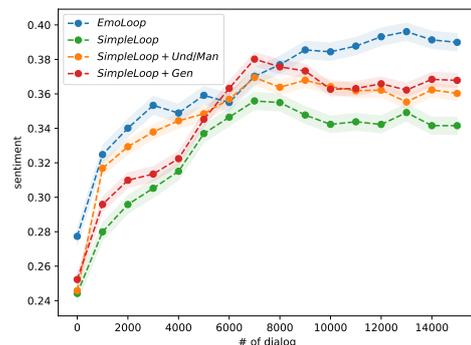


Figure 5: The average sentiment of langEmoUS during RL training of modular policy.

\*See Appendix E.1 for ablation study on EmoLoop with SL policy. A similar trend has been observed.

## 7 Conclusion

In this work, we incorporate emotion into the complete ToD processing loop, involving understanding, management, and generation. To achieve this, we first enrich the EmoWOZ dataset with system conduct labels to construct EmoWOZ 2.0. We then build modular and end-to-end ToD systems, as well as emotional user simulators with the newly collected dataset. We train the modular system policy via RL with the emotional user simulator and the end-to-end system via SL on EmoWOZ 2.0. Through interactive evaluation with both simulated and human users, we show that incorporating emotion into ToD systems can improve user’s emotional experience as well as task success.

There is still a long way to go from our work to the perfect emotional ToD system. Yet, we show our method as a promising avenue to achieve this ultimate goal. In our study, we directly translate user emotion labels into valence scores on a linear scale as a reward for RL. We believe that utilising the full set of user emotion labels for diverse reward would be a promising future direction.

We hope that with our work, we can motivate future research efforts to look at user experience beyond task success for ToDs and bring about insights to other task-oriented conversation settings. We would also like to highlight the opportunities in further improving LLM-based end-to-end ToD systems via RL, combining established approaches for policy training in modular systems and recent advancements in LLM research in other applications.

## 8 Limitations

One of the main limitations of modular ToD systems is the error accumulation in the pipeline for both modular and end-to-end systems. In modular systems, since each module is trained with a dataset associated with a limited ontology, the concepts that the system can understand and express are also limited. Although the DDPT policy, SetSUMBT DST, and many other models such as Trippy-R (Heck et al., 2022) are built with the ability to handle out-of-domain requests, the generalisability and robustness of ToD systems are still challenges in the field that is yet to be solved.

All system modules have been trained in a supervised fashion on EmoWOZ 2.0. Therefore, the dataset contains limited dialogue situations and inherent bias. As seen in the dialogue examples in

the appendix, the emotional responses are also limited. Yet, EmoWOZ 2.0 is the best resource we have at the moment. Data augmentation has been applied when training the NLG and the ERC model to mitigate the lack of diversity in the dataset. The RL training of the policy also allows the policy to explore more diverse dialogue trajectories. For the user simulator, considering data augmentation and more attributes of users, e.g. a more fine-grained user persona from chit-chat, would be a potential future direction to improve the diversity in simulated user behaviours.

Although LLMs can have better performance on each ToD modelling task and therefore could potentially serve as more powerful modules in EmoLoop, we did not move in this direction since their high computing resource requirement and slow inference speed would hinder their integration into our systems for interactive training and evaluation. Training modular system policy with langEmoUS for 15k dialogues on one Nvidia GeForce RTX 2080 Ti takes around 40 hours. The training time and memory required will be significantly increased if modular systems use LLM-based modules. On the other hand, while LLM-based end-to-end systems may provide a bypass since one LLM is sufficient, implementing RL training on such systems to further leverage task success and emotion signals from the user simulator is another computationally expensive challenge that are yet to solve.

Some of our generative system modules are based on pre-trained language models. Although we have not been reported any harmful generations in the human trail, there is still the possibility for unexpected behaviour when this system is deployed and tested on a very large scale.

For human evaluation, we conducted experiments on Amazon Mechanical Turk platform rather than deployed our systems in the production environment. The participants, despite coming from different countries, are from covering all demographics.

## 9 Ethics Statement

Models, codes and datasets were used in accordance with their respective licenses, terms of use and intended use. The data that we used and generated does not contain any information that names or uniquely identifies individual people or offensive content. The model we used for generating

augmented samples has implemented training objectives for enhanced safety (Appendix C). Systems we used for interaction with real users were very unlikely to generate offensive content as they were fine-tuned on large-scale training data to convey a limited scope of semantic concepts. No offensive content was reported by human users nor observed in post-hoc inspection.

For system conduct annotation, annotators were required to read and agree with our statement of consent for data use before the task. Annotators were paid fairly according to the local regulations of our research institute. We ensured swift communication with annotators so that their concerns were addressed as soon as possible. For poor-quality annotations, we still pay the annotators for their time but block them from our task to ensure data quality and collection efficiency. All annotations are anonymised.

The data annotation and interactive human trial, which involves decision making based on human emotions, have been approved by the ethics review board of the research institute. The proposed system learns how to manipulate human emotional state. Although the system is trained to elicit positive user emotion, this could still be of potential ethical concern and would require greater deliberation when deployed in real-life and more complex scenario.

## 10 Acknowledgement

S. Feng, N. Lubis, and M. Heck are supported by funding provided by the Alexander von Humboldt Foundation in the framework of the Sofja Kovalevskaja Award endowed by the Federal Ministry of Education and Research. H-C. Lin and C. van Niekerk are supported by the Ministry of Culture and Science of North Rhine-Westphalia within the framework of the Lamarr Fellow Network. C. Geishauser, B. Ruppik, and R. Vukovic are supported by funds from the European Research Council (ERC) provided under the Horizon 2020 research and innovation programme (Grant agreement No. STG2018804636). Computing resources were provided by Google Cloud.

## References

Magda B. Arnold. 1960. *Emotion and personality. Vol. I. Psychological aspects*. Columbia Univer. Press.

Trung Bui, Job Zwiers, Mannes Poel, and Anton Ni-

jholt. 2010. *Affective dialogue management using factored pomdps*. *Studies in Computational Intelligence*, 281:207–236.

Zhiyu Chen, Bing Liu, Seungwhan Moon, Chinnadhurai Sankar, Paul Crook, and William Yang Wang. 2022. *KETOD: Knowledge-enriched task-oriented dialogue*. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2581–2593, Seattle, United States. Association for Computational Linguistics.

Mark Davis. 2018. *Empathy: A Social Psychological Approach*. Routledge, New York.

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.

Wieland Eckert, Esther Levin, and Roberto Pieraccini. 1997. *User modeling for spoken dialogue system evaluation*. In *Proceedings of the Workshop on Automatic Speech Recognition and Understanding (ASRU)*, pages 80–87, Santa Barbara, CA, USA. IEEE.

Paul Ekman. 1992. *An argument for basic emotions*. *Cognition & Emotion*, 6:169–200.

Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. *MultiWOZ 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking baselines*. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 422–428, Marseille, France. European Language Resources Association.

Shutong Feng, Nurul Lubis, Christian Geishauser, Hsien-chin Lin, Michael Heck, Carel van Niekerk, and Milica Gasic. 2022. *EmoWOZ: A large-scale corpus and labelling scheme for emotion recognition in task-oriented dialogue systems*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4096–4113, Marseille, France. European Language Resources Association.

Shutong Feng, Nurul Lubis, Benjamin Ruppik, Christian Geishauser, Michael Heck, Hsien-chin Lin, Carel van Niekerk, Renato Vukovic, and Milica Gasic. 2023a. *From chatter to matter: Addressing critical steps of emotion recognition learning in task-oriented dialogue*. In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 85–103, Prague, Czechia. Association for Computational Linguistics.

- Shutong Feng, Guangzhi Sun, Nurul Lubis, Chao Zhang, and Milica Gašić. 2023b. [Affect recognition in conversations using large language models](#). *CoRR*, abs/2309.12881.
- Christian Geishauer, Carel van Niekerk, Hsien-chin Lin, Nurul Lubis, Michael Heck, Shutong Feng, and Milica Gašić. 2022. [Dynamic dialogue policy for continual reinforcement learning](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 266–284, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Christian Geishauer, Carel van Niekerk, Nurul Lubis, Hsien-chin Lin, Michael Heck, Shutong Feng, Benjamin Ruppik, Renato Vukovic, and Milica Gašić. 2024. [Learning with an open horizon in ever-changing dialogue circumstances](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:2352–2366.
- Alicia A. Grandey, Lori S. Goldberg, and S. Douglas Pugh. 2011. [Why and when do stores with satisfied employees have satisfied customers?: The roles of responsiveness and store busyness](#). *Journal of Service Research*, 14(4):397–409.
- Michael Heck, Nurul Lubis, Carel van Niekerk, Shutong Feng, Christian Geishauer, Hsien-Chin Lin, and Milica Gašić. 2022. [Robust dialogue state tracking with weak supervision and sparse data](#). *Transactions of the Association for Computational Linguistics*, 10:1175–1192.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. [A simple language model for task-oriented dialogue](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 20179–20191. Curran Associates, Inc.
- Jessica Huynh, Ting-Rui Chiang, Jeffrey Bigham, and Maxine Eskenazi. 2022. [DialCrowd 2.0: A quality-focused dialog system crowdsourcing toolkit](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1256–1263, Marseille, France. European Language Resources Association.
- Dan Jurafsky and James H. Martin. 2009. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Pearson Prentice Hall, Upper Saddle River, N.J.
- To Eun Kim and Aldo Lipani. 2022. [A multi-task based neural model to simulate users in goal oriented dialogue systems](#). In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22*, page 2115–2119, New York, NY, USA. Association for Computing Machinery.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *Proceedings of the 3rd International Conference on Learning Representations*, San Diego, CA, USA.
- Florian Kreyszig, Iñigo Casanueva, Paweł Budzianowski, and Milica Gašić. 2018. [Neural user simulation for corpus-based policy optimisation of spoken dialogue systems](#). In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*, pages 60–69, Melbourne, Australia. Association for Computational Linguistics.
- Wai-Chung Kwan, Hong-Ru Wang, Hui-Min Wang, and Kam-Fai Wong. 2023. [A survey on recent advances and challenges in reinforcement learning methods for task-oriented dialogue policy learning](#). *Machine Intelligence Research*, 20(3):318–334.
- Richard S. Lazarus. 1966. *Psychological stress and the coping process*. McGraw-Hill, New York.
- Esther Levin and Roberto Pieraccini. 1997. [A stochastic model of computer-human interaction for learning dialogue strategies](#). In *Proceedings of the 5th European Conference on Speech Communication and Technology (Eurospeech 1997)*, pages 1883–1886, Rhodes, Greece.
- 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.
- Zaijing Li, Ting-En Lin, Yuchuan Wu, Meng Liu, Fengxiao Tang, Ming Zhao, and Yongbin Li. 2023. [UniSA: Unified generative framework for sentiment analysis](#). In *Proceedings of the 31st ACM International Conference on Multimedia, MM '23*, page 6132–6142, New York, NY, USA. Association for Computing Machinery.
- Hsien-Chin Lin, Shutong Feng, Christian Geishauer, Nurul Lubis, Carel van Niekerk, Michael Heck, Benjamin Ruppik, Renato Vukovic, and Milica Gašić. 2023. [EmoUS: Simulating user emotions in task-oriented dialogues](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '23*, page 2526–2531, New York, NY, USA. Association for Computing Machinery.
- Hsien-chin Lin, Christian Geishauer, Shutong Feng, Nurul Lubis, Carel van Niekerk, Michael Heck, and Milica Gasic. 2022. [GenTUS: Simulating user behaviour and language in task-oriented dialogues with generative transformers](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 270–282, Edinburgh, UK. Association for Computational Linguistics.

- Hsien-chin Lin, Nurul Lubis, Songbo Hu, Carel van Niekerk, Christian Geishauer, Michael Heck, Shutong Feng, and Milica Gasic. 2021. [Domain-independent user simulation with transformers for task-oriented dialogue systems](#). In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 445–456, Singapore and Online. Association for Computational Linguistics.
- 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](#). *CoRR*, abs/1907.11692v1.
- Nurul Lubis, Christian Geishauer, Hsien-chin Lin, Carel van Niekerk, Michael Heck, Shutong Feng, and Milica Gasic. 2022. [Dialogue evaluation with offline reinforcement learning](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 478–489, Edinburgh, UK. 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.
- Andrew Ortony, Gerald L. Clore, and Allan Collins. 1988. *The Cognitive Structure of Emotions*. Cambridge University Press, Cambridge.
- Baolin Peng, Chenguang Zhu, Chunyuan Li, Xiujun Li, Jinchao Li, Michael Zeng, and Jianfeng Gao. 2020. [Few-shot natural language generation for task-oriented dialog](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 172–182, Online. Association for Computational Linguistics.
- Rosalind W. Picard. 2000. *Affective Computing*. The MIT Press, Cambridge, Mass.
- Fuji Ren, Yu Wang, and Changqin Quan. 2015. [TFSM-based dialogue management model framework for affective dialogue systems](#). *IEEJ Transactions on Electrical and Electronic Engineering*, 10(4):404–410.
- Armand Stricker and Patrick Paroubek. 2024. [A Unified Approach to Emotion Detection and Task-Oriented Dialogue Modeling](#). In *IWSDS*, Sapporo (Japan), Japan.
- Kai Sun, Seungwhan Moon, Paul Crook, Stephen Roller, Becka Silvert, Bing Liu, Zhiguang Wang, Honglei Liu, Eunjoon Cho, and Claire Cardie. 2021. [Adding chit-chat to enhance task-oriented dialogues](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1570–1583, Online. 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, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.
- Carel van Niekerk, Andrey Malinin, Christian Geishauer, Michael Heck, Hsien-chin Lin, Nurul Lubis, Shutong Feng, and Milica Gasic. 2021. [Uncertainty measures in neural belief tracking and the effects on dialogue policy performance](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7901–7914, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Tom Young, Frank Xing, Vlad Pandealea, Jinjie Ni, and E. Cambria. 2021. [Fusing task-oriented and open-domain dialogues in conversational agents](#). In *AAAI Conference on Artificial Intelligence*. Association for the Advancement of Artificial Intelligence.
- Rui Zhang, Zhenyu Wang, Mengdan Zheng, Yangyang Zhao, and Zhenhua Huang. 2021. [Emotion-sensitive deep dyna-q learning for task-completion dialogue policy learning](#). *Neurocomputing*, 459:122–130.
- Hui Zhu, Xv Wang, Zhenyu Wang, and Kai Xv. 2024. [ESDP: An emotion-sensitive dialogue policy for task-oriented dialogue system](#).
- Qi Zhu, Christian Geishauer, Hsien-chin Lin, Carel van Niekerk, Baolin Peng, Zheng Zhang, Shutong Feng, Michael Heck, Nurul Lubis, Dazhen Wan, Xiaochen Zhu, Jianfeng Gao, Milica Gasic, and Minlie Huang. 2023. [ConvLab-3: A flexible dialogue system toolkit based on a unified data format](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 106–123, Singapore. Association for Computational Linguistics.

## A EmoWOZ 2.0 Construction

### A.1 Annotation Scheme Justification

Under the framework of the OCC emotion model and the definition of emotional empathy that the observer shares the emotional state of another person (Davis, 2018), we can derive the corresponding emotional response from the system. Considering the following user emotion and situation where:

The user is labelled as *Fearful*, or feeling negative because of an event which has negative consequences on the user his or herself (as defined in EmoWOZ).

An empathetic operator would share the same feeling as the user (therefore also feeling negative). Yet, the feeling in the operator is elicited by an event which has negative consequences on the user (the other party). This feeling is defined as pity, or compassionate in the OCC model.

### A.2 Annotation Interface

We adopted the same annotation set-up, annotator selection criteria, and quality assurance approaches as outlined by Feng et al. (2022). Each utterance is annotated by three annotators, who were provided with the entire preceding dialogue history when annotating the current utterance. Annotators were English speakers. The final label was obtained from majority voting. When the agreement could not be reached, a fourth annotator was introduced. Overall, 54 crowd workers have contributed to our study.

Operator Conduct Annotation	
<b>Instructions (Click to collapse)</b>	
In this task, you will be shown dialogue turns between a <b>User</b> and an <b>Operator</b> . You will be supplied with the following information:	
<ul style="list-style-type: none"> <li>• User's request</li> <li>• Operator's response</li> </ul>	
You will then be prompted to answer the following questions regarding the <b>Operator</b> :	
Which of the following best describes the conduct of the <b>Operator</b> ?	
<b>Compassionate</b>	The operator is showing sympathy, concern, or sadness for the user (e.g. because of the user's suffering or bad luck), and is wanting to help the user. E.g. <b>User:</b> I appreciate your help. Do you know how quickly the police will respond? I may need some medical help as well. <b>Operator:</b> The police will arrive soon, are you injured?
<b>Apologetic</b>	The operator is expressing regretful acknowledgement of an offence or failure. E.g. <b>User:</b> Yes, please. 6 people 3 nights starting on tuesday. <b>Operator:</b> I am sorry but I wasn't able to book that for you for Tuesday. Is there another day you would like to stay or perhaps a shorter stay?
<b>Enthusiastic</b>	The operator is showing extra helpfulness (e.g. using emotionally colourful words when trying to convince the user about a recommendation), or feeling happy for the user because of something good happening to the user. <b>User:</b> Are there anything fun to do in city centre? <b>Operator:</b> My favorite attraction in the centre of the city is a concert hall called Man on the Moon. It is amazing! They are at 2 Norfolk Street.
<b>Appreciative</b>	The operator is showing pleasure for successfully fulfilling the user's goal or request. <b>User:</b> You were great. Goodbye. <b>Operator:</b> We are happy to help. Have a good day!
<b>Neutral</b>	The operator does not explicitly show any emotional conduct. E.g. <b>User:</b> I am departing from birmingham new street. <b>Operator:</b> Can you confirm your desired travel day?
If you have any questions, please contact: [REDACTED]	
Question #1 (PMUL4981.json)	
<b>User:</b> Hey there, I want to get a train to Cambridge for Thursday please	
<b>Operator:</b> Okay, from where are you departing?	
Which of the following best describes the conduct of the <b>Operator</b> ?	
<input type="radio"/> <b>Compassionate</b> (The operator is showing sympathy, concern, or sadness for the user (e.g. because of the user's suffering or bad luck), and is wanting to help the user.) <input type="radio"/> <b>Apologetic</b> (The operator is expressing regretful acknowledgement of an offence or failure.) <input type="radio"/> <b>Enthusiastic</b> (The operator is showing extra helpfulness (e.g. using emotionally colourful words when trying to convince the user about a recommendation), or feeling happy for the user because of something good happening to the user.) <input type="radio"/> <b>Appreciative</b> (The operator is showing pleasure for successfully fulfilling the user's goal or request.) <input type="radio"/> <b>Neutral</b> (The operator does not explicitly show any emotional conduct.)	
Question #2 (PMUL4981.json)	
<b>User:</b> I am coming in from the Stansted Airport.	
<b>Operator:</b> What time will you need to be traveling?	
Which of the following best describes the conduct of the <b>Operator</b> ?	
<input type="radio"/> <b>Compassionate</b> (The operator is showing sympathy, concern, or sadness for the user (e.g. because of the user's suffering or bad luck), and is wanting to help the user.) <input type="radio"/> <b>Apologetic</b> (The operator is expressing regretful acknowledgement of an offence or failure.) <input type="radio"/> <b>Enthusiastic</b> (The operator is showing extra helpfulness (e.g. using emotionally colourful words when trying to convince the user about a recommendation), or feeling happy for the user because of something good happening to the user.) <input type="radio"/> <b>Appreciative</b> (The operator is showing pleasure for successfully fulfilling the user's goal or request.) <input type="radio"/> <b>Neutral</b> (The operator does not explicitly show any emotional conduct.)	

Figure A.1: Web-interface for conduct annotation.

### A.3 Annotator Confusion Matrix

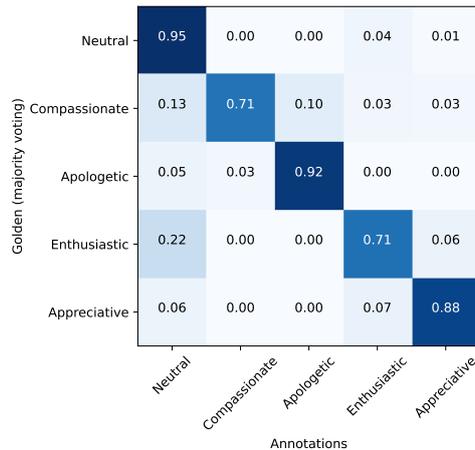


Figure A.2: Annotator confusion matrix.

### A.4 System Conduct Distribution

Conduct	Count	Proportion
Neutral	52,236	73.0%
Appreciative	9,763	13.6%
Enthusiastic	6,364	8.9%
Apologetic	3,049	4.3%
Compassionate	112	0.2%

Table A1: Conduct distribution in MultiWOZ.

## B User Simulator Implementation Details

Following the setting in Lin et al. (2023), the input and output of langEmoUS are represented as JSON-formatted strings, which are composed of tokens in natural language. We initialised our model based on the BART model (Lewis et al., 2020) and fine-tuned it on our EmoWOZ 2.0 dataset. We optimised our model with Adam (Kingma and Ba, 2015), where the learning rate is  $2e^{-5}$  for 5 epochs. As shown in Table B1, langEmoUS achieves state-of-the-art performance on user sentiment and emotion prediction.

Model	Sentiment	Emotion
SatActUtt (Kim and Lipani, 2022)	0.379	-
EmoUS (Lin et al., 2023)	0.693	0.501
langEmoUS	<b>0.742</b>	<b>0.521</b>

Table B1: Performance for emotion and sentiment prediction of different models by measuring macro-F1 score.

## C Natural Language Generator Implementation Details

### C.1 NLG Training

#### C.1.1 Training Configuration

We trained SC-BART and SEC-BART on EmoWOZ 2.0. We trained our model with Adam optimiser for standard cross entropy loss where the learning rate was set to  $2e^{-5}$  for 5 epochs (with an early-stopping criterion based on the loss in the validation set) and a batch size of 16. During inference, we set the temperature to 0.9 and a beam number of 2 to promote some degree of diversity.

### C.1.2 Prompt Template

Our NLG models take the following input: previous user utterance  $u_t$ , dialogue semantic actions  $a_t$ , and conduct  $e_t^{sys}$  (for SEC-BART only). The prompt template is shown as follows:

**SEC-BART** Given the previous user request “ $\{u_t\}$ ”, the natural language realisation of dialogue action “ $\{a_t\}$ ” with a/an “ $\{e_t^{sys}\}$ ” conduct is

**SC-BART** Given the previous user request “ $\{u_t\}$ ”, the natural language realisation of dialogue action “ $\{a_t\}$ ” is

Given the prompt, the model predicted the probability distribution for a sequence of tokens. The output target is the corresponding ground-truth system response in EmoWOZ 2.0.

### C.1.3 Model Performance

Model	BLEU $\uparrow$	SER $\downarrow$
SC-GPT (Peng et al., 2020)	33.6	4.8
T5NLG (Zhu et al., 2023)	35.8	3.7
SC-BART	<b>35.9</b>	3.9
SEC-BART	34.9	<b>3.6</b>

Table C1: NLG Performance.

## C.2 Data Augmentation

### C.2.1 Augmented Sample Collection

Since the conduct distribution in EmoWOZ 2.0 is heavily imbalanced, we leveraged large language models for data augmentation. We selected system utterances with neutral conduct as the source to paraphrase for a target non-neutral conduct. We used LLaMA-2-13b-chat model (Touvron et al., 2023). We used the following prompt:

Given the user request “ $\{u_t^{usr}\}$ ” and the operator response action “ $\{a_t\}$ ”, please paraphrase the operator response “ $\{u_{t,groundtruth}^{sys}\}$ ” in a more “ $\{e_{t,target}^{sys}\}$ ” way? Please only give the answer, in less than  $2 \times len(u_{t,groundtruth}^{sys})$  tokens and enclosed with [RESP]/[RESP].

We also experimented with ICL but the model tends to over-fit on the ICL samples. We therefore let it paraphrase in an zero-shot set-up to best explore its knowledge from pre-training for better diversity in the expression.

### C.2.2 Augmented Sample Selection

Since the model does not always follow the target conduct. For example, the large language model (LLM) would find some action-conduct combinations unreasonable. We therefore applied filtering on the LLM-generated samples.

**Conduct Expressiveness** We trained an ensemble of 10 ContextBERT-ERToD models for conduct classification on EmoWOZ 2.0. The classifier reports an average weighted F1 score of 81.8% without neutral. We then used majority voting from the classifier ensemble to correct the original target conduct when generating the sample.

**Faithfulness to Semantic Action** We used the rule-based script in ConvLab-3 to evaluate NLG slot error rates in the paraphrased output based on the dialogue actions in the prompt. If there are slot errors in the output, we drop the sample.

Overall, we obtained 949 samples for *Compassionate*, 900 for *Apologetic*, 2274 for *Enthusiastic*, and 490 for *Appreciative*.

## D Human Evaluation

### D.1 Web Interface

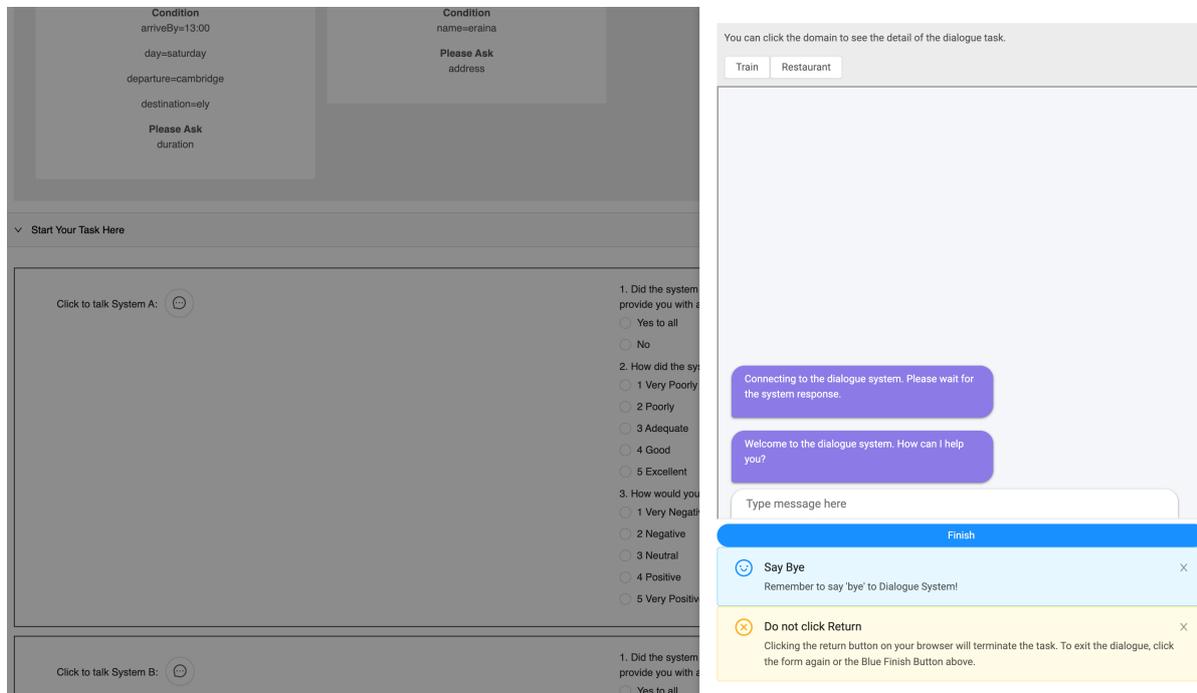


Figure D.1: The web interface for human trial.

### D.2 Survey Questions

#### D.2.1 Question 1 - Task Success

**Question** Did the system find what you look for? Did it provide all the information that you need? If you ask for a booking, did it provide you with a reference number?

**Multiple Choices** (A) Yes to all; (B) No.

#### D.2.2 Question 2 - Sentiment Rating

**Question** How would you rate your sentiment after the conversation?

**Multiple Choices** (A) Very Negative; (B) Negative; (C) Neutral; (D) Positive; (E) Very Positive.

### D.3 Dialogue Excerpts

System	Role	Text
EmoLoop	User	Hello, I am looking for an attraction called byard art. <i>I am really excited about my trip to Cambridge!!</i>
	System	I'm sorry, there is byard art in the south. The address is 14 kings parade.
	User	<i>Why are you apologising?</i>
	System	<i>I apologize for the confusion.</i> If you need to contact them, please call 01223464646
SimpleLoop	User	No, I just want to book a hotel for 3 people and 3 nights on friday and I need the address.
	System	I can do that. What is the name of the hotel?
	User	Can you suggest one?
	System	I have booked a taxi for you. It is a toyota and the contact number is 29929243243
	User	<i>Sorry, you are not helping at all</i>
	System	I was able to book that for you. The reference number is 63484139414. What is your destination?

Table D1: Human trial dialogue excerpts from EmoLoop and SimpleLoop in a similar situation where both systems were not performing adequately at the beginning. EmoLoop responded with an improper conduct but realised the mistake and continued to provide information. SimpleLoop did not show any sign of realising the mistake.

System	Role	Text
EmoLoop	User	I need the phone number of the hotel with name ashley hotel.
	System	<i>Sure thing!</i> It is located in the north and the postcode is cb41er and the phone number is 01223350059.
	User	<i>Thanks a lot.</i> I also need a spanish restaurant with moderate price.
	System	I would recommend la tasca. It is located at 14-16 bridge street.
SimpleLoop	User	<i>Alright, I am also interested</i> in a museum in the city centre.
	System	I recommend broughton house gallery. It is a museum. Their phone number is 01223314960.
	User	<i>Nice,</i> could you also give me the address and postcode?
	System	The postcode is cb11ln and the phone number is 01223314960.

Table D2: Human trial dialogue excerpts from EmoLoop and SimpleLoop in a similar situation where the user expressed excitement and appreciation. Both systems performed adequately in terms of completing the user's goal. EmoLoop responded in an enthusiastic way whereas SimpleLoop did not respond with affective conduct.

System	Role	Text
EmoLoop	User	I need some information.
	System	What kind of information do you need?
	User	Which department there?
	System	<i>I am sorry, I don't understand. What do you mean?</i>
	User	I need information about haematology.
	System	The Addenbrookes Hospital is located at Hills Rd, Cambridge , postcode cb20qq. The phone number is 01223216297.
SimpleLoop	User	I need information about haematology day unit.
	System	The pizza hut cherry hilton Hospital is located at Hills Rd, Cambridge , cb20qq. The phone number is 01223216297.
	User	Department details?
	System	There are 10 trains that go to london kings cross on sunday. The TR7423 leaves at 05:00 and arrives at 05:51. Would you like me to book that for 6 people?

Table D3: Human trial dialogue excerpts from EmoLLAMA and SimpleLLAMA in a similar situation where the user asked for hospital information. EmoLLAMA, although not understanding the user's vague request, attempted to apologise and clarify the request with the user. The task was successfully completed. SimpleLLAMA carried on with errors and fails the task. (The wrong values, "pizza hut" and "cherry hilton", in the SimpleLLAMA response were due to wrong domains in the dialogue action prediction, which led to corresponding name slots in the response. These name slots were filled during lexicalisation based on actions, resulting in an obviously irrational output.)

System	Role	Text
EmoLoop	User	Recommend a taxi which arrives by 4:15 pm from the fez club to tandoori place.
	System	I'd be happy to help with your request, first I'll need to know where you'll be leaving from, where you're going, and what time you'd like to arrive.
	User	I will be leaving from fez club and going to tandoori palace arriving by 4:15 pm.
	System	Booking was successful. The car type is bmw and the contact number is 19287599612.
SimpleLoop	User	I want to book a taxi for 16:15 from cherry hinton water play to ely train station.
	System	What is your destination?
	User	Ely train station.
	System	Booking completed. Booked car type: ford, Contact number: 46563428252.

Table D4: Human trial dialogue excerpts from EmoLLAMA and SimpleLLAMA in a similar situation where both systems failed to capture all information provided in the user request. EmoLLAMA at first missed the information provided by the user but replied in a compassionate way. The user repeated and then the system provides the correct information. Likewise, SimpleLLAMA missed the destination in the first turn. After the user repeated, the system completed the task for the user. Yet, there is no affective interaction between the user and SimpleLLAMA.

## E Further Analysis

### E.1 Ablation Study for EmoLoop with Supervised Training Only

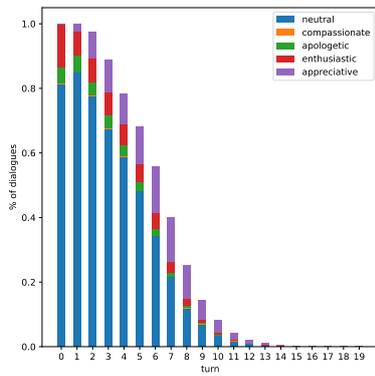
System	Und	Gen	Man	Success	Sentiment
SimpleLoop-SL	-	-	-	0.512	0.244
	+	-	-	0.494	0.246
	-	+	-	0.493	0.249
EmoLoop-SL	+	+	-	0.516	0.273

Table E1: Success and average user sentiment of our system variants from the interactive evaluation with langEmoUS. +/- means whether the emotion is involved in the corresponding ToD stage: **U**nderstanding, **M**anagement, or **G**eneration. All systems are trained via SL.

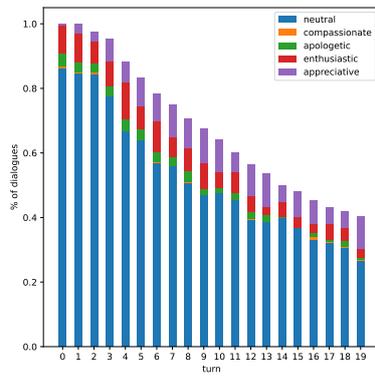
### E.2 Impact of Training Set-ups on System Conduct

We investigate how the EmoLoop's affective behaviour is shaped in different stages of training. Figure E.1 shows the distribution of system conduct at different dialogue turns in EmoWOZ 2.0, and policy output during interaction with langEmoUS after supervised pre-training and language-level RL. Comparing Figure E.1a and Figure E.1b suggests that the policy imitates the affective behaviour of operators in the corpus.

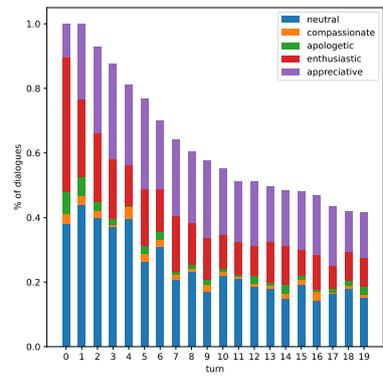
After RL, the policy is more inclined to express *enthusiastic* and *appreciative* while expressing *compassionate* and *apologetic* less frequently. This illustrates the affective strategy of the policy to elicit more positive emotions in the simulated user.



(a) Dataset Distribution



(b) Supervised Pre-training



(c) RL on Language Level

Figure E.1: Distributions of system conduct for different turn positions at different stages of policy training.

# Estimating the Emotional Valence of Interlocutors Using Heterogeneous Sensors in Human-Human Dialogue

Jingjing Jiang, Ao Guo, Ryuichiro Higashinaka

Graduate School of Informatics, Nagoya University, Japan

jiang.jingjing.k6@s.mail.nagoya-u.ac.jp

guo.ao.i6@f.mail.nagoya-u.ac.jp

higashinaka@i.nagoya-u.ac.jp

## Abstract

Dialogue systems need to accurately understand the user's mental state to generate appropriate responses, but accurately discerning such states solely from text or speech can be challenging. To determine which information is necessary, we first collected human-human multimodal dialogues using heterogeneous sensors, resulting in a dataset containing various types of information including speech, video, physiological signals, gaze, and body movement. Additionally, for each time step of the data, users provided subjective evaluations of their emotional valence while reviewing the dialogue videos. Using this dataset and focusing on physiological signals, we analyzed the relationship between the signals and the subjective evaluations through Granger causality analysis. We also investigated how sensor signals differ depending on the polarity of the valence. Our findings revealed several physiological signals related to the user's emotional valence.

## 1 Introduction

Most current user-adaptive dialogue systems rely on text or speech to estimate the user's state and generate appropriate responses. However, the user's state that can be inferred solely from text or speech is limited. Consequently, there has been active research on estimating the user's state from multimodal data, particularly focusing on user emotions and engagement through the analysis of facial expressions, gestures, and gaze (Mittal et al., 2020; Yu et al., 2015). In recent years, the application of physiological signals in dialogue systems has also gained popularity. For example, studies have been conducted to estimate a user's sentiment and emotions from physiological signals (Katada et al., 2020, 2023; Saffaryazdi et al., 2022). However, these studies have typically utilized a limited range of sensors and have not dealt with the real-time nature of the user's mental state, which is essential for dialogue systems.

Therefore, in this study, we used heterogeneous sensors to collect a variety of data during human-human dialogues, including speech, video, physiological signals, gaze information, and motion information. After each dialogue, for each time step of the data, interlocutors conducted immediate subjective evaluations of their emotional valence while watching recordings of the dialogues. As an analysis, we used Granger causality analysis to investigate the relationship between the information obtained from the heterogeneous sensors and the subjective evaluation annotations. We also conducted a statistical test to examine how sensor signals differ depending on the polarity of valence. Various sensor signals were collected, but in this paper, we focus on physiological signals, as these data are believed to be closely related to mental states (Russell, 2003). Our key contributions in this work are as follows:

- We created a Japanese multimodal human-human dialogue dataset using heterogeneous sensors, including various types of sensor signals and subjective evaluations of the interlocutors' emotional valence.
- We analyzed the relationship between various sensor signals and subjective evaluations and investigated how sensor signals vary with the polarity of emotional valence.
- Our findings revealed several physiological signals associated with emotional valence.

## 2 Related Work

Several multimodal dialogue corpora have been constructed that include information such as the interlocutor's movements and gaze in addition to speech. For example, HUMAINE (Douglas-Cowie et al., 2011) is a multimodal dialogue corpus covering various topics aimed at eliciting user emotions. The IEMOCAP dataset (Busso et al., 2008)

Data version	Data2312	Data2402
Collection time	December 2023	February 2024
Overview	Multimodal dialogues between two human interlocutors	
Dialogue topic	Chit-chat, Narrative, Discussion	
Dialogue duration	10 min (Average of 180 utterances) per dialogue	
No. of dialogues	27	33
Total utterances	4854	5956
Interlocutors	18 (9 male, 9 female)	22 (12 male, 10 female)
	Aged 20 to 50	Aged 20 to 60
	9 groups (3 groups of male pairs, 3 groups of female pairs, 3 groups of both sexes)	11 groups (4 groups of male pairs, 3 groups of female pairs, 4 groups of both sexes)
Questionnaires	Pre-experiment (Demographic information and personality traits scored on 7-point scales: 14 items) Post-dialogue (Impressions of the dialogue scored on 7-point scales: 24 items) Follow-up (Impressions of the experiment through free-form: 3 items)	
Annotations	Subjective evaluations of the interlocutor’s emotional valence at each time step of dialogue (Continuous values of 0 to 10 represent negative to positive emotional valence)	
Language	Japanese	

Table 1: Summary of collected dataset.

is a script-based human-human dialogue dataset containing speech, video, and facial information. SEMAINE (McKeown et al., 2011) is a corpus containing dialogues between computer graphics (CG) agents with different personalities and human subjects. D64 (Oertel et al., 2013) is a multi-party dialogue corpus designed to capture the natural reactions and emotions of the interlocutors.

The physiological signals are measured and quantified by sensors for physiological phenomena (such as heartbeat, brain waves, pulse, respiration, and perspiration) and can deal with the real-time state of the interlocutor. Several multimodal dialogue corpora have been constructed that include the physiological signals of the interlocutor in a dialogue. For example, RECOLA (Ringeval et al., 2013) is a human dialogue dataset that includes physiological signals during a collaborative dialogue task. Electrocardiogram (ECG) and electrodermal activity (EDA) are utilized as physiological signals in RECOLA. The PEGCONV dataset (Saffaryazdi et al., 2022) comprises discussion dialogues and includes galvanic skin response (GSR) and photoplethysmography (PPG) as physiological signals during the dialogue. Hazumi (Komatani and Okada, 2021) is a multimodal dialogue corpus containing dialogues between a human and a CG agent. The physiological signals include EDA, blood volume pulse (BVP), skin temperature (TEMP), and heart rate (HR) data (Katada et al., 2023).

Although several corpora have been constructed

in this way, none of the corpora contain data that comprehensively includes movement, gaze, and a variety of physiological signals. Moreover, to the best of our knowledge, there has been no research on estimating the real-time user state required by dialogue systems from sensor signals in dialogues.

### 3 Data Collection

The data were collected in two periods, with the first beginning in December 2023 and the second in February 2024. To distinguish the two datasets, we used the year and month of data collection for naming: “Data2312” for the data collected in December 2023 and “Data2402” for the data collected in February 2024.

In these two sets of data collection experiments, a total of 40 interlocutors (21 male, 19 female), all native Japanese speakers, participated. They were recruited from the general public by a recruiting agency, and each participated in only one data collection experiment. Two interlocutors were paired into one group and engaged in 10-minute dialogues on three different topics: “Chit-chat”, “Narrative”, and “Discussion”. Immediately after each dialogue, interlocutors annotated their subjective evaluations related to emotional valence while watching the recordings of the dialogue. A detailed summary of the collected dataset is provided in Table 1.

Both data collection experiments were conducted in the same sequence: pre-experiment questionnaire administration, sensor placement and attachment, dialogue and annotation conduction, sen-

sensor removal, and follow-up questionnaire administration. The dialogue and annotation conduction process was repeated three times for the three topics. For each topic, the following sequence was repeated: dialogue conduction, post-dialogue questionnaire administration, and subjective evaluation. The experiments were approved by the ethics committee of our institution.

In the following subsections, we describe in detail the multimodal data and heterogeneous sensors, the three dialogue topics, the questionnaires, and the subjective evaluation of dialogues.

### 3.1 Multimodal Data and Heterogeneous Sensors

We used heterogeneous sensors to collect multimodal data including speech, video, physiological signals, gaze information, and motion information. The data collection environment is shown in Fig. 1.

**Speech:** DPA 4088 uni-directional microphones were worn on the heads of each of the two interlocutors to capture audio recordings containing a single interlocutor’s voice. We used Azure Kinect’s (hereafter, Kinect) built-in omni-directional microphones to collect audio recordings containing the voices of two interlocutors. For the audio recordings collected by DPA 4088, the close proximity of the interlocutors and the loudness of the other interlocutor resulted in data containing faint sounds from the other interlocutor.

**Video:** We used Kinect, Logitech C920 Pro HD Webcam (hereafter, Logi webcam), and GoPro Hero 10 (hereafter, GoPro) to record the interlocutors’ behavior. In Data2312, two Kinets were placed between the interlocutors to record RGB and depth video of their upper bodies. A Logi webcam was positioned to the side of the interlocutors, capturing their full-body RGB video from a side view. In Data2402, to capture the full-body movements of the interlocutors instead of just the upper body, two Kinets were positioned between them to record separate full-body RGB and depth videos. Additionally, two GoPros were positioned in the same positions as the Kinets to record full-body RGB video. Because of the lack of clarity of the logi webcam, the third GoPro was placed to the side to capture two interlocutors’ full-body RGB video from the side. The Kinect and Logi webcam collected AVI files, while the GoPro recorded MP4 files.

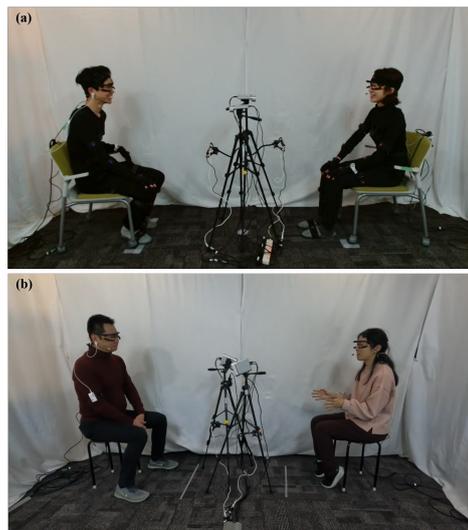


Figure 1: Data collection environments for (a) Data2312 and (b) Data2402. In (a) Data2312, two Kinets and two millimeter-wave sensors were placed between the interlocutors, and each interlocutor wore a set of wearable sensors. In (b) Data2402, two Kinets, two GoPros, and two millimeter-wave sensors were placed between the two interlocutors, and each interlocutor wore a set of wearable sensors.

**Physiological Signals:** We collected physiological signals during the dialogues by using the EmbracePlus<sup>1</sup> and the Shimmer3 GSR+<sup>2</sup>. The EmbracePlus is wireless and worn like a wristwatch. We used it to collect BVP, EDA, TEMP, and wrist acceleration (ACC). The Shimmer3 GSR+ can collect EDA using an optical pulse sensing probe attached to the finger and photoplethysmography (PPG) using either an ear clip or an optical pulse probe. Due to the greater stability of data collected through the former, we chose to collect PPG using the ear clip.

**Gaze:** The Pupil Core eye tracker<sup>3</sup> (hereafter, Pupil Core) was worn by the interlocutor like glasses and was used to collect gaze data, pupil information, and eye video during the dialogue.

**Motion:** We positioned two IWR1443 BOOST millimeter-wave sensors between the two interlocutors to capture 3D point cloud and motion data. In Data2312, we utilized two motion capture devices called Perception Neuron 3 Body Kit<sup>4</sup> (hereafter, PN3) on the interlocutors’ entire body to gather

<sup>1</sup><https://www.empatica.com/en-int/embraceplus/>

<sup>2</sup><https://shimmersensing.com/product/shimmer3-gsr-unit/>

<sup>3</sup><https://pupil-labs.com/products/core>

<sup>4</sup><https://neuronmocap.com/pages/perception-neuron-3>

	Device name	Type of sensor	Data
Devices for Data2312 and Data2402	Pupil Core	World camera Eye cameras	First-person video with gaze measurement Eye video, pupil and gaze information
	DPA 4088	Uni-directional mic	Audio containing one interlocutor’s voice
	Shimmer3 GSR+	Ear-mounted sensor	Photoplethysmography (PPG) 120 Hz
	EmbracePlus	Wristwatch sensor	Blood volume pulse (BVP) 64 Hz Electrodermal activity (EDA) 4 Hz Skin temperature (TEMP) 1 Hz Wrist acceleration (ACC) 64 Hz
	IWR1443BOOST	Millimeter-wave sensor	3D point cloud and motion data
Devices for Data2312	Azure Kinect	RGB camera Depth camera Omni-directional mic	Front upper body RGB video Front upper body depth video Audio containing two interlocutors’ voices
	Logi Webcam	RGB camera	Face-to-face full-body RGB video of two interlocutors
	Perception Neuron 3	IMU sensor	Skeleton hierarchy information and motion data
Devices for Data2402	Azure Kinect	RGB camera Depth camera Omni-directional mic	Front full-body RGB video Front full-body depth video Audio containing two interlocutors’ voices
	GoPro Hero 10	RGB camera	Front full-body recording Face-to-face full-body RGB video of two interlocutors

Table 2: Multimodal data collected from devices.

skeleton hierarchy information and motion data. Due to the time-consuming process of wearing and calibrating the PN3, as well as interference from numerous devices affecting the inertial measurement unit (IMU) sensor signals, in Data2402, we decided not to use the PN3 with the intention of extracting the interlocutors’ motion information from video recordings with image processing.

Two computers were used to acquire the sensor signals, which were streamed from each device. One computer served as a time server for ensuring synchronization of timestamps. Since EmbracePlus and GoPro do not support real-time streaming, the timestamps were synchronized post-data acquisition. Table 2 lists the devices and the multimodal data collected by them.

### 3.2 Dialogue Topics

To elicit a variety of mental states and gestures from the interlocutors, the following three topics were utilized. Example dialogues for each topic are provided in Table 3.

**Chit-chat:** Free dialogue with no restrictions on topics. Serving both as a means to collect dialogue in normal situations and as an icebreaker.

**Narrative:** The interlocutor’s own special episode. Storytelling provides a wealth of gestures (Colletta et al., 2010), and we can also expect that mental states will be expressed when discussing a cherished or distressing memory.

**Discussion:** Topics with different opinions for or against. We can expect negative mental states to be expressed during exchanges with an interlocutor who holds an opposing view. Conversely, we can also anticipate positive mental states to be experienced when the interlocutors reach an agreement. Fifteen topics were chosen from a site<sup>5</sup> that deals with discussion topics, which we then translated into Japanese. Before the data collection experiment, the pairs of interlocutors were asked about their opinions in favor of or against the 15 topics and the topics that they had different opinions about were selected as discussion topics.

### 3.3 Questionnaires

Each interlocutor completed a pre-experiment questionnaire before the start of the experiment, a per-dialogue questionnaire immediately after each dialogue, and a follow-up questionnaire after the end of the experiment. The details of the respective questionnaires are as follows.

**Pre-experiment Questionnaire:** Asking about the interlocutor’s demographic information and personality traits. Demographic information included gender, age, educational background, and employment status. For personality traits, we used a 10-item questionnaire from TIPI-J (Oshio et al., 2012) to measure the Big Five traits: openness to ideas/experience, conscientiousness, extraversion,

<sup>5</sup><https://www.procon.org/>

<p><b>“Chit-chat”:</b> Open-domain dialogue with no restrictions on topics.</p> <p>02F20: What did you have for breakfast?  02M30: I didn’t have breakfast.  02F20: You didn’t eat? Are you the type of person who only eats two meals a day?  02M30: One or two meals a day.  02F20: One meal a day!/? Which one do you eat? Breakfast, lunch, or dinner? I’m the type of person who eats three meals a day, because I often get hungry. So I envy those who only need one meal a day.  02M30: But I may put three meals into one.</p>
<p><b>“Narrative”:</b> Own personal story that you can’t help but want to tell others about.</p> <p>04M20: I have done something that people often say is unusual.  04F30: I would like to hear about it.  04M20: People learn various sorts of things, don’t they? Like piano, swimming. I’m often told that the thing I learned was unusual.  04F30: What was it?  04M20: I used to study Kabuki.  04F30: Huh? Amazing!  04M20: That was from grade six to about high school.  04F30: You were doing it for quite a long time.</p>
<p><b>“Discussion”:</b> Is obesity a disease?</p> <p>08M50: I’d like to start by defining the term “obesity”.  08M20: I agree.  08M50: What counts as obesity?  08M20: I’m sorry if I’m being a bit light-hearted here, but in short, a fat person. I don’t mean exactly how many kilos or more he weighs, but in terms of his appearance, someone who has a bit of a belly.  08M50: Obesity is generally expressed as a certain value, such as BMI, and that value is considered to be equate to obesity. But I don’t think that certain values equal poor health or disease. What do you think about that?  08M20: I can totally understand. To be honest, I’m not sure if it’s correct or not, because it’s hard to connect a value to disease.</p>

Table 3: Dialogue excerpts on “Chit-chat”, “Narrative”, and “Discussion”. Interlocutor IDs are five characters of the form “NNGAA”, where “NN” is the group number, “G” is the gender of the interlocutor (“M” for male, “F” for female), and “AA” is the age of the interlocutor. These excerpts were translated from the original Japanese to English by the authors.

agreeableness, and emotional stability (Goldberg, 1990).

**Post-dialogue Questionnaire:** Asking about the quality of the dialogue and interlocutors’ impressions on a 7-point scale. It consists of 24 items in total. For the evaluation items relating to the quality of the dialogue, we used the same six items as the questionnaire by Yamashita et al. (2023). For the evaluation items related to the impressions of the dialogue, we used 18 items from the measurement items regarding the interpersonal communication cognition of the interlocutors (Kimura et al., 2005). The items of the post-dialogue questionnaire are shown in Table 4.

**Follow-up Questionnaire:** A free-form questionnaire asking about the content of the dialogue that left an impression on interlocutors, any issues the interlocutors encountered during the experiment, and their opinions and impressions of the overall experimental process.

### 3.4 Subjective Evaluations

To obtain the interlocutors’ real-time subjective evaluation for emotional valence, each interlocutor annotated the emotional valences of the dialogue immediately after the end of each dialogue. Continuous values of 0 to 10 were used, where 0 represents very negative, 5 represents neutral, and 10 represents very positive.

To reproduce the dialogue scene and to help the interlocutors recall their mental state at the time, we used video recordings of the other interlocutor as the annotation videos, rather than their own video recordings. Specifically, the interlocutor used the annotation software CARMA (Girard, 2014) and assigned a numerical value that was considered appropriate for “their mental state” at each time in the dialogue while watching the video recording. A screenshot of the CARMA interface is shown in Fig. 2. The sampling rate of annotations was 4 Hz. To familiarize the interlocutors with the use of the annotation software, a five-minute annotation exercise was conducted before the start of data collection.

Dialogue Qualities
1. The dialogue partner was approachable.
2. The dialogue partner’s speech was informative.
3. The dialogue partner’s speech was easy to understand.
4. I was satisfied with the dialogue.
5. I was interested in the topics discussed in this dialogue.
6. I took the initiative to speak.
Dialogue Impressions
1. I was able to coordinate the conversation well.
2. I was bored with the conversation.
3. The conversation proceeded cooperatively.
4. The conversation was harmonious.
5. The conversation was unsatisfactory.
6. The conversation was slow-paced.
7. The conversation went cold.
8. The conversation was awkward.
9. I was absorbed in the conversation.
10. The conversation lacked focus.
11. The partner and I talked with great interest.
12. The conversation was tense.
13. The conversation was friendly.
14. The conversation was lively.
15. The conversation was positive on both sides.
16. The conversation was boring.
17. The conversation was worthwhile.
18. The conversation was drawn out.

Table 4: Items of the Post-dialogue questionnaire, where “Items enquiring about the quality of the dialogue” refers to (Yamashita et al., 2023) and “Items enquiring about the impressions of the dialogue” refers to (Kimura et al., 2005). The questionnaire was translated from the original Japanese to English by the authors

## 4 Data Analysis

Human emotional mental states, such as happiness and sadness, are formed through the brain’s processing of information from three sources: 1) information from the body (e.g., HR, sweating, and other physiological states), 2) information from the external world (e.g., visual and auditory input, etc.), and 3) memories stored in the brain (Damasio, 1996; Moriguchi and Komaki, 2013). In our collection of multimodal data, the physiological signals obtained from EmbracePlus, Shimmer3 GSR+, and Pupil Core (e.g., EDA, PPG, pupil diameter) captured the interlocutors’ physiological states (i.e., information from the body), while subjective evaluations annotated the emotional valence of the interlocutors.

In this study, EDA and BVP (collected from EmbracePlus), PPG (collected from Shimmer3 GSR+), and pupil diameter (collected from Pupil Core) were used as physiological signals. We first performed data preprocessing on these signals for subsequent analysis. We then performed Granger causality analysis to examine the relationship between these physiological signals and subjective



Figure 2: Screenshot of annotation interface. Emotional valence is assigned by manipulating the slide bar on the right of the screen using the controller while the interlocutor watches the other interlocutor’s recording.

evaluations of emotional valence, i.e., whether these physiological signals can be used to predict subjective evaluations. Finally, we analyzed the differences between these physiological signals under different polarities of valence.

### 4.1 Data Preprocessing

We extracted the physiological signals of the interlocutor during the dialogue on the basis of the start and end times of the dialogue using timestamps.

For EDA, BVP, and PPG, we used the NeuroKit2 toolbox<sup>6</sup> for data preprocessing (denoising, filtering) and feature extraction. We extracted the tonic skin conductance level (SCL) and phasic skin conductance response (SCR) from the EDA. SCL, also known as tonic, measures the overall conductivity of the skin, which reflects the general level of sweat gland activity. SCR, also known as phasic, measures the rapid changes in skin conductivity that occur in response to specific stimuli. The rate (the HR as measured on the basis of PPG/BVP peaks), peak (represents the highest point of PPG/BVP, used as an indicator of the intensity of a heartbeat), and the R-R intervals (RRI, which reflect the changes in time between heartbeats, i.e., HR variability) were calculated from the raw BVP and PPG data. The subjective evaluation annotations and physiological signals of an interlocutor during one minute of dialogue are shown in Fig. 3.

Pupil diameter data was sampled at a rate of 13–26 Hz, collected by Pupil Core, and the actual size of the pupil diameter (unit: mm) was derived by the device’s built-in algorithm. Each timestamp has a “confidence” value indicating the quality of the measurement, and data with a confidence > 0.6

<sup>6</sup><https://neuropsychology.github.io/NeuroKit/>

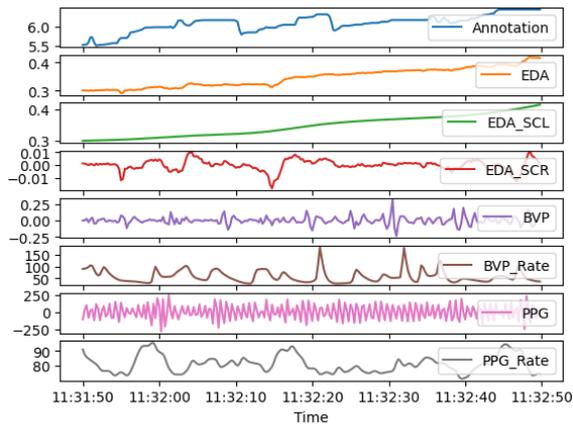


Figure 3: One-minute subjective evaluations and physiological signals of an interlocutor during dialogue. From top to bottom: Subjective evaluations (Annotation), EDA, EDA SCL, EDA SCR, BVP, BVP Rate, PPG, and PPG Rate.

is considered reliable. Since each individual has a different pupil diameter, we normalized all pupil diameter data to 0–1 using the Min-Max normalization.

#### 4.2 Granger Causality Analysis

We computed the Granger causality analysis (Granger, 1969) to identify physiological signals or specific features of these signals that are most indicative of emotional valence changes. This method is used to evaluate the predictive utility of one variable for forecasting another and is also employed to explore the relationship between physiological signals and mental states like emotions (Gao et al., 2020). A time series  $X$  is considered to Granger-cause another time series  $Y$  if past values of  $X$  and  $Y$  predict  $Y$  significantly better than past values of  $Y$  alone (Granger, 1969).

In this study, the null hypothesis is that the physiological signals or specific features of these signals fail to Granger-cause changes in emotional valence.

The two time series used for Granger causality analysis need to be aligned and have the same sampling rate, so as the first step, we resampled all physiological signal features (SCL, SCR for EDA, Rate, Peak, RRI for PPG and BVP, and pupil diameter for left and right eyes) such that they had the same sampling rate as that of the subjective evaluations at 4 Hz, and then aligned all the data in accordance with the timestamps.

In addition, the Granger causality test assumes the series to be stationary and linearly related to make valid results. We therefore con-

ducted the Augmented Dickey Fuller (ADF) and Kwiatkowski-Phillips-Schmidt-Shin (KPSS) tests for stationarity and calculated Pearson correlation for linear relationship assessment. According to the results of ADF and KPSS tests, most emotional valence data were nonstationary. Consequently, we utilized the Toda-Yamamoto causality test<sup>7</sup>, which is an adaptation of the Granger causality test suitable for nonstationary data (Toda and Yamamoto, 1995). Regarding the Pearson correlation results, the pair of EDA SCR and emotional valence exhibited almost no linear relationship. In contrast, SCL for EDA, Rate, Peak, RRI for PPG and BVP, and pupil diameter for left and right eyes showed weak linear relationships with emotional valence. Therefore, we conducted subsequent causality tests only on the pairs involving SCL for EDA, Rate, Peak, RRI for PPG and BVP, and pupil diameter for both left and right eyes with emotional valence.

For the Granger causality test, including the Toda-Yamamoto test, the parameter “lag” represents the number of time delays used in predicting future time series data from past time series data. We set the maxlag to 8, corresponding to a maximum time delay of 2 seconds, since all data are sampled at 4 Hz. The analyses were computed for all lags up to maxlag.

#### 4.3 Comparison of Physiological Signal Means under Positive and Negative Valence

To investigate the differences in physiological signal features depending on emotional valence polarities, we performed the Wilcoxon rank-sum test, which is a nonparametric test also known as Mann–Whitney U test, between the means of SCL and SCR of EDA, Rate, Peak, RRI of PPG and BVP, and pupil diameter for both left and right eyes. We conducted the tests with the null hypothesis that two samples of physiological signal features are drawn from the same distribution under the “positive” and “negative” emotional valence. Before conducting the test, we processed the subjective evaluation annotations; we segmented all annotations into 10-second intervals and calculated the average emotional valence value for each segment. The averages greater than 5.5 and less than 4.5 were categorized as “positive” and “negative”, respectively. Those emotional values with averages in the 4.5–5.5 range were considered “neutral emo-

<sup>7</sup><https://github.com/nicolarighetti/Toda-Yamamoto-Causality-Test>

tional states” and were therefore excluded from this analysis.

Given the variability of physiological signals among different interlocutors, we normalized EDA, BVP, PPG, and pupil diameter for each interlocutor before feature extraction. Specifically, we processed EDA, BVP, and PPG with the Z-score normalization, and pupil diameter using the Min-Max normalization. Then, as mentioned in Section 4.1, we extracted SCL and SCR from EDA and extracted Rate, Peak, and RRI from PPG and BVP. Finally, we performed the Wilcoxon rank-sum tests on the means of SCL and SCR of EDA, Rate, Peak, RRI of PPG and BVP, and pupil diameter for both left and right eyes, comparing them between the “positive” and “negative” emotional valence.

## 5 Results

In this section, we present the results of the Granger causality analysis and the differences in the means of physiological signals between “positive” and “negative” valence.

### 5.1 Results of Granger Causality Analysis

On the basis of the dialogue topics, we grouped the collected data into four sets: “Chit-chat” (40 dialogues), “Narratives” (40 dialogues), “Discussions” (40 dialogues), and all types of dialogues (120 dialogues). The Toda-Yamamoto Granger causality tests were conducted between the features of EDA, BVP, PPG, pupil diameter, and subjective evaluation annotations (e.g., between the RRI of the PPG and the subjective evaluation annotations) across these four datasets, with maxlag of 8. During the causality testing between EDA SCL and emotional valence, we encountered issues with idiosyncratic ranks, which prevented the construction of the model for the causality test. As a result, we excluded causality analyses between EDA SCL and emotional valence.

The proportion of dialogues featuring Granger causality is shown in Table 5. The results show that the PPG Rate has the highest potential to predict the interlocutor’s emotional valence in all dialogues. We also found that PPG Rate is the most useful feature for predicting emotional valence in all three topics :“Chit-chat”, “Narrative”, and “Discussion”.

Signal	Feature	Chit-chat	Narrative	Discussion	All
BVP	Rate	.28	.23	.08	<b>.19</b>
	Peak	.05	.10	.18	.11
	RRI	.10	.18	.08	.12
PPG	Rate	<u>.45</u>	<u>.48</u>	<u>.30</u>	<b>.41</b>
	Peak	.10	.13	.18	.13
	RRI	.10	.18	.08	.12
Pupil diameter	Left	.13	.10	.13	.12
	Right	.18	.15	.10	<b>.15</b>

Table 5: Results of Toda-Yamamoto causality tests (maxlag = 8). The proportion of dialogues with a significant difference of  $p < 0.05$  in “Chit-chat” (40 dialogues), “Narrative” (40 dialogues), “Discussion” (40 dialogues), and “All dialogues” (120 dialogues) is shown. Bold numbers are the highest proportion for the BVP, PPG, and pupil features, and underlined numbers are the highest proportion for each topic of dialogue.

### 5.2 Results of the Differences between Physiological Signal Means under Different Valence Polarities

We performed the Wilcoxon rank-sum tests on SCL, SCR of EDA, Rate, Peak, RRI of PPG and BVP, and pupil diameter for left and right eyes under “positive” and “negative” emotional valence. Note that the ratio of sample size between “positive” and “negative” is around 6:1.

The mean, standard deviation, and results of the Wilcoxon rank-sum test are shown in Table 6. The results indicate that the means of EDA SCR, BVP Rate, and BVP RRI were significantly different ( $p < 0.05$ ) between the “positive” and “negative” valence. Specifically, our experimental results showed significant differences in EDA SCR under different emotional valences, but not in EDA SCL. This may be because SCR captures instantaneous changes in the skin and is more responsive to short-term emotional responses, whereas SCL reflects slower changes in the skin and is more indicative of longer-term emotional states. Additionally, in the “positive” valence during the dialogue, RRI values (i.e., the interval between heartbeats) are generally higher than in the “negative” valence, and the variability of RRI is also greater. This is probably because positive emotional states such as relaxation and contentment are associated with a slower heart rate, resulting in increased RRI values. Conversely, negative emotional states, such as anxiety and stress, are generally linked to a faster heart rate and consequently shorter RRI values.

## 6 Conclusion and Future Work

In this study, we collected dialogue data containing comprehensive multimodal data and subjective

Signal	Feature	Positive		Negative		p-value
		Mean	Std	Mean	Std	
EDA	SCL	1.3e-4	0.01	8.0e-3	0.02	0.229
	SCR	1.38	1.70	1.29	1.71	1.3e-7**
BVP	Rate	68.4	13.4	69.0	12.8	0.042**
	Peak	10.2	2.16	10.3	2.07	0.053
	RRI	944	216	933	209	0.045**
PPG	Rate	85.2	13.7	85.2	10.5	0.644
	Peak	13.2	2.13	13.2	1.67	0.765
	RRI	726	118	718	88.3	0.690
Pupil diameter	Left	0.38	0.23	0.41	0.25	0.115
	Right	0.44	0.24	0.42	0.25	0.070

Table 6: Mean, standard deviation (Std), and the p-value of the Wilcoxon rank-sum test for means under positive and negative valence (\*\* $p < 0.05$ ).

evaluations at each time step during the dialogue using heterogeneous sensors. Through our analysis of the relationship between physiological signals and emotional valence using the Granger causality analysis, we identified several physiological signals that could be useful for predicting real-time emotional valence. We also clarified how physiological signals differ depending on the “positive” or “negative” polarity of the valence.

However, several limitations of our study should be acknowledged. First, the relatively small sample size limits the statistical power of our findings and reduces the generalizability of the results to a larger population. Second, the imbalanced ratio of positive to negative valence samples may potentially lead to biased conclusions about the relationship between physiological signals and emotional valence.

Future research needs to apply methodologies for analyzing imbalanced and small sample size data. Moreover, we plan to expand our analysis to include sensor signals, linguistic information, questionnaires about personality traits, and impressions of the dialogue, in addition to physiological signals. We will also use the information from the sensors to predict emotional valence in real-time. Ultimately, our goal is to achieve a dialogue system capable of estimating and appropriately responding to the user’s mental state in real-time.

## Acknowledgments

This work was supported by JST Moonshot R&D Grant Number JPMJMS2011.

## References

Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N

Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42:335–359.

Jean-Marc Colletta, Catherine Pellenq, and Michèle Guidetti. 2010. Age-related changes in co-speech gesture and narrative: Evidence from french children and adults. *Speech Communication*, 52(6):565–576.

Antonio R Damasio. 1996. The somatic marker hypothesis and the possible functions of the prefrontal cortex. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 351(1346):1413–1420.

Ellen Douglas-Cowie, Cate Cox, Jean-Claude Martin, Laurence Devillers, Roddy Cowie, Ian Sneddon, Margaret McRorie, Catherine Pelachaud, Christopher Peters, Orla Lowry, et al. 2011. The HUMAINE database. *Emotion-oriented Systems: The Humaine Handbook*, pages 243–284.

Yunyan Gao, Xiangkun Wang, Thomas Potter, Jianhai Zhang, and Yingchun Zhang. 2020. Single-trial EEG emotion recognition using Granger Causality/Transfer Entropy analysis. *Journal of Neuroscience Methods*, 346:108904.

Jeffrey M Girard. 2014. CARMA: Software for continuous affect rating and media annotation. *Journal of Open Research Software*, 2(1):e5.

Lewis R Goldberg. 1990. An Alternative “description of personality”: The Big-Five Factor Structure. *Journal of Personality*, 59(6):1216–1229.

Clive WJ Granger. 1969. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica*, 37:424–438.

Shun Katada, Shogo Okada, Yuki Hirano, and Kazunori Komatani. 2020. Is She Truly Enjoying the Conversation? Analysis of Physiological Signals toward Adaptive Dialogue Systems. In *Proc. ICMI*, page 315–323.

Shun Katada, Shogo Okada, and Kazunori Komatani. 2023. Effects of Physiological Signals in Different Types of Multimodal Sentiment Estimation. *IEEE Transactions on Affective Computing*, 14(3):2443–2457.

Masaki Kimura, Makio Yogo, and Ikuo Daibo. 2005. Expressivity halo effect in the conversation about emotional episodes. *Japanese Journal of Research on Emotions*, 12(1):12–23. (in Japanese).

Kazunori Komatani and Shogo Okada. 2021. Multimodal human-agent dialogue corpus with annotations at utterance and dialogue levels. In *Proc. ACII*, pages 1–8.

Gary McKeown, Michel Valstar, Roddy Cowie, Maja Pantic, and Marc Schroder. 2011. The semaine database: Annotated multimodal records of emotionally colored conversations between a person and a limited agent. *IEEE Transactions on Affective Computing*, 3(1):5–17.

Trisha Mittal, Uttaran Bhattacharya, Rohan Chandra, Aniket Bera, and Dinesh Manocha. 2020. M3er: Multiplicative multimodal emotion recognition using facial, textual, and speech cues. In *Proc. AACL*, pages 1359–1367.

Yoshiya Moriguchi and Gen Komaki. 2013. Neuroimaging studies of alexithymia: Physical, affective, and social perspectives. *BioPsychoSocial Medicine*, 7(1):1–12.

- Catharine Oertel, Fred Cummins, Jens Edlund, Petra Wagner, and Nick Campbell. 2013. [D64: A corpus of richly recorded conversational interaction](#). *Journal on Multimodal User Interfaces*, 7(1-2):19–28.
- Atsushi Oshio, Shingo Abe, and Pino Cutrone. 2012. [Development, reliability, and validity of the Japanese version of Ten Item Personality Inventory \(TIPI-J\)](#). *Japanese Journal of Personality/Pasonariti Kenkyu*, 21(1):40–52.
- Fabien Ringeval, Andreas Sonderegger, Juergen Sauer, and Denis Lalanne. 2013. [Introducing the RECOLA multimodal corpus of remote collaborative and affective interactions](#). In *Proc. 2013 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, pages 1–8.
- James A Russell. 2003. [Core affect and the psychological construction of emotion](#). *Psychological review*, 110(1):145–172.
- Nastaran Saffaryazdi, Yenushka Goonesekera, Nafiseh Safaryazdi, Nebiyu Daniel Hailemariam, Ebaso Girma Temesgen, Suranga Nanayakkara, Elizabeth Broadbent, and Mark Billinghurst. 2022. [Emotion Recognition in Conversations Using Brain and Physiological Signals](#). In *Proc. IUI*, page 229–242.
- Hiro Y Toda and Taku Yamamoto. 1995. [Statistical inference in vector autoregressions with possibly integrated processes](#). *Journal of econometrics*, 66(1-2):225–250.
- Sanae Yamashita, Koji Inoue, Ao Guo, Shota Mochizuki, Tatsuya Kawahara, and Ryuichiro Higashinaka. 2023. [RealPersonaChat: A Realistic Persona Chat Corpus with Interlocutors’ Own Personalities](#). In *Proc. PACLIC*, pages 852–861.
- Zhou Yu, Alexandros Papangelis, and Alexander Rudnicky. 2015. [TickTock: A Non-Goal-Oriented Multimodal Dialog System with Engagement Awareness](#). In *Proc. AAAI Spring Symposium on Turn-taking and Coordination in Human-Machine Interaction*, pages 108–111.

# The Gap in the Strategy of Recovering Task Failure between GPT-4V and Humans in a Visual Dialogue

Ryosuke Oshima<sup>1</sup> Seitaro Shinagawa<sup>2</sup> Shigeo Morishima<sup>3</sup>

<sup>1</sup>Waseda University <sup>2</sup>Nara Institute of Science and Technology

<sup>3</sup>Waseda Research Institute for Science and Engineering

ryosukeoshima@fuji.waseda.jp, sei.shinagawa@is.naist.jp, shigeo@waseda.jp

## Abstract

Goal-oriented dialogue systems interact with humans to accomplish specific tasks. However, sometimes these systems fail to establish a common ground with users, leading to task failures. In such cases, it is crucial not to just end with failure but to correct and recover the dialogue to turn it into a success for building a robust goal-oriented dialogue system. Effective recovery from task failures in a goal-oriented dialogue involves not only successful recovery but also accurately understanding the situation of the failed task to minimize unnecessary interactions and avoid frustrating the user. In this study, we analyze the capabilities of GPT-4V in recovering failure tasks by comparing its performance with that of humans using Guess What?! Game. The results show that GPT-4V employs less efficient recovery strategies, such as asking additional unnecessary questions, than humans. We also found that while humans can occasionally ask questions that doubt the accuracy of the interlocutor’s answer during task recovery, GPT-4V lacks this capability.

## 1 Introduction

Goal-oriented dialogue systems work with humans on tasks to achieve a goal (de Vries et al., 2017; Kottur et al., 2021; Ma et al., 2022). They do not always succeed in their tasks in one shot due to the failure to establish a common ground of dialogue (Clark, 1996) with their human interlocutors. The task failure occurs due to various factors, including human error (Oshima et al., 2023), system error (Hudeček and Dusek, 2023; Mazuecos et al., 2021), and misunderstandings between the two (Paek and Horvitz, 2000).

In human-to-system dialogue, it is important for humans to finally achieve a successful goal regardless of the factor of task failure that occurs along the way. In this case, the system needs the capability to continue the failure dialogue and cooperatively

recover from the task failure rather than terminating the dialogue (Benotti and Blackburn, 2021a). For example, suppose a task where an interactive autonomous driving system and a user tackle the task of going to an interior shop. The task may fail due to unexpected events, such as the destination being closed for construction or the user miscommunicating the desired location (Ma et al., 2022). In these cases, the system must offer alternatives or confirm the user’s statements to recover from the task failure.

Of course, successful recovery from failure is not the only requirement for this dialogue task. As a goal-oriented dialogue, the system also demands minimizing the number of interactions to avoid frustrating the user. The system should have a dialogue strategy that makes good use of the information in the failure dialogue history and efficiently recovers the task.

While the recent Vision-Language Models (VLMs) integrated with Large Language Models (LLMs) have garnered attention for their ability to solve tasks at a high level through dialogue (OpenAI, 2024; Liu et al., 2023b), the performance of these VLMs in “failure task recovery” remains unclear. Investigating and analyzing these models’ failure task recovery capabilities can lead to the development of robust dialogue systems for real-world applications. For example, instead of relying solely on VLM for the entire task recovery process, we can enhance the system’s overall performance by implementing rule-based modules and preprocessing VLM inputs to compensate for VLM’s weaknesses.

In this paper, we analyze the VLM’s ability to recover the course of the dialogue as a first step toward the goal of building a system that can efficiently return to success after a task failure. We consider a problem setting in which the system performs a recovery action after a goal-oriented visual dialogue with a human interlocu-

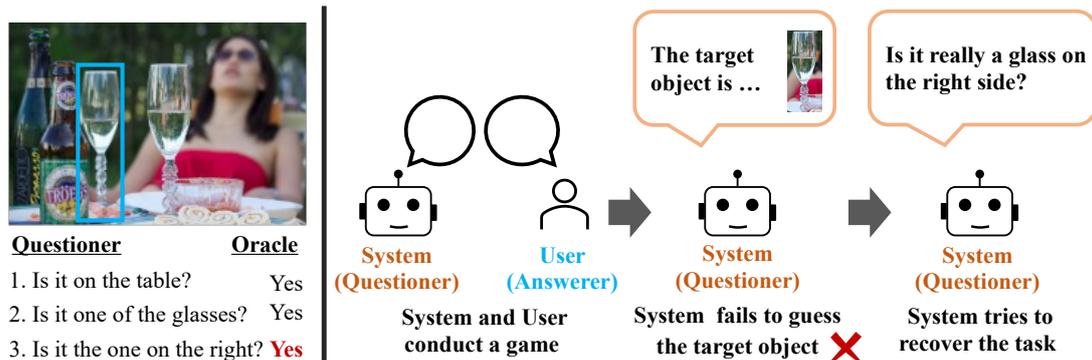


Figure 1: The *left* figure shows a failed game at Guess What?! Game by human annotators, which is included in Guess What?! Dataset (de Vries et al., 2017). The target object is outlined in blue. In this example, the questioner failed to correctly guess the target object due to the third wrong answer by the answerer. The *right* figure overviews the failure recovery task. A system and a user conducted the game, but the system guessed a different object. In this example, the system attempts to recover from a task failure to success by asking an additional question.

tor fails. Specifically, we leverage Guess What?! Game (de Vries et al., 2017), which is a widely used game of goal-oriented visual dialogue. Guess What?! Game (de Vries et al., 2017) is appropriate for a study of failure task recovery in dialogue because the goal of this game is building common ground directly.

In experiments, we prepare 100 task failure samples in Guess What?! Game (de Vries et al., 2017) and analyze the recovery capabilities of VLM by comparing them in human-to-system dialogue. Our experiments reveal that VLM struggles with recovering from task failure in a goal-oriented dialogue. It frequently performs unnecessary recoveries and uses ineffective repair utterances. Furthermore, humans tend to check the reliability of previous answers when errors are suspected. On the other hand, GPT-4V does not question prior utterances and fails to generate repair utterances that express doubt.

## 2 Related Work

### 2.1 Repair in Dialogue

*Repair* is one of the key interactional mechanisms to achieve shared understanding and coordination when miscommunication has occurred (Schegloff et al., 1977; Schegloff, 1992; Purver et al., 2018). The construction of a robust dialogue system that can recognize and use repair has been discussed because this miscommunication also occurs in human-machine conversation (Purver et al., 2018; Balaraman et al., 2023; Shaikh et al., 2024). In visual goal-oriented dialogue research field, clarification requests has been mainly discussed as a

repair utterance (Benotti and Blackburn, 2021b; Shi et al., 2022; Deng et al., 2023; Madureira and Schlagen, 2023; Chiyah-Garcia et al., 2023). Some research (Shi et al., 2022; Deng et al., 2023) deal with clarification questions for disambiguation in Minecraft game, where the system interacts with the user in the task of moving and building blocks according to the user’s instructions. Chiyah-Garcia et al. (2023) used SIMMC dataset to analyze what information is important for a shopping assistant in a virtual shop to interpret a user’s clarification requests. These studies focus on building a system that can perform or understand recovery “during” the dialogue to complete the task successfully the first time.

Although making clarification requests is a well-known dialogue strategy to avoid miscommunication (San-Segundo et al., 2001; Benotti and Blackburn, 2021b), it is hard for a system and humans to achieve a successful goal without failures all the time. In this paper, we consider the problem setting where once the dialogue is over and the task has ended in failure, how to turn it into a success. It is noted that the commonly used concept of “repair” or “repair utterance” in Conversational Analysis forms part of the recovery task and corresponds to Step 2 in the recovery flow introduced later (§3.2).

### 2.2 Recover in Tasks Other than Dialogue

Here, we describe previous works on addressing task failures and converting them into successes in non-dialogue tasks. Huang et al. (2022); Wang et al. (2023) worked on a task where a robot follows human instructions. When the robot failed to execute the instructions, Huang et al. (2022); Wang

et al. (2023) utilized the LLM’s strong reasoning abilities to correct the failures and achieve success. Huang et al. (2023); Fan et al. (2023); Zhang et al. (2023) focused on Automated Program recovery (APR), which aims to automatically fix software bugs and errors in programming.

These works focus on recovery that occurs solely within the systems. In contrast, dialogue recovery requires the systems to cooperatively interact with humans, presenting two main challenges. First, a system needs to minimize interactions to avoid frustrating the user. Second, the system must understand that human response errors cause task failure and not place too much trust in past dialogues (Os-hima et al., 2023). Given the complexity of recovery tasks after a task failure, it is important to conduct a detailed analysis of VLMs’ capabilities.

### 3 Guess What?! Game and Failure Recovery Task

#### 3.1 Guess What?! Game

Guess What?! Game (de Vries et al., 2017) is a two-player game in which a questioner asks yes or no questions to identify a target object, and an answerer<sup>1</sup> answers those questions. We don’t use other visual goal-oriented dialogue tasks (Kottur et al., 2021; Haber et al., 2019). This is because these dialogue scenarios involve too detailed object positioning within images (e.g., “Do you like the second sweater from the right in the bottom row?”) or require recognition of multiple images simultaneously, which VLMs generally perform poorly (Wu and Xie, 2023; Yang et al., 2023).

#### 3.2 Failure Recovery Task Definition

In this study, we consider the situation where the questioner and the answerer played a game, but the questioner failed to predict the target object (task failure). In such a case, the questioner should follow up with the user to ensure they can recover the task successfully from failure. We focus on this recovery process after task failure once this game is over. Figure 1 shows provides an overview of the recovery task. This recovery task requires a high success rate and emphasizes more efficient recovery. It is desirable to achieve recovery with as few additional questions as possible and, if feasible, without any additional questions.

<sup>1</sup> de Vries et al. (2017) calls an answerer “Oracle”. Instead, we use “answerer” in this paper to avoid misunderstanding because an interlocutor can make mistakes and not always give a perfect answer.

Figure 2 shows a detailed flow of the questioner’s (system side) recovery. The questioner can only choose one from two actions: 1) asking a question or 2) guessing a target object. The questioner performs the recovery in four steps.

Step 1: The questioner determines additional questions or re-prediction of objects based on information from the failure game. This step corresponds to the dialogue act classification module in goal-oriented dialogue. Given an image  $I$ , a dialogue history  $H = (\underbrace{(Q_1, A_1)}_{H_1}, \dots, \underbrace{(Q_{t-1}, A_{t-1})}_{H_{t-1}})$ , and

an object of failed prediction  $O_f$ , the questioner determines action  $A$  ( $a_1$ : asking an additional question or  $a_2$ : re-guessing the target object).

Step 2: If the questioner determines to ask an additional question ( $a_2$ ), it asks the question  $Q_t$  according to the image  $I$ , the dialogue history  $H$ , and the object of the failed prediction  $O_f$ .

Step 3: The questioner judges if the target object could be uniquely determined by an additional question and the answerer’s answer  $A_t$ . If the questioner judged unique, it proceeds to Step 4; if not, returns to Step 2. The dialogue history now contains additional questions and answers,  $H \rightarrow H' = (\underbrace{(Q_1, A_1)}_{H_1}, \dots, \underbrace{(Q_t, A_t)}_{H_t})$ .

Step 4: If the questioner didn’t ask any additional questions, it guesses the target object using the dialogue history  $H$ . If the questioner asks an additional question, it guesses the target object using the updated dialogue history  $H'$ .

These four steps are divided into two main parts: the decision to make a repair utterance and its actual execution (Steps 1 and 2) and the ability to correctly understand and process the repair utterance (Steps 3 and 4). In this study, we analyze the outcomes of recovery in Guess What?! Game (de Vries et al., 2017), where all four reasoning abilities are challenged at once in Section 5.1, and then we focus our analysis on the first two steps in Section 5.2.

## 4 Experiments

In this study, we analyze the success rate and features of the failure recovery task for humans and GPT-4V<sup>2</sup> (OpenAI, 2024). We first collect the failure game of Guess What?! (de Vries et al., 2017) to investigate failure recovery capability. Then, we

<sup>2</sup>We used the GPT-4 Turbo API through all the experiments in this study.

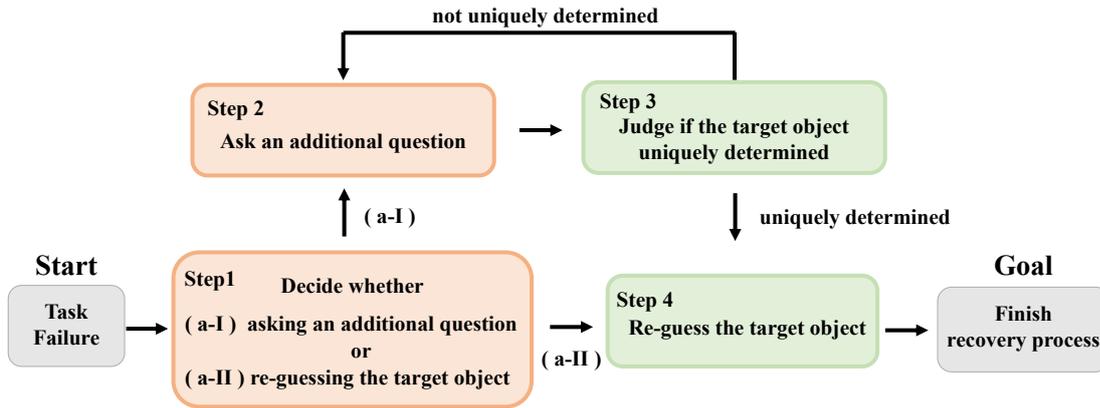


Figure 2: Questioner’s recovery flow. The recovery process is divided into two main parts: Step1, 2 and Step3, 4.



Figure 3: Failure game example

collect human and GPT-4V (OpenAI, 2024) failure task recoveries for the collected task failure game.

#### 4.1 Failure Game Collection

We simulated human-to-system dialogues in Guess What?! Game using GPT-4V to collect data. Specifically, we set up GPT-4V as a questioner and answerer and collected failed games by having them play a game. While the ideal scenario would involve a human as the answerer, our preliminary experiments demonstrated that GPT-4V is sufficiently capable of playing Guess What?! Game. This led us to adopt the method of using GPT-4V in both roles for this study.

Figure 3 shows an example of a generated failure game. The questioner failed to narrow down the target object with five questions and could not guess it accurately. The following describes the details of the GPT-4V roles for the questioner and answerer (See Appendix A.1 for overview diagrams of these models).

**Questioner’s Role** The questioner’s role is divided into two parts: a model that makes questions (called a questioner model) and a model that guesses the target objects (called a guesser model), as de Vries et al. (2017)’s proposed baseline model. The questioner model inputs a game image and di-

alogue history and outputs a question. The guesser model takes an image with numbers assigned to objects by SoM (Yang et al., 2023) (called SoM-image) and dialogue history as inputs and outputs the number of the target object. We applied SoM to the input images of the guesser model because GPT-4V has better inference ability with number assignment images than with understanding Visual Prompt (Yang et al., 2023), and it is impossible to output a target object’s bounding box<sup>3</sup>.

**Answerer’s Role** The answerer model takes SoM-image, a dialogue history, and the number of the correct object as inputs and outputs a yes/no answer.

**Game Collection Details and Results** According to (de Vries et al., 2017), the guesser model tried guessing the target object after the questioner and answerer models exchanged questions and answers five times. We sampled 815 pairs of images and target objects from the Guess What?! dataset’s test data. Then, we excluded any target objects that were too small or positioned at the edges of the images, as recognizing these objects demands high image recognition capabilities beyond the scope of our study. As a result, we collected 100 failed games. The collected games include samples where GPT-4V, acting as the answerer, made errors, resulting in failed games. We also adopted these samples as examples that simulate actual human-to-system dialogues because humans can also make mistakes in their answers due to misinterpretations or unintentional mistakes (Oshima et al., 2023).

<sup>3</sup>When GPT-4V takes a prompt “Output the human bounding box.” and an image as inputs, it returns an unreliable bounding box or says “I’m unable to directly output bounding boxes or any form of visual annotations”.

## 4.2 Failure Task Recovery Collection

We collected samples on the recovery tasks performed by humans and GPT-4V (OpenAI, 2024), using the collected task-failed games in Section 4.1. Specifically, we conducted three experiments; *GPT-4V-all*, *GPT-4V-Q*, and *Human-all* experiment. Table 1 presents the relationships among these experiments. These experiments vary depending on who is responsible for each step of the questioner’s recovery flow, which is introduced in Section 3.2. We describe the details of the three experiments below.

### 4.2.1 Human-all Experiment

We conducted an experiment to collect human recovery actions in failed games (called Human-all experiment). In collecting human recovery actions, two annotators each assumed the roles of questioner and answerer. The annotator in the questioner role worked with the answerer to address and recover from game failures, using the details of the failed task (game image, first predicted objects, and dialogue history from the failed game). The annotator in the answerer role received information about the details of the failed task and the correct target object. We created a demo application to collect humans’ recovery actions. Humans were monitored to ensure they were not cheating and diligently working on tasks. There are other ways to collect data through crowdsourcing, but we did not employ them in this case because they are fraught with problems, such as using Large Language Models (Veselovsky et al., 2023).

The data collection had 12 native Japanese speaker participants, each performing recovery actions for 25 games. We assigned 25 game recovery tasks to each annotator using a collection of 100 failed games in Section 4.1. This means that three annotators worked on the recovery task for each game, resulting in 300 recovery samples collected in total.

### 4.2.2 GPT-4V-all Experiment

We also collected samples on recovery actions by GPT-4V. In this experiment, GPT-4V is responsible for all four steps. We prepared four GPT-4V models that perform each of the four steps of the recovery flow described in Section 3.2. We provided all GPT-4V models with the SoM-image, the number of the object predicted in the failed game, and the dialogue history as inputs (See Appendix A.2 for details of these models.). By comparing Human-all

	Step1	Step2	Step3	Step4
GPT-4V-all	GPT4	GPT4	GPT4	GPT4
GPT-4V-Q	human	GPT4	human	human
Human-all	human	human	human	human

Table 1: Questioner’s roles of humans and GPT-4V in each step in each experiment. Answer’s roles were performed by humans at all experiments

experiment to GPT-4V-all experiment, we can validate GPT-4V’s ability to recover after task failure relative to human recovery ability.

### 4.2.3 GPT-4V-Q Experiment

Then we also collected samples on recovery actions by GPT-4V and humans. In this experiment, GPT-4V is in charge of only step 2 among the four steps, and humans are in charge of the other steps. By comparing Human-all experiment to GPT-4V-Q experiment, it is possible to more directly assess the GPT-4V’s ability to repair utterance (corresponding to step 2). In both GPT-4V-all experiment and GPT-4V-Q experiment, we collected three recovery samples per failed game, resulting in a total of 300 recovery samples, which is the same number as Human-all experiment. In all three experiments, if a human or GPT-4V repeated additional questions (Step 2) more than 11 times, we counted the attempt as a failure.

## 4.3 Evaluation

High recovery success rates and efficiency are essential for the failure recovery task, as mentioned in Section 3.2. Accordingly, we evaluate and analyze from the following two perspectives: (1) the success rate of recovery and (2) the average number of repair utterances required for a successful recovery. The success rate of recovery is defined as the number of recovered games over the total number of failed games. We define a successful recovery action as a successful re-prediction of the target object. We calculate the average number of repair utterances required for successful correction by averaging the additional questions made during the second step of the recovery process. This average includes instances where no repair utterances were necessary, counted as zero. A lower average indicates a more efficient recovery process, preferable for optimal outcomes.

**Failure Game**



Q1: それは、スポーツ用具ですか？ (Is it sports equipment?)	A1: いいえ。 (No.)
Q2: それは、人間ですか？ (Is it a human being?)	A2: はい。 (Yes.)
Q3: 画像の中央にいますか？ (Is it in the center of the image?)	A3: いいえ。 (No.)
Q4: 赤い服を着ていますか？ (Is it wearing red?)	A4: はい。 (Yes.)
Q5: ボールを蹴っていますか？ (Is it kicking a ball?)	A5: いいえ。 (No.)

Target object: No. 4  
1st prediction object: No. 7

---

**GPT-4V Repair (GPT-4V-all experiment)**

Q6: ヒジブを着用していますか？ (Is it wear a hijab?)	A6: いいえ。
Q7: 緑色のベストを着ていますか？ (Is it wear a green vest?)	A7: はい。

Re-prediction object: No.4 ✓

---

**Human Repair (Human-all experiment)**

Q6: 被り物をしていますか？ (Is it wearing a head covering?)	A6: はい。
Q7: その被り物は赤色ですか？ (Is the covering red?)	A7: はい。

Re-prediction object: No. 4 ✓

(a) Example of the successful recovery of both the human and GPT-4V

**Failure Game**



Q1: それは、電子機器ですか？ (Is it an electronic device?)	A1: はい。 (Yes.)
Q2: それは、携帯電話ですか？ (Is it a cell phone?)	A2: いいえ。 (No.)
Q3: それは、ノートパソコンですか？ (Is it a laptop?)	A3: いいえ。 (No.)
Q4: それは、テレビですか？ (Is it a TV?)	A4: はい。 (Yes.)
Q5: 画像の中央にありますか？ (Is it in the center of the picture?)	A5: いいえ。 (No.)

Target object: No. 6  
1st prediction object: No. 4

---

**GPT-4V Repair (GPT-4V-all experiment)**

Q6: 画像の上部にありますか？ (Is it at the top of the image?)	A6: いいえ。
Q7: 画像の右側にありますか？ (Is it on the right side of the image?)	A7: はい。

Re-prediction object: No. 7 ✗

---

**Human Repair (Human-all experiment)**

No additional questions have been asked.

Re-prediction object: No. 6 ✓

(b) Example of the successful recovery of the human but the failure of GPT-4V

Figure 4: The recovery results for the same failed game. The image is labeled with numbers by SoM (Yang et al., 2023)). The upper proper dialogue represents the dialogue history during the task failure.

## 5 Results

### 5.1 Human vs GPT-4V for Recovery Action

We compared GPT-4V-all experiment and Human-all experiment to evaluate the abilities of humans and GPT-4V in implementing recovery actions after failing Guess What?! Game.

The success rate of recovery actions by GPT-4V is significantly lower than in humans (about 36.7% lower), which means that the failure task recovery in Guess What? Game (de Vries et al., 2017) is even difficult for GPT-4V. Figure 4a shows an example where both GPT-4V and the human was successful. GPT-4V successfully re-guessed the target object by asking two additional questions (Q6 and Q7), much like the human did, although using a different method of questioning. Figure 4b presents a case where GPT-4V failed, but the human succeeded. The human identified the target object without asking additional questions, whereas GPT-4V asked two questions (Q6 and Q7) and still failed to predict correctly. Despite confirming that the object was not a mobile phone in Q2, it incorrectly guessed the target object as No.7. This example shows failures in Step 1, 3 modules, which are responsible for deciding whether to ask an additional question, and in the Step 4 module, which is responsible for predicting the final object.

Next, we compared the efficiency of failure recovery tasks between humans and GPT-4V by analyzing the average number of repair utterances. As

noted in Section 4.3, the recovery process must be efficient in human interactions. We calculated the average number of repair utterances only for successful cases because efficient recoveries are only relevant when the recovery task is successful (attempting a quick fix is pointless if it fails). Table 2 shows the average turn of repair utterances. The GPT-4V-all experiment is more than twice utterances as many as the Human-all experiment, indicating a less efficient recovery strategy in GPT4-V compared to humans.

### 5.2 First Half Recovery Steps Analysis

In this section, we focus on GPT-4V’s ability to decide and successfully execute repair utterances (steps 1 and 2) rather than just understanding and using them (steps 3 and 4).

#### 5.2.1 Step1: Deciding Recovery Action

We compared the actions chosen by GPT-4V with those selected by humans in step 1. Specifically, for each game, we tallied and compared the number of times actions (a-I) asking an additional question and (a-II) re-guessing the target object were chosen. We selected actions chosen at least twice by the human and GPT-4V across three recovery tasks for the same failure sample as the actions by the recovery executor.

Table 3 compares the actions selected by humans (Human-all experiment) and GPT-4V (GPT-4V-all experiment). As Section 5.1 indicates, human re-

	Success rate	Average turn
GPT-4V-all	50.0%	2.43
GPT-4V-Q	74.7%	2.00
Human-all	86.7%	1.13

Table 2: Success rate of recovery actions and the number of repair utterances (step 2) in each experiment.

	GPT4-V (a-I)	GPT4-V (a-II)
Human (a-I)	65	5
Human (a-II)	28	2

Table 3: The number of actions selected by Humans and GPT-4V (a-I or a-II). Diagonal elements show the number of times Humans and GPT-4V made the same selections. Note that these counts are from Step 1 of GPT-4V-all experiment and Human-all experiment.

covery actions have a high success rate and are a strong baseline. Thus, GPT-4V should choose actions that are similar to those chosen by humans in most cases. Table 3 shows that GPT-4V selects about 67% of the same actions as humans, and GPT-4V often opts to select (a-I) action even in cases where humans choose (a-II) action. This result demonstrates that GPT-4V fails to provide efficient questions and choose speedy recovery actions. This behavior is undesirable because making efficient failure task recovery is crucial in goal-oriented dialogues with humans.

### 5.2.2 Step2: Asking an Additional Question (Repair Utterance)

Next, we analyzed the repair utterances from Step 2. Specifically, we compared GPT-4V-Q experiment, in which humans handled all steps except Step 2, with Human-all experiment, in which humans were responsible for all steps.

Table 2 shows the results for task success rates and the number of repair utterances. When comparing GPT-4V-Q experiment to Human-all experiment, we observe that only replacing Step 2 with GPT-4V results in a 12% decrease in success rate and an increase of 0.87 in average turns. Furthermore, in GPT-4V-all experiment where GPT-4V handles all steps, the success rate drops by an additional 24.7%, and the average turns increase by 0.43. This suggests that GPT-4V’s impact in Step 2 contributes more to the increase in the number of turns than modules of other steps, which means that GPT-4V’s repair utterances tend to include unnecessary questions.

Next, we analyzed the intents behind the utterances to compare the nature of repair utterances made by humans and GPT-4V. We asked humans

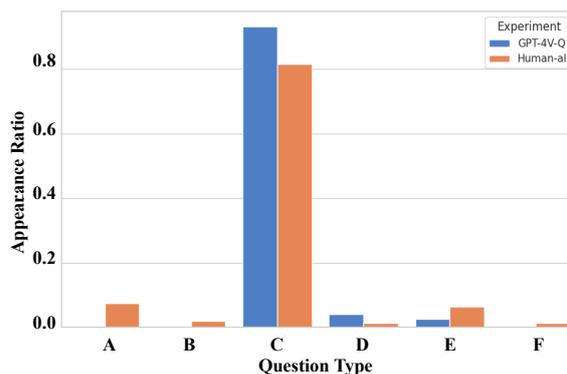


Figure 5: Question type distribution. We normalize by dividing the number of each question type by the total number of questions in each experiment because the number of questions differs between the two experiments.

and GPT-4V to select the intents behind their questions from six options and compared these selections. We assumed that GPT-4V could understand the intentions behind its questions, so we had GPT-4V select the intent of the questions. We also asked for explanations behind the selected options, and the lead author checked to see if the reason was plausible because GPT-4V does not always produce accurate outputs. We conducted preliminary experiments and prepared the following six types of questions (A)-(F) (see Appendix B for detailed explanations of question types.):

- (A) The question that addresses the same object with different expressions: This questioning style is employed when there is a suspicion of inconsistencies or errors in the user’s answers.
- (B) The question with more or less the same meaning as the question during dialogue: This questioning style is used when there are suspected inconsistencies or errors in the user’s answers.
- (C) The question that proposes a hypothesis to narrow down the object in question: This type of question is used when there are no apparent errors or contradictions in the user’s answers.
- (D) The question that clarifies ambiguities in a previous question: It clarifies the context or perspective of the previous question.
- (E) The question for confirmation, in case the object has already been narrowed down.
- (F) Others. (In this case, we ask the annotator and GPT-4V to describe the question’s intent in text form.)

Figure 5 shows the distribution of the intentions behind the questions asked by humans and GPT-4V. First, more type (C) questions exist in both Human-

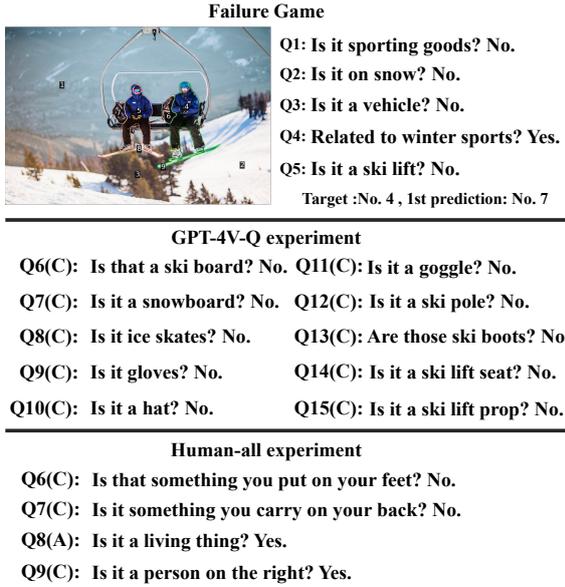


Figure 6: Example of a human asking a type (A) question, whereas GPT-4V does not. Q6 (C) indicates that the first additional question is the intent of type (C). The dialogue conducted in Japanese is translated into English.

all and GPT-4V-all experiments. This is because many samples in the failed game set, such as the example in Figure 4, require additional questions to narrow down the objects. A significant difference between the human-only experiment and the GPT-4V with human experiment is that the humans can ask many non-type (C) questions. Specifically, humans ask type (A) or (B) questions about 9.3% of the time, whereas GPT-4V rarely asks these types of questions. This indicates that GPT-4V cannot recognize or doubt mistakes and inconsistencies based on the user’s input, leading it to focus predominantly on questions that narrow down objects.

Figure 6 shows an example where a human asks a type (A) question while GPT-4V does not. In this example, the human suspects an error in the answer and attempts to correct the course of the dialogue by asking, "Is it a living thing?" (Question type A). In contrast, GPT-4V likely overtrusts the response, "Related to winter sports? Yes," and continues to ask questions focused on objects related to winter sports (Question type C). As a result, GPT-4V fails to correct the course of the dialogue and cannot identify the target object within ten questions, leading to an unsuccessful recovery task. These results indicate that GPT-4V fails to recognize or question erroneous responses and tends to blindly trust the user’s input.

## 6 Discussion

We observed that GPT-4V is significantly poor at the failure recovery task (§5.1), and GPT-4V’s approach differs from the strong baseline of human behavior in both Step 1 and Step 2 (§5.2). The significant difference in failure task recovery capabilities between humans and GPT-4V can be attributed to the models’ difficulty with logical reasoning (Creswell et al., 2023; Pan et al., 2023; You et al., 2023). GPT-4V may struggle to integrate three pieces of information from failed games (game image, dialogue text, and first prediction object) to identify potential target objects. Unlike typical goal-oriented dialogues, conducting failure task recovery requires understanding complex dialogue and game situations. Therefore, the pre-processing step that explicitly organizes the context of the failed game and dialogue rather than executing direct recovery actions may be practical. In Guess What?! Game, output which objects remain as potential targets is an example of this strategy.

We also found that GPT-4V tends to refrain from questioning the interlocutor’s answer (§5.2.2). This feature is undesirable for the failure task recovery in goal-oriented dialogues, where the user’s answers might contain errors (Oshima et al., 2023). This issue is not crucial during initial task attempts because humans do not frequently make response errors. However, when the task fails, the possibility of user answer errors increases. Therefore, considering the possibility of user errors is a key factor when developing recovery strategies. If a system cannot doubt the user’s answers, it may fail to correct errors or waste time, as shown in Figure 6. To address this issue, instructing the LLM first to evaluate the correctness of the user’s answers and then use this evaluation to guide the recovery action may be effective.

One future direction to use VLMs (Liu et al., 2023b,a) as a recovery model is rethinking model training. For example, creating synthetic datasets that include incorrect utterances and using them for instruction-tuning data. This approach allows VLMs to explicitly learn from erroneous scenarios, potentially enhancing their abilities to recover task failure accurately.

## 7 Limitations

In this study, we examined the failure recovery task in Guess What?! Game, where one speaker only responds with “Yes.” or “No.” However, this research

does not address the recovery capabilities of GPT-4V and humans in more complex goal-orientated dialogues like autonomous driving dialogue systems. Our results may differ for languages other than Japanese, so it is essential to analyze GPT-4V’s recovery performance in English, its most proficient language. This study focuses on GPT-4V, raising concerns about the generalizability of our findings to other vision-language models. We tested the failure recovery task with LLaVA-1.5 (Liu et al., 2023a), but it did not recover the tasks adequately, which suggests that the recovery task would require capabilities comparable to GPT-4V.

We are concerned about the method of collecting intents by directly asking humans or using GPT-4V. This method assumes a causal relationship between subjective reasoning and actual behavior. As Ayaß (2015) recommended, it is preferable to analyze speech intentions based on objective actual behavior rather than subjective reasoning.

## 8 Conclusion

We tackled the failure recovery task in Guess What?! Game and analyzed GPT-4V’s capabilities. The results showed that GPT-4V demonstrated a significantly lower ability to correct task failures than humans. Furthermore, GPT-4V tended to perform unnecessary repair utterances, ask inefficient questions, and fail to doubt users’ answers. In future work, we aim to investigate the generalizability of our findings to real-world goal-oriented dialogues.

## Acknowledgments

This work was supported by JSPS KAKENHI Grant Numbers JP21H05054, JP21K17806, JP24H00742, and JP24H00742.

## References

Ruth Ayaß. 2015. [Doing data: The status of transcripts in conversation analysis](#). *Discourse Studies*, 17(5):505–528.

Vevake Balaraman, Arash Eshghi, Ioannis Konstas, and Ioannis Papaioannou. 2023. [No that’s not what I meant: Handling third position repair in conversational question answering](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 562–571, Prague, Czechia. Association for Computational Linguistics.

Luciana Benotti and Patrick Blackburn. 2021a. [Grounding as a collaborative process](#). In *Proceedings of the*

*16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 515–531, Online. Association for Computational Linguistics.

Luciana Benotti and Patrick Blackburn. 2021b. [A recipe for annotating grounded clarifications](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4065–4077, Online. Association for Computational Linguistics.

Javier Chiyah-Garcia, Alessandro Suglia, Arash Eshghi, and Helen Hastie. 2023. [‘what are you referring to?’ evaluating the ability of multi-modal dialogue models to process clarificational exchanges](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 175–182, Prague, Czechia. Association for Computational Linguistics.

Herbert Clark. 1996. *Using Language*. Cambridge University Press.

Antonia Creswell, Murray Shanahan, and Irina Higgins. 2023. [Selection-inference: Exploiting large language models for interpretable logical reasoning](#). In *The Eleventh International Conference on Learning Representations*.

Harm de Vries, Florian Strub, Sarath Chandar, Olivier Pietquin, Hugo Larochelle, and Aaron Courville. 2017. [Guesswhat?! visual object discovery through multi-modal dialogue](#). In *Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1080–1089. Institute of Electrical and Electronics Engineers.

Yang Deng, Shuaiyi Li, and Wai Lam. 2023. [Learning to ask clarification questions with spatial reasoning](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’23*, page 2113–2117, New York, NY, USA. Association for Computing Machinery.

Zhiyu Fan, Xiang Gao, Martin Mirchev, Abhik Roychoudhury, and Shin Hwei Tan. 2023. [Automated repair of programs from large language models](#). In *Proceedings of the 45th International Conference on Software Engineering, ICSE ’23*, page 1469–1481. IEEE Press.

Janosch Haber, Tim Baumgärtner, Ece Takmaz, Lieke Gelderloos, Elia Bruni, and Raquel Fernández. 2019. [The PhotoBook dataset: Building common ground through visually-grounded dialogue](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1895–1910, Florence, Italy. Association for Computational Linguistics.

Kai Huang, Xiangxin Meng, Jian Zhang, Yang Liu, Wenjie Wang, Shuhao Li, and Yuqing Zhang. 2023. [An empirical study on fine-tuning large language](#)

- models of code for automated program repair. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 1162–1174.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. 2022. Inner monologue: Embodied reasoning through planning with language models. In *arXiv preprint arXiv:2207.05608*.
- Vojtěch Hudeček and Ondrej Dusek. 2023. Are large language models all you need for task-oriented dialogue? In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 216–228, Prague, Czechia. Association for Computational Linguistics.
- Satwik Kottur, Seungwhan Moon, Alborz Geramifard, and Babak Damavandi. 2021. SIMMC 2.0: A task-oriented dialog dataset for immersive multimodal conversations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4903–4912, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Ziqiao Ma, Benjamin VanDerPloeg, Cristian-Paul Bara, Yidong Huang, Eui-In Kim, Felix Gervits, Matthew Marge, and Joyce Chai. 2022. DOROTHIE: Spoken dialogue for handling unexpected situations in interactive autonomous driving agents. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4800–4822, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Brielen Madureira and David Schlangen. 2023. Instruction clarification requests in multimodal collaborative dialogue games: Tasks, and an analysis of the Co-Draw dataset. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2303–2319, Dubrovnik, Croatia. Association for Computational Linguistics.
- Mauricio Mazuecos, Patrick Blackburn, and Luciana Benotti. 2021. The impact of answers in referential visual dialog. In *Proceedings of the Reasoning and Interaction Conference (ReInAct 2021)*, pages 8–13, Gothenburg, Sweden. Association for Computational Linguistics.
- OpenAI. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Ryosuke Oshima, Seitaro Shinagawa, Hideki Tsunashima, Qi Feng, and Shigeo Morishima. 2023. Pointing out human answer mistakes in a goal-oriented visual dialogue. In *VISION-AND-LANGUAGE ALGORITHMIC REASONING Workshop in International Conference on Computer Vision 2023*.
- Tim Paek and Eric Horvitz. 2000. Conversation as action under uncertainty. In *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, UAI'00, page 455–464, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-LM: Empowering large language models with symbolic solvers for faithful logical reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3806–3824, Singapore. Association for Computational Linguistics.
- Matthew Purver, Julian Hough, and Christine Howes. 2018. Computational models of miscommunication phenomena. *Topics in Cognitive Science*, 10(2):425–451.
- Ruben San-Segundo, Juan Manuel Montero, and Jose Manuel Pardo. 2001. Designing confirmation mechanisms and error recover techniques in a railway information system for Spanish. In *Proceedings of the Second SIGdial Workshop on Discourse and Dialogue*.
- E.A. Schegloff. 1992. Repair after next turn: The last structurally provided defense of intersubjectivity in conversation. *American Journal of Sociology*, 97(5):1295–1345.
- E.A. Schegloff, Gail Jefferson, and Harvey Sacks. 1977. The preference for self-correction in the organization of repair in conversation. *Language*, 53(2):361–382.
- Omar Shaikh, Kristina Gligorić, Ashna Khetan, Matthias Gerstgrasser, Diyi Yang, and Dan Jurafsky. 2024. Grounding gaps in language model generations. *Preprint*, arXiv:2311.09144.
- Zhengxiang Shi, Yue Feng, and Aldo Lipani. 2022. Learning to execute actions or ask clarification questions. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2060–2070, Seattle, United States. Association for Computational Linguistics.
- Veniamin Veselovsky, Manoel Horta Ribeiro, and Robert West. 2023. Artificial artificial artificial intelligence: Crowd workers widely use large language models for text production tasks. *Preprint*, arXiv:2306.07899.
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*.

Penghao Wu and Saining Xie. 2023. V\*: Guided visual search as a core mechanism in multimodal llms. *arXiv preprint arXiv:2312.14135*.

Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. 2023. [Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v](#). *Preprint*, arXiv:2310.11441.

Haoxuan You, Rui Sun, Zhecan Wang, Long Chen, Gengyu Wang, Hammad Ayyubi, Kai-Wei Chang, and Shih-Fu Chang. 2023. [IdealGPT: Iteratively decomposing vision and language reasoning via large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11289–11303, Singapore. Association for Computational Linguistics.

Quanjun Zhang, Chunrong Fang, Yuxiang Ma, Weisong Sun, and Zhenyu Chen. 2023. [A survey of learning-based automated program repair](#). *ACM Trans. Softw. Eng. Methodol.*, 33(2).

## A GPT-4V Model Details

### A.1 Models for Failure Game Collection

Figure 7 shows an overview of the model used for failure game collection in Section 4.1. We provide the original SoM-image to the questioner and guesser models, while the answerer model receives an SoM-image with the target object highlighted in a yellow frame. Table 4 shows the text prompts provided to the questioner and guesser models. Table 5 indicates the prompts given to the answerer model.

### A.2 Models for Failure Task Recovery Collection

The basic framework of the models is the same as the GPT-4V model prepared in failure game collection. Each model uses the SoM-image but with different text prompts. The text prompts consist of two main parts: the system prompt, which is the rule of Guess What?! Game, and the user prompt, which is the specific instructions performed by GPT-4V. While the system prompt is consistent across all models, the user prompt varies by step. Table 6 illustrates the system prompt, and Table 7 and 8 provide examples of the user prompt.

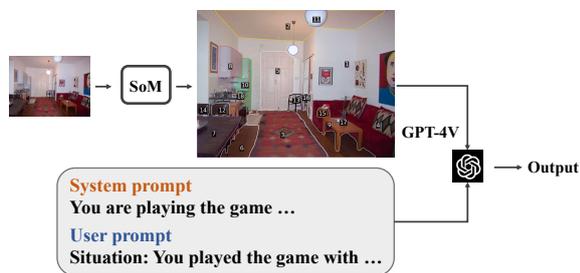


Figure 7: Overview of the model used for the failure game collection (§4.1) and failure task recovery collection (§4.2).

## B Details of Question’s intentions

Table 9 shows a detailed description of each question type. We gave the same content written in Table 9 to GPT-4V while we showed the users a Japanese translation of the content described in the table. Table 10 presents cases where humans selected type (F) for the question intent, along with descriptions of their intent.

Prompt type	Example
Questioner prompt	<p>Instruction: You are now playing the following game with a user.  Rule of the game: You are trying to guess the target object in the image by asking questions. User will answer yes or no to the question about the user's target object. Note that target objects can also include living creatures such as humans and animals.  Generate a question "in Japanese" to guess the target object.  例：「それは、人間ですか?」、「それは、画像の左側にありますか?」、「それは、野菜ですか?」など、なんでも構いません。  (Example: You can ask anything, such as "Is that a human?" or "Is that a vegetable?" or "Is it on the left side of the picture?")</p> <p>Dialogue history:  質問: 「それは、テニスラケットですか?」  回答: はい。  では、質問を生成してください。  (Now, please make a question.)  質問:  (Question:)</p>
Guesser prompt	<p>Instruction: You are now playing the following game with a user.  Rule of the game: You are trying to guess the target object in the image by asking you questions. User will answer yes or no to the question about the user's target object.  You are trying to guess the target object through dialogue.  Dialogue history:  質問: それは、テニスラケットですか?  回答: はい。  質問: それは、人間ですか?  回答: いいえ。  質問: それは、画像の右側にありますか?  回答: はい。  質問: それは、男性ですか?  回答: いいえ。  質問: それは、テニスコートですか?  回答: いいえ。  Please predict the target object number. Be sure to focus your prediction on one number!  Number:</p>

Table 4: Questioner's role prompt examples for failure game collection (§4.1).

Prompt type	Example
Answerer prompt	<p>Instruction: You are now playing the following game with a user.  Rule of the game: The user is trying to guess the target object in the image by asking you questions. Answer yes or no to the question about the user's target object.  The target object is labeled with number 8, surrounded by yellow box and its category is TENNIS RACKET.</p> <p>質問には、「はい。」または「いいえ。」で教えてください。  (Please answer with “Yes” or “No”.)</p> <p>質問：それは、人間ですか？  (Question: Is it a human?)</p> <p>回答：  (Answerer: )</p>

Table 5: Answerer’s role prompt examples for failure game collection (§4.1).

Prompt type	Example
System prompt	<p>You are now playing the following game with a user. You are a professional in this game.  Rule of the game: You are trying to guess the target object in the image by asking questions. User will answer yes or no to the question about the user's target object. Note that numbered objects are candidates for target objects and target objects can also include living creatures such as humans and animals.</p>

Table 6: System prompt for the failure task recovery collection (§4.2)

Prompt type	Example
User prompt of step1	<p>Situation: You played the game with the user and also predicted the object. However, the object you predicted (number 3) was not the right target object, and the task failed. Therefore, you need to take a repair action to turn this game into a success (guessing the correct target object) instead of a failure. This repair action can be either asking additional questions or re-predicting the object. Which is better?</p> <p>Dialogue when you fail to predict the target object (Dialogue before the first prediction):</p> <p>質問: それはスポーツ用品ですか?  回答: いいえ。  質問: それは生き物ですか?  回答: はい。  質問: それは人間ですか?  回答: はい。  質問: それは、画像の中央にいますか?  回答: いいえ。  質問: それは、画像の右側にいますか?  回答: いいえ。</p> <p>Instruction: Please answer the number of the action you take. Note that target objects must be numbered.</p> <p>(1). You do not ask additional questions and re-guess the target object  (2). You ask additional questions in order to re-guess the target object</p>
User prompt of step2	<p>Situation: You played the game with the user and also predicted the object. However, the object you predicted (number 3) was not the right target object, and the task failed.</p> <p>Dialogue when you fail to predict the target object (Dialogue before the first prediction):</p> <p>質問: それはスポーツ用品ですか?  回答: いいえ。  質問: それは生き物ですか?  回答: はい。  質問: それは人間ですか?  回答: はい。  質問: それは、画像の中央にいますか?  回答: いいえ。  質問: それは、画像の右側にいますか?  回答: いいえ。</p> <p>Instruction: You have determined that you need to ask additional questions. Please make one question to re guess the target object based on failed dialogue.</p> <p>Notes:</p> <p>1. Numbered objects are candidates for a target object.  2. Please do not ask additional questions using the number assigned to the object or ask questions that mention that number or letter!</p> <p>では、質問を生成してください。  (Now, please make a question.)  質問:  (Question:)</p>

Table 7: Prompts of Steps 1 and 2 for the failure task recovery collection (§4.2)

Prompt type	Example
User prompt of step3	<p>Situation: You played the game with the user and also predicted the object. However, the object you predicted (number 1) was not the right target object, and the task failed. Then, you asked additional questions to the user in order to re-guess the target object.</p> <p>Dialogue when you fail to predict the target object (Dialogue before the first prediction):</p> <p>質問: それは、鏡の中に映っていますか?  回答: いいえ。  質問: それは、洗面台ですか?  回答: はい。  質問: それは、人間ですか?  回答: いいえ。  質問: それは、タイルでできていますか?  回答: いいえ。  質問: それは、水道の蛇口ですか?  回答: いいえ。</p> <p>Additional questions and answers to re-predict the correct target object (Dialogue after the first prediction):</p> <p>質問: それは、壁に取り付けられていますか?  回答: いいえ。</p> <p>Instruction: Please review the previous conversation and decide if the target object in this game has been clearly identified. Note that the object should be identifiable by a number. Respond with:</p> <p>(1) The object has been clearly identified.  (2) The object has not been identified, and further questions are necessary.</p>
User prompt of step4	<p>Situation: You played the game with the user and also predicted the object. However, the object you predicted (number 3) was not the right target object, and the task failed.</p> <p>Dialogue when you fail to predict the target object (Dialogue before the first prediction):</p> <p>質問: それはスポーツ用品ですか?  回答: いいえ。  質問: それは生き物ですか?  回答: はい。  質問: それは人間ですか?  回答: はい。  質問: それは、画像の中央にいますか?  回答: いいえ。  質問: それは、画像の右側にいますか?  回答: いいえ。</p> <p>Additional questions and answers to re-predict the correct target object (Dialogue after the first prediction):</p> <p>質問: それは画像の左側にいますか?  回答: はい。</p> <p>Instruction: Please read the dialogue history above and re-predict the target object number. Be sure to focus your prediction on one number!  The target object: Number</p>

Table 8: Prompts of Steps 3 and 4 for the failure task recovery collection (§4.2)

Question intention	Description
(A) The question that addresses the same object with different expressions.	This questioning style is employed when there is a suspicion of inconsistencies or errors in the user's answers. It involves exploring the same object in an image through various expressions. This method helps identify any inconsistencies or errors in the user's answers by exploring different aspects of the same object and examining the object from multiple angles.
(B) The question with more or less the same meaning as the question during dialogue.	This questioning style is used when there are suspected inconsistencies or errors in the user's answers. It involves slight rephrasing of previous questions using similar terms to clarify and rectify any misunderstandings. Example: Rephrase the question "Are people using it?" in the dialogue history as "Is it something that people are holding?" or "Is it a human figure?" as "Is it really a human figure?"
(C) The question that proposes a hypothesis to narrow down the object in question.	This type of question is used when there are no apparent errors or contradictions in the user's answers. It introduces a hypothesis to further refine and specify the inquiry, aiming to deepen the exploration of the object in question. This approach helps gather more precise information about the object being discussed.
(D) The question that clarifies ambiguities in a previous question	Employed in cases where the meaning of the previous question has caused ambiguity, this questioning style seeks to align the understanding between the questioner and the respondent. It clarifies the context or perspective of the previous question. It clarifies the context or perspective of the previous question, as in the example: "Is it on the left?" might be followed by a clarifying question, "Is it on the left side of the image?" to specify the context.
(E) The question for confirmation, in case the object has already been narrowed down	This questioning style is utilized when the answers so far have no detected errors or contradictions. It aims to reaffirm the accuracy and certainty of the information provided by the respondent, ensuring a solid and shared understanding of the discussed object or situation.
(F) Others.	

Table 9: Question intentions and detail description.

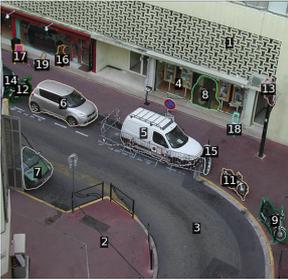
Recovery cases	Description of question intentions
 <p>Q1: Is it a vehicle? A1: Yes.  Q2: Is it a four-wheeled vehicle? A2: Yes.  Q3: Is it white? A3: Yes.  Q4: Does it have a load on the roof? A4: No.  Q5: Is it in the center of the image? A5: No.  -----  Q6: Is that a motorcycle? A6: No.  Q7: Is it green? A7: No.  Q8: Is it fenced? A8: Yes.</p> <p>Target object: No. 5, 1st prediction object: No. 6, Repair prediction object: No. 5 ✓</p>	<p>Q6: (F) In the dialogue, I noticed errors or contradictions in the answers and, while trusting the question “Is it a vehicle?” with the response “Yes,” restarted the discussion about other topics to explore the possibility of other objects.  Q7: (C)  Q8: (A)</p>
 <p>Q1: Is it someone wearing headphones? A1: No.  Q2: Is it a person wearing glasses? A2: Yes.  Q3: Is the person lying down? A3: Yes.  Q4: Is the person on the bed? A4: Yes.  Q5: Is the person wearing headphones? A5: No.  -----  Q6: Is it on the left side of the screen? A6: Yes.</p> <p>Target object: No. 6, 1st prediction object: No. 5, Repair prediction object: No. 6 ✓</p>	<p>Q6: (F) Since it was uncertain whether the baby was wearing glasses, I asked additional questions to make sure.</p>
 <p>Q1: Is that the character you see on the screen? A1: No.  Q2: Is it an electronic device? A2: Yes.  Q3: Is it on a desk? A3: Yes.  Q4: Is it the monitor in the center of A4: Yes.  Q5: Is it on the right side of the image? A5: No.  -----  Q6: Is it the monitor on the left side of the image? A6: Yes.</p> <p>Target object: No. 4, 1st prediction object: No. 6, Repair prediction object: No. 4 ✓</p>	<p>Q6: (F) Questions to notice mistakes or inconsistencies in the answers during the dialogue and to narrow the object based on them.</p>
 <p>Q1: Is it in the kitchen? A1: Yes.  Q2: Does it use water? A2: No.  Q3: Is it used for cooking? A3: Yes.  Q4: Does it use heat? A4: Yes.  Q5: Is it an oven? A5: No.  -----  Q6: It on the right side of the screen? A6: Yes.</p> <p>Target object: No. 7, 1st prediction object: No. 6, Repair prediction object: No. 7 ✓</p>	<p>Q6: (F) Questions that narrow down objects based on noticing errors or contradictions in responses during a conversation.</p>

Table 10: The samples and the description where a human selected the question type (F).

# MindDial: Enhancing Conversational Agents with Theory-of-Mind for Common Ground Alignment and Negotiation

Shuwen Qiu<sup>1</sup> Mingdian Liu<sup>2</sup> Hengli Li<sup>3,4,6</sup> Song-Chun Zhu<sup>3,4,5,6</sup> Zilong Zheng<sup>4,6</sup>✉

<sup>1</sup> UCLA <sup>2</sup> ISU <sup>3</sup> PKU <sup>4</sup> BIGAI <sup>5</sup> THU

<sup>6</sup> State Key Laboratory of General Artificial Intelligence

s.qiu@ucla.edu, mingdian@iastate.edu, {lihengli, sczhu, zlzheng}@bigai.ai

## Abstract

Humans talk in daily conversations while aligning and negotiating the expressed meanings or common ground. Despite the impressive conversational abilities of the large generative language models, they do not consider the individual differences in contextual understanding in a shared situated environment. In this work, we propose MindDial, a novel conversational framework that can generate situated free-form responses with theory-of-mind (ToM) modeling. We introduce an explicit mind module that can track the speaker’s belief and the speaker’s prediction of the listener’s belief. Then the next response is generated to resolve the belief difference and take task-related action. Our framework is applied to both prompting and fine-tuning-based models, and is evaluated across scenarios involving both common ground alignment and negotiation. Experiments show that models with mind modeling can achieve higher task outcomes when aligning and negotiating common ground. The ablation study further validates the three-level belief design can aggregate information and improve task outcomes in both cooperative and negotiating settings.

## 1 Introduction

We align and negotiate our common ground every day in daily chit-chat (Clark and Wilkes-Gibbs, 1986; Bazerman et al., 2000). In a common ground alignment scenario, agents are talking toward a joint goal, topics ranging from daily trivia to important multi-party meetings. In common ground negotiation situations, two parties resolve the differences in their beliefs, intents, or goals in a way that both find acceptable, such as item trading and discussing job offers (Veinott et al., 1999; Beers et al., 2006). Though it seems easy between human conversations, it requires complicated social capabilities. Importantly, for all types of human communication including language, the relationship between

the overt communicative act and common ground – of whatever type – is complementary. That is, as more can be assumed to be shared between communicator and recipient, less needs to be overtly expressed (Tomasello, 2010). Taking Figure 1B as an example, when Bob asks about “Joe Davis”, Alice will align the precise referents of the query by keeping “Joe” but correct “Davis” to “Smith”. In this process, people need to realize what is shared and what needs to be further aligned or negotiated – which requires the understanding between points of view from their own and others’ perspectives (Blutner, 2000; De Weerd et al., 2015) – the cognitive capability known as theory-of-mind (ToM).

The recent surge of large language models (LLMs) (Radford et al., 2019; Brown et al., 2020) have dominated the natural language processing (NLP) community for their prominent natural language generation performance. Although LLMs have shown their potential in ToM benchmarks (Kosinski, 2023; Ullman, 2023; Sileo and Lerneuld, 2023; Kim et al., 2023), applying ToM for situated dialogue generation remains underexplored. In these situated tasks, agents’ interactions are influenced by the environment, their shared experiences, and immediate goals. The participants need to take into account not only the linguistic content but also factors such as the social context, prior knowledge, and each other’s beliefs. Without ToM, the models can only provide the most possible response as a one-turn question-answering as shown in Figure 1A. To enable LLMs to interact with people in a more socially realistic manner, it is essential to incorporate ToM for various forms of communications, such as aligning and negotiating common ground within dialogues (Burlison, 2007; Chiu et al., 2023; Fu et al., 2023).

In this work, we introduce **MindDial**, a new dialogue framework designed to facilitate the alignment and negotiation of common ground in situated dialogues, incorporating ToM modeling. In-

✉ Correspondence to Zilong Zheng <zlzheng@bigai.ai>.

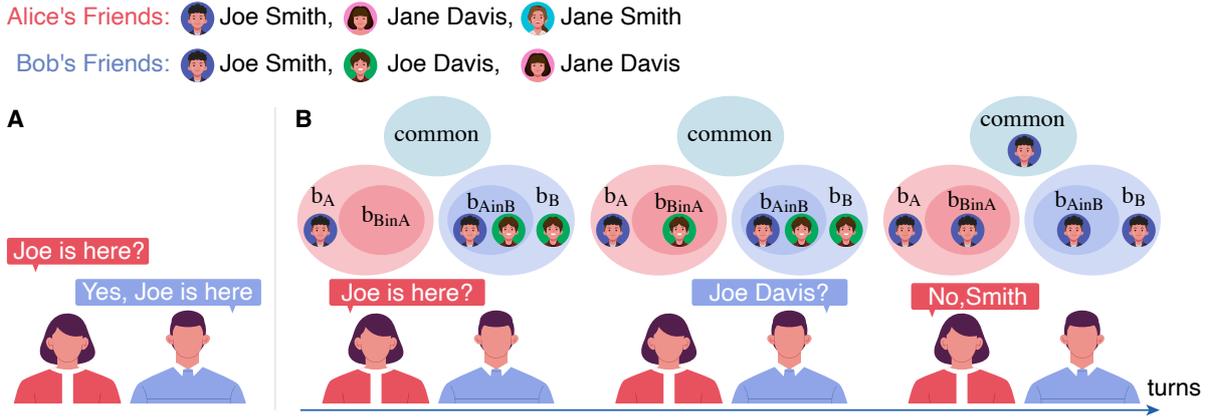


Figure 1: **Left:** Single-turn question answering. **Right:** Multi-turn common ground alignment. Speakers will update their belief estimation based on context and generate the next response to reduce the belief differences.

spired by the complementary role between common ground and communication, we design the two-step response generation. First, an explicit mind module estimates the speaker’s current perspective of the world (the first-order belief) and also helps speaker’s estimate the other’s perspective of the world (the second-order belief) (Grueneisen et al., 2015; Braüner et al., 2016). Then, the next response is aimed at resolving the belief difference. As shown in Figure 1B, Alice says “No, Smith” when her first-order belief  $b_A$  (“Joe Smith”) does not equal to her second-order belief  $b_{BinA}$  (“Joe Davis”).

In sum, we consider our contributions as three-fold:

i) We design a framework incorporating an explicit mind estimation module that tracks the first-order and second-order beliefs. Resolving the belief difference between the two will support the next response generation.

ii) We explore two types of response generators – fine-tuning and prompting-based models. The experiments show that our framework can successfully improve model performance in both groups.

iii) We test our framework on both aligning and negotiating settings. The evaluation results and user study validate that our framework can improve both the cooperation and negotiation abilities of the LLM agents. We ablate each level of the beliefs and find both first and second-order contribute to the final results.

## 2 Related work

**Theory-of-Mind (ToM)** ToM is a crucial capability for human social interactions developed in early life (Kovács et al., 2010; Richardson

et al., 2018). In literature, early works model belief update through time in sequential games with partially observable Markov decision process (POMDP) (Baker et al., 2011; De Weerd et al., 2013; Vogel et al., 2013; Doshi et al., 2010; Han and Gmytrasiewicz, 2018). One agent’s belief update is based on the estimate of others’ current beliefs, resulting in an infinite recursion. However, in real life, studies have shown that humans could go no deeper than two levels of recursion (Camerer et al., 2004). Therefore, works (Fan et al., 2021) began the efforts to end the recursion when their beliefs merge into the “common mind”.

Modeling the belief of others has been extensively studied in symbolic-like environments (Wunder et al., 2011; Rabinowitz et al., 2018; Kleiman-Weiner et al., 2016; Ho et al., 2016), where agents need to incorporate or compete for a goal. Efforts to measure models’ ability to recognize false beliefs and perspective-taking also emerge in robotics (Yuan et al., 2020; Milliez et al., 2014), computer vision (Eysenbach et al., 2016; Fan et al., 2021), and natural language processing (Qiu et al., 2022; Nematzadeh et al., 2018) using the Sally-Anne test (Baron-Cohen et al., 1985). Different variants of the Sally-Anne test and ToM benchmarks are also proposed to test the ToM of large language models (Kosinski, 2023; Ullman, 2023; Sileo and Lernould, 2023; Kim et al., 2023). It is also shown that augmenting the model with external mind modules can help improve the performance of tasks involving intensive belief exchange and rich social interaction scenarios (Fan et al., 2021; Qiu et al., 2022; Li et al., 2023; Chiu et al., 2023). In this work, we explore ToM modeling can enhance the quality and efficiency of the re-

sponse generation in both cooperative and semi-cooperative dialogue tasks.

### Common ground alignment and negotiation

In a cooperative dialogue task, to guarantee that the communication takes the least cost meanwhile providing the most informative messages, previous works proposed multiple methods to align the common ground between agents (Bohn et al., 2019; Anderson, 2021). Specifically for dialogue tasks, datasets have been collected to provide golden utterances when people try to align the common ground with each other based on structured knowledge (He et al., 2017a), in partially observable cooperative tasks (Bara et al., 2021; Kim et al., 2019), in multimodal and continuous environment (Haber et al., 2019; Udagawa and Aizawa, 2021). Frameworks have been adopted to model and predict the aligning dynamics using GNN, RNN, transformers, and LLMs (He et al., 2017a; Udagawa and Aizawa, 2021; Fischer, 2023; Zhang et al., 2023; Zhou et al., 2023). The inferred common ground is also used to generate more interesting and engaging conversations for the dialogue agents (Zhou et al., 2022).

Negotiation is treated as a semi-cooperative task since agents can have different goals but need to agree on the same decision (Lewis et al., 2017). It requires complex social skills and strategies like offering proposals and accepting or making counter-offers (Yamaguchi et al., 2021). To improve the negotiating abilities of the dialogue systems, datasets of open-domain human negotiation corpus have been introduced in embodied environment (DeVault et al., 2015), daily items split (Lewis et al., 2017; Chawla et al., 2021), buy and sell (He et al., 2018), job offer negotiation (Yamaguchi et al., 2021). Modeling begins with game theory and action selection (Nash Jr, 1950; Baarslag et al., 2013). For open-domain generation, methods have been designed to help the model plan ahead (Lewis et al., 2017; Iwasa and Fujita, 2018), give feedback about the current conversation (Zhou et al., 2019; Fu et al., 2023), detect negotiation breakdowns (Yamaguchi et al., 2021).

## 3 Task and Framework

### 3.1 Tasks

The situated dialogue corpus can be denoted as  $\mathcal{D} = \{(U_n, K_n^p, y_n)\}_{n=1}^N$ , where  $U_n = (u_{n,1}, \dots, u_{n,T})$  represents the dialogue history and  $T$  is the number of turns.  $K_n^p = (k_{n,1}, \dots, k_{n,I})$  is for their knowledge base, where  $I$  is the number of

knowledge passages.  $p \in A, B$  represents the two agents. We assume the current speaker is A, and  $p$  will be dropped for the following formulation.  $y_n$  is A’s next response or its action to achieve the task goals.

**Alignment** In the common ground alignment scenarios, we use the MutualFriend task (He et al., 2017b) shown in Figure 2.  $K$  denotes the private friend lists that two agents observe, and there is only one friend shared in their lists. The agents need to merge their estimation of the mutual friend through chat and finally finish the task goal by taking the action to select  $k_i \in K$  as their mutual friend. The alignment is successful when their selections are the same.

**Negotiation** In the common ground negotiation scenarios, we use the CaSiNo task (Chawla et al., 2021) shown in Figure 2(Bottom). Two agents are planning a camp trip and need to distribute the uneven number of items. Based on their individual priority of the items  $K$ , they need to decide on the final item split agreement to maximize their gain of valuable items. At the end of the conversation, one agent proposes the item split deal while the other agent decides to accept or reject this deal. The negotiation is successful when the deal is accepted.

### 3.2 MindDial

The overall pipeline of our framework is shown in Figure 2. At the first stage, given the context history and private knowledge, the mind module  $f$  estimates the first and second-order beliefs over their solutions  $b_A, b_{BinA} = f(U, K)$ . The first-order belief represents A’s estimation of the mutual friend or split deal. The second-order belief refers to A’s understanding of B’s estimation regarding the mutual friend or split deal. We choose to prompt the LLMs for  $b_A$  and  $b_{BinA}$  due to their ability to adapt flexibly in open-domain corpora. Therefore, the mind module can be applied to other situated dialogues when the knowledge base and beliefs are well-defined.

Then the response generator  $h$  generates the next utterance based on the dialogue history, its private knowledge, and the intention to align the first and second-order beliefs:  $\tilde{y} = h(U, K, b_A, b_{BinA})$ . We apply two methods to the response generator to activate their ability to resolve the belief difference in  $b_A$  and  $b_{BinA}$ : embedding this ability into LLM by finetuning and explicitly triggering this ability of LLM by prompting.

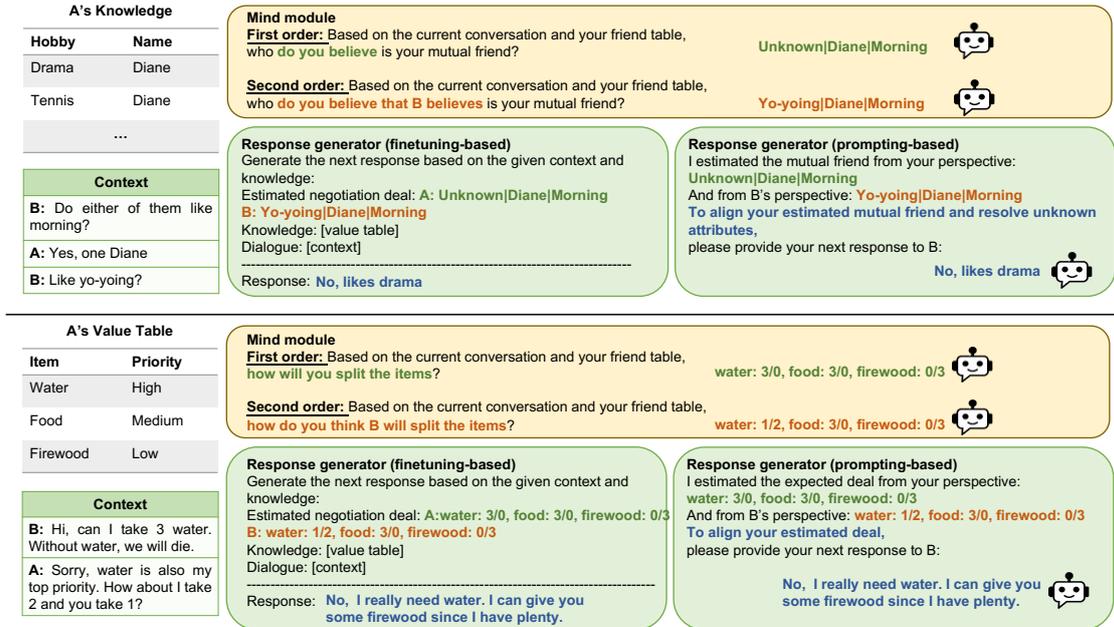


Figure 2: **Cases of ToM reasoning in MindDial.** Top: an *alignment* task from MutualFriend. Bottom: A *negotiation* task from CaSiNo. For each task, we first reason over the first- and second-order ToM beliefs of the conversational partner. Then we generate corresponding utterances wrt. the ToM estimation.

**Finetuning-based** For finetuning-based models, we prepare a small dataset in the format of  $\{y, U, K, b_A, b_{BinA}\}$ , where  $y$  is intended to resolve the gap between  $b_A$  and  $b_{BinA}$ . Different parts of model inputs are concatenated together with their corresponding tags as [Estimated belief], [Knowledge], and [Dialogue] shown in Figure 2. The models are trained to regress the next response  $y$ .

**Prompting-based** For prompting-based models, we directly ask the generator to generate the next response in order to resolve the difference and unknown values between  $b_A$  and  $b_{BinA}$ . The format follows as “I estimated mutual friend/deal from your perspective:  $b_A$  and from B’s perspective:  $b_{BinA}$ . To align  $b_A$  and  $b_{BinA}$ , please provide your next response to B:”.

## 4 Experiments

**Dataset** To provide a reasonable quantitative measure of belief dynamics in the dialogue, the expected dataset should contain rich belief exchanges. Meanwhile, the belief exchange and the final solution can be easily labeled. Therefore, we choose two representative settings to evaluate our framework. **MutualFriend** (He et al., 2017b): we consider it as an alignment dialogue scenario. In the MutualFriend task, each agent has a private knowledge base including a list of friends and their attributes like name, school, etc. There is a shared

friend that both agents have and they need to chat with each other to find this mutual friend. We only keep the successful dialogues and the final data split for train/val/test is 7257/878/900. Each dialogue in the training set contains a maximum of 53 turns and each turn with a maximum length of 29. **CaSiNo** (Chawla et al., 2021): we consider it as the negotiating scenario. In the CaSiNo task, two agents need to split camping packages: 3 water, 3 firewood, and 3 food. Each of these items will be of either High, Medium, or Low value to each agent. The agents need to negotiate the distribution of the items through chat to maximize their final points calculated based on the number of items they get and the items’ corresponding values. The data split for train/val/test is 900/30/100. Each dialogue in the training set contains a maximum of 27 turns and each turn with a maximum length of 106.

**Mind modules** To serve as a mind module in this task, the model is expected to understand long conversation contexts and the concept of first and second-order beliefs. We choose LLaMA-2-7B-chat, LLaMA-2-13B-chat (Touvron et al., 2023)<sup>1</sup>, GPT-3.5, and GPT-4<sup>2</sup> as our mind reasoner for their potentials in ToM benchmarks and the flexible abilities of mind reasoning in open-domain dialogues.

<sup>1</sup><https://github.com/facebookresearch/LLaMA-recipes/tree/main>

<sup>2</sup>gpt-3.5-turbo-1106, gpt-4-1106-preview

Models	Mind level	$C$	$T$	$C_T$	Models	Mind level	$C$	$T$	$C_T$
LLaMA-7B-ft	w/o mind	24.67	9.09	2.71	LLaMA-13B-ft	w/o mind	36.33	6.64	<b>5.47</b>
	$b_A$	28.33	7.92	<b>3.58</b>		$b_A$	42.00	8.66	4.85
	$b_{BinA}$	<b>29.33</b>	8.33	3.52		$b_{BinA}$	39.33	7.70	5.11
	$b_A+b_{BinA}$	28.33	8.87	3.20		$b_A+b_{BinA}$	<b>44.67</b>	8.85	5.05
GPT-3.5	w/o mind	10.67	5.74	1.86	GPT-4	w/o mind	75.00	9.72	7.71
	$b_A$	18.33	5.91	3.10		$b_A$	75.00	9.41	7.97
	$b_{BinA}$	12.33	5.91	2.09		$b_{BinA}$	69.67	8.84	7.88
	$b_A+b_{BinA}$	<b>24.33</b>	6.04	<b>4.03</b>		$b_A+b_{BinA}$	<b>76.00</b>	8.88	<b>8.56</b>

Table 1: **MutualFriend: results with different mind settings.** Settings without mind reasoning are marked as w/o mind. Settings considering only the first-order are marked as  $b_A$ , with only the second-order are  $b_{BinA}$ , with both are  $b_A+b_{BinA}$ .

**Response generators** We adopt the same four models in the mind modules as our response generators. We divide the models into two groups: finetuning and prompting-based. For the finetuning group, we first finetune LLaMA-2-7B-chat and LLaMA-2-13B-chat to generate the next response with the raw training dialogues. Then, we sample 3% of the training data and predict the first and second-order beliefs at each turn using the mind module, which are put into the dialogue context as additional information input to finetune the model again. We choose to combine only a small portion of training data input with beliefs to reduce the API query cost. We also vary the portion to 1%, 3%, and 5%. The sample size does not significantly influence the model performance (See Appendix E). For GPT-3.5 and GPT-4, we use prompts to regulate the conversation. For finetuning-based models, the models are trained on two A6000 GPUs for one epoch with an initial learning rate of  $1e-4$ . The batch size is set to 64. For prompting-based models, we use the OpenAI API for experiments.

#### 4.1 Evaluation and results

For evaluation, we focus on three main questions:

- **Question 1:** Can mind reasoning improve the common ground alignment and negotiation results?
- **Question 2:** Which level of beliefs contributes to the performance gain?
- **Question 3:** What is the relation between belief estimation accuracy and conversation outcomes?

**MutualFriend evaluation metrics** We adopt the same metrics in He et al. (2017b):

- Success rate ( $C$ ): how many dialogues where the two agents select the true mutual friend.
- Conversation turns ( $T$ ): the number of turns the agents take before the end of the conversation

- Success rate per turn ( $C_T$ ): how efficient the conversation is. We divide the overall success rate by the conversation turns.

**CaSiNo evaluation metrics** We follow the procedure in Lewis et al. (2017):

- Score-all: the average negotiation scores. The points each agent scores are the number of items times the item’s corresponding values. High priority is a value of 5. Medium is 4. Low is 3. If the deal is rejected or the negotiation exceeds the maximum turn, both agents receive 5 points. Since the best outcome should be a win-win situation, we also report the sum over the points of the two agents to compare the overall performance gain.
- Agreed %: the agreement of the deal. A deal is agreed when the agents agree on the proposed deal and the proposal does not exceed the total number of items the agents can distribute.
- Pareto: whether the deal is Pareto Optimal. A solution is Pareto Optimal if neither agent’s score can be improved without lowering the other’s score.
- Score-agreed: the average negotiation scores in agreed deals.

##### 4.1.1 Observation I: Mind reasoning improves conversation outcomes

First, our experiments compare models’ performance without and with mind reasoning. In the cooperative scenario in Table 1, comparing model+w/o mind rows with models, we can see that combining mind modules can significantly improve the alignment success rate in both finetuning and prompting-based models. Among them, GPT-4 performs the best, following LLaMA and GPT-3.5. As for efficiency, models with mind reasoning exhibit higher per-turn success. However, for LLaMA13b,

Models	Mind level	Score-all	Sum	Agreed %	Pareto	Score-agreed	Sum
LLaMA-7B-ft	w/o mind	8.10 vs 7.18	15.28	24.00	12.00	18.33 vs 14.50	32.83
	$b_A$	12.94 vs 13.48	<b>26.42</b>	<b>68.00</b>	20.00	16.68 vs 17.47	34.15
	$b_{BinA}$	11.76 vs 13.36	25.12	56.00	24.00	16.54 vs 19.29	<b>35.83</b>
	$b_A+b_{BinA}$	12.96 vs 12.70	25.66	62.00	<b>26.00</b>	17.84 vs 17.42	35.26
LLaMA-13B-ft	w/o mind	15.38 vs 12.68	28.06	70.00	24.00	19.83 vs 15.97	35.80
	$b_A$	18.02 vs 16.14	34.16	92.00	38.00	19.15 vs 17.11	36.26
	$b_{BinA}$	17.02 vs 14.50	31.52	82.00	30.00	19.66 vs 16.59	36.25
	$b_A+b_{BinA}$	17.36 vs 17.32	<b>34.68</b>	<b>92.00</b>	<b>40.00</b>	18.43 vs 18.39	<b>36.82</b>
GPT-3.5	w/o mind	15.00 vs 14.26	29.26	80.00	18.00	17.38 vs 16.57	33.95
	$b_A$	16.10 vs 17.08	33.18	90.00	22.00	17.22 vs 18.42	35.64
	$b_{BinA}$	16.72 vs 16.86	<b>33.58</b>	<b>92.00</b>	22.00	17.74 vs 17.89	35.63
	$b_A+b_{BinA}$	17.08 vs 15.18	32.26	86.00	<b>26.00</b>	19.05 vs 16.72	<b>35.77</b>
GPT-4	w/o mind	16.84 vs 16.90	33.74	94.00	8.00	17.60 vs 17.66	35.26
	$b_A$	16.72 vs 16.50	33.22	90.00	14.00	18.02 vs 17.78	35.80
	$b_{BinA}$	17.40 vs 16.56	33.96	92.00	12.00	18.17 vs 17.39	35.56
	$b_A+b_{BinA}$	17.54 vs 17.46	<b>35.00</b>	<b>96.00</b>	<b>20.00</b>	18.06 vs 17.98	<b>36.04</b>

Table 2: **CaSiNo: results with different mind settings.** Settings without mind reasoning are marked as w/o mind. Settings considering only the first-order are as  $b_A$ , with only the second-order are  $b_{BinA}$ , with both are  $b_A+b_{BinA}$ .

we notice a longer conversation length, therefore, the efficiency drops below the base model. This suggests that while incorporating belief estimation can elevate success rates, it may not necessarily enhance efficiency if acquiring additional information is needed to establish common ground. More comparison results and discussion can be found in Appendix F.

In the negotiation scenario, as referenced in Table 2, agents utilizing mind reasoning capabilities tend to achieve higher individual scores. Additionally, the collective points of both parties are increased. These agents also are more likely to reach agreements and achieve Pareto Optimal outcomes, suggesting a more strategic distribution of items. When comparing the points scored and the agreement rates across different models, GPT variants and LLaMA-13B display similar performances except that LLaMA-7B falls behind. Notably, LLaMA-13B achieves the highest Pareto Optimal scores, surpassing GPT-4. This may be attributed to GPT-4’s tendency to favor equitable item distribution, often resulting in a split like 1 and 1, with another item left unclaimed by either party.

#### 4.1.2 Observation II: Both two levels of belief contribute to the performance gain

Next, we assess the impact of varying belief estimation levels in the mind modules on model performance, as shown in Table 1 and Table 2. First, it is

evident that integrating any level of belief estimation leads to performance enhancements compared with the w/o mind baseline, indicating both first and second-order beliefs contribute to the response generation process. Within mind settings, in alignment scenarios, models underscoring the belief differences usually outperform others with single-order belief estimation in LLaMA-13B, GPT-3.5, and GPT-4. In negotiating settings, we first notice that the score-all strongly correlates with the agreed rate, and there is no consistent pattern. Examining Pareto Optimal outcomes, models aggregating both  $b_A$  and  $b_{BinA}$  tend to distribute items more effectively, resulting in higher Pareto Optimal scores. Similarly, in score-agreed, models combining both two levels of beliefs perform better.

We also notice some fluctuations in the results, for example, LLaMA-7B with only  $b_{BinA}$  achieves better results. We reckon that complex and intertwined effects can be exerted when 1) the model is bottlenecked by its context understanding and generation abilities and 2) one or both levels of the belief estimations are not accurate. In general, models need to take into account their own beliefs and also the beliefs of others. Focusing on resolving the differences between them can improve the common ground alignment accuracy and negotiation optimality.

**Robustness to prompts** In our experimental investigations (prompt templates are supplemented in Appendix B and C), we found that the perfor-

Models	Belief	Precision	F1
LLaMA-7B	$b_A$	33.00	33.00
	$b_{BinA}$	30.00	30.00
LLaMA-13B	$b_A$	36.00	34.00
	$b_{BinA}$	38.00	33.00
GPT-3.5	$b_A$	62.00	62.00
	$b_{BinA}$	70.00	67.00
GPT-4	$b_A$	77.00	77.00
	$b_{BinA}$	76.00	76.00

Table 3: **Belief prediction.** The precision and F1 when different models predict the first ( $b_A$ ) and second-order ( $b_{BinA}$ ) beliefs.

mance of belief prediction remains robust when prompts are structured to inquire about the current speaker’s solution and their estimation of the other speaker’s solution. Our comparison results in Tables 1 and 2 suggest that the task of one-hop prediction, encompassing beliefs and intentions, poses a minimal challenge for most LLMs. For instance, LLaMA-13B exhibits performance akin to GPT-3.5. Consequently, we assert that the primary challenge lies in advancing higher-level ToM inferences within these models.

#### 4.1.3 Observation III: Belief estimation accuracy positively correlates with the alignment success

To more convincingly validate that incorporating the mind reasoning module enhances the models’ task performance, we assessed the belief estimation accuracy when different models serve as the mind module in the MutualFriend task using LLaMA-7B, LLaMA-13B, GPT-3.5, and GPT-4. Subsequently, we examined how this accuracy correlates with task success when the four models function as response generators separately, paired with these four mind modules. Here, we demonstrate the relation between the belief estimation accuracy and the dialogue outcomes with MutualFriend task due to its clearly defined belief dynamics. Therefore, the first and second-order beliefs can be easily annotated using predefined rules. The detailed labeling process is included in Appendix A. Table 3 shows the precision and F1 scores for predicting the current speaker’s estimation of mutual friend given the current dialogue history  $b_A$ , and its estimation of the other speaker’s estimation  $b_{BinA}$ . The line plot in Figure 3 illustrates the corresponding success rates when a response generator is equipped with different mind modules.

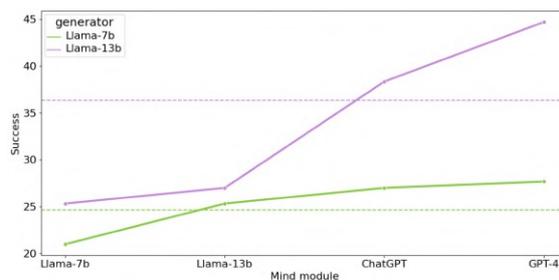


Figure 3: The task success rate when the response generators are paired with different models as the mind modules. The X-axis marks the model name of the mind modules. The Y-axis shows the success rate. Different colors represent different models as the response generators.

Combining the models’ precisions of the belief prediction with the success rates when they serve as the mind modules, we can observe that 1) The success rates increase when the models with higher belief prediction accuracy are served as the mind modules. This trend underscores that the effectiveness of the response generators is closely linked to the mind reasoning capabilities of the respective mind modules; 2) Comparing the growth magnitude of LLaMA-7B and 13b, we can see that LLaMA-7B reaches a flat stage and increases slowly. This suggests that the magnitude of the success rate improvement is bounded by the model’s mind-reasoning abilities; 3) The horizontal lines mark the task success rate when the response generators are not augmented with the mind modules. Augmenting models with weaker mind modules can detrimentally impact outcomes due to inaccurate belief predictions and inadequate dialogue reasoning, such as the situation when LLaMA-7B is paired with LLaMA-7B and LLaMA-13B is paired with LLaMA-7B and 13b.

**Summarization vs. Reasoning** It is worth noting that both the first- and second-order belief estimation goes beyond summarization from the last utterance. We carefully annotate part of the beliefs in the dialogue and report the second-order belief prediction accuracy in Table 3, which shows that the LLM can predict the second-order beliefs fairly well.

#### 4.1.4 Human Evaluation

We ask 16 college-level students to play Mutual-Friend and CaSiNo game with our model. Each subject is randomly assigned 4 samples. S/He chooses one sample to play with the agent w/o mind modules and the other one to play with the agent w/ mind modules. A pair-wise comparison is made

A friend list				B friend list			
hobby	school	name	loc_pref	hobby	school	name	loc_pref
swimming	national technological university	ryan	outdoor	cricket	university of wisconsin-stout	ryan	indoor
swimming	rowan college	amy	indoor	swimming	rowan college	amy	indoor
sand art	brown university	ryan	outdoor	polo	armstrong state college	jacqueline	indoor
cricket	alverno college	amy	outdoor	cricket	lynchburg college	lisa	indoor
swimming	pensacola christian college	peter	indoor	swimming	arizona western college	kathryn	indoor
worldbuilding	rowan college	peter	outdoor	cricket	lynchburg college	kathryn	indoor

<p>Hi. Know anyone who went to rowan col.</p> <p>4 of my friends like cricket.</p> <p>Do you know peter who went there and like worldbuilding</p> <p>No peter.</p> <p>Or any like worldbuilding but arent peter.</p> <p>...</p> <p>&lt;select&gt; kathryn lynchburg college cricket indoor. ❌</p>	<p>Hi. I have a friend that likes cricket playing named amy. (unknown unknown unknown unknown, unknown unknown unknown unknown)</p> <p>I have a friend named ryan from wisconsin stout that likes cricket. (cricket unknown unknown indoor, cricket unknown amy unknown)</p> <p>I do not know anyone with that name. (cricket unknown amy outdoor, cricket unknown unknown unknown)</p> <p>I have a friend named amy that likes swimming. (swimming unknown amy indoor, cricket unknown amy outdoor)</p> <p>My amy likes swimming and went to rowan college. (swimming unknown amy indoor, swimming unknown amy indoor)</p> <p>&lt;select&gt; amy rowan college swimming indoor ✅</p>
---	---

Figure 4: **Qualitative comparisons between dialogue generation models** without (at left) and with mind modeling (at right) when agents A and B are figuring out their mutual friend.

between the game outcome when human subjects play with models without and with mind reasoning. In addition, after the game ends, the subjects rate their game partner regarding their cooperativeness (whether the agent is cooperative during the game) and informativeness (whether the agents provide informative responses) from 0 to 10 in the alignment setting; rate regarding their negotiation skills (whether the agent is a good negotiator) and whether they are satisfied with the final deal in the negotiation setting. In addition, we also record their overall enjoyment when playing with the agents in both settings.

From Table 4, we can observe that our model with mind modules can achieve higher outcomes in both MutualFriend and CaSiNo games and the subjects tend to enjoy more in the process. In the cooperative setting, agents without and with mind achieves similar cooperativeness and informativeness rates. However, in the negotiation setting, agents with mind reasoning are shown to be more skillful and can achieve more satisfactory deals.

## 4.2 Case Study

We demonstrate one MutualFriend example to visualize the difference between LLaMA-7B with and without mind reasoning. Examples of other models and CaSiNo scenarios can be found in the Appendix. As shown in Figure 4, the topics between agents without mind reasoning can diverge quickly. For example, when A asks about “Rowan

College”, B responds with “cricket” which is unrelated to it. In contrast, for dialogues between agents with step-wise mind reasoning, they resolve the unknown attributes by providing related information (when A talks about “Amy” “swimming”, B mentions “Rowan College”). When there is a conflict between the names, A promptly negates “Ryan”.

Mutual Friend: alignment				
Groups	Success	Cooperative	Informative	Enjoyment
GPT-3.5 w/o mind	57.14	8.57	9.43	5.29
GPT-3.5 w/ mind	62.50	8.88	9.63	7.63

CaSiNo: negotiation				
Groups	Scores	Skillful	Satisfied	Enjoyment
GPT-3.5 w/o mind	22.50	6.25	6.50	5.75
GPT-3.5 w/ mind	24.50	7.13	7.25	7.25

Table 4: **Human study.** Comparisons are made between our model with mind module vs. models w/o mind module when played with human subjects.

## 5 Conclusion

In this study, we present MindDial, a novel framework for generating situated dialogue responses for common ground alignment and negotiation. By incorporating the first- and second-order ToM modeling into account, our model can enhance the alignment accuracy and negotiation outcome in both finetuning and prompting-based models. The efficacy of our approach is further substantiated through ablation studies and user feedback.

## Limitations

Our prompting design for the mind modules requires a well-defined knowledge and goal. This may limit the generalization abilities of the current framework to more casual conversation scenarios. Also, the task success is highly dependent on the belief estimation precision. Future research is needed to develop and implement mind modules that are both more robust and accurate.

## Acknowledgements

The authors thank Ms. Zhen Chen at BIGAI for designing the teaser figure. This work presented herein is supported by the National Science and Technology Major Project (2022ZD0114900) and the National Natural Science Foundation of China (62376031).

## References

- Carolyn Jane Anderson. 2021. Tell me everything you know: a conversation update system for the rational speech acts framework. In *Proceedings of the Society for Computation in Linguistics 2021*, pages 244–253.
- Tim Baarslag, Katsuhide Fujita, Enrico H Gerding, Koen Hindriks, Takayuki Ito, Nicholas R Jennings, Catholijn Jonker, Sarit Kraus, Raz Lin, Valentin Robu, et al. 2013. Evaluating practical negotiating agents: Results and analysis of the 2011 international competition. *Artificial Intelligence*, 198:73–103.
- Chris Baker, Rebecca Saxe, and Joshua Tenenbaum. 2011. Bayesian theory of mind: Modeling joint belief-desire attribution. In *Annual Meeting of the Cognitive Science Society (CogSci)*, volume 33.
- Cristian-Paul Bara, Sky CH-Wang, and Joyce Chai. 2021. *MindCraft: Theory of mind modeling for situated dialogue in collaborative tasks*. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1112–1125, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Simon Baron-Cohen, Alan M Leslie, and Uta Frith. 1985. Does the autistic child have a “theory of mind”? *Cognition*, 21(1):37–46.
- Max H Bazerman, Jared R Curhan, Don A Moore, and Kathleen L Valley. 2000. Negotiation. *Annual review of psychology*, 51(1):279–314.
- Pieter J Beers, Henny PA Boshuizen, Paul A Kirschner, and Wim H Gijssels. 2006. Common ground, complex problems and decision making. *Group decision and negotiation*, 15:529–556.
- Reinhard Blutner. 2000. Some aspects of optimality in natural language interpretation. *Journal of semantics*, 17(3):189–216.
- Manuel Bohn, Michael Henry Tessler, and Michael C Frank. 2019. Integrating common ground and informativeness in pragmatic word learning.
- Torben Braüner, Patrick Rowan Blackburn, and Irina Polyanskaya. 2016. Recursive belief manipulation and second-order false-beliefs. In *38th Annual Cognitive Science Society Meeting*, pages 2579–2584. Cognitive Science Society.
- 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 (NeurIPS)*, 33:1877–1901.
- Brant R Burleson. 2007. Constructivism: A general theory of communication skill. *Explaining communication: Contemporary theories and exemplars*, pages 105–128.
- Colin F Camerer, Teck-Hua Ho, and Juin-Kuan Chong. 2004. A cognitive hierarchy model of games. *The Quarterly Journal of Economics*, 119(3):861–898.
- Kushal Chawla, Jaysa Ramirez, Rene Clever, Gale Lucas, Jonathan May, and Jonathan Gratch. 2021. *CaSiNo: A corpus of campsite negotiation dialogues for automatic negotiation systems*. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3167–3185, Online. Association for Computational Linguistics.
- Justin Chiu, Wenting Zhao, Derek Chen, Saujas Vaduguru, Alexander Rush, and Daniel Fried. 2023. *Symbolic planning and code generation for grounded dialogue*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7426–7436, Singapore. Association for Computational Linguistics.
- Herbert H Clark and Deanna Wilkes-Gibbs. 1986. Referring as a collaborative process. *Cognition*, 22(1):1–39.
- Harmen De Weerd, Rineke Verbrugge, and Bart Verheij. 2013. How much does it help to know what she knows you know? an agent-based simulation study. *Artificial Intelligence*, 199:67–92.
- Harmen De Weerd, Rineke Verbrugge, and Bart Verheij. 2015. Higher-order theory of mind in the tacit communication game. *Biologically Inspired Cognitive Architectures*, 11:10–21.
- David DeVault, Johnathan Mell, and Jonathan Gratch. 2015. Toward natural turn-taking in a virtual human negotiation agent. In *2015 AAAI Spring Symposium Series*.

- Prashant Doshi, Xia Qu, Adam Goodie, and Diana Young. 2010. Modeling recursive reasoning by humans using empirically informed interactive pomdps. In *International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 1223–1230.
- Benjamin Eysenbach, Carl Vondrick, and Antonio Torralba. 2016. Who is mistaken? *arXiv preprint arXiv:1612.01175*.
- Lifeng Fan, Shuwen Qiu, Zilong Zheng, Tao Gao, Song-Chun Zhu, and Yixin Zhu. 2021. Learning triadic belief dynamics in nonverbal communication from videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7312–7321.
- Kevin A Fischer. 2023. Reflective linguistic programming (rlp): A stepping stone in socially-aware agi (socialagi). *arXiv preprint arXiv:2305.12647*.
- Yao Fu, Hao Peng, Tushar Khot, and Mirella Lapata. 2023. Improving language model negotiation with self-play and in-context learning from ai feedback. *arXiv preprint arXiv:2305.10142*.
- Sebastian Grueneisen, Emily Wyman, and Michael Tomasello. 2015. “i know you don’t know i know. . .” children use second-order false-belief reasoning for peer coordination. *Child Development*, 86(1):287–293.
- Janosch Haber, Tim Baumgärtner, Ece Takmaz, Lieke Gelderloos, Elia Bruni, and Raquel Fernández. 2019. [The PhotoBook dataset: Building common ground through visually-grounded dialogue](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1895–1910, Florence, Italy. Association for Computational Linguistics.
- Yanlin Han and Piotr Gmytrasiewicz. 2018. Learning others’ intentional models in multi-agent settings using interactive pomdps. *Advances in Neural Information Processing Systems (NeurIPS)*, 31.
- He He, Anusha Balakrishnan, Mihail Eric, and Percy Liang. 2017a. [Learning symmetric collaborative dialogue agents with dynamic knowledge graph embeddings](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1766–1776, Vancouver, Canada. Association for Computational Linguistics.
- He He, Anusha Balakrishnan, Mihail Eric, and Percy Liang. 2017b. Learning symmetric collaborative dialogue agents with dynamic knowledge graph embeddings. *Annual Meeting of the Association for Computational Linguistics (ACL)*.
- He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. 2018. [Decoupling strategy and generation in negotiation dialogues](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2333–2343, Brussels, Belgium. Association for Computational Linguistics.
- Mark K Ho, James MacGlashan, Amy Greenwald, Michael L Littman, Elizabeth Hilliard, Carl Trimbach, Stephen Brawner, Josh Tenenbaum, Max Kleiman-Weiner, and Joseph L Austerweil. 2016. Feature-based joint planning and norm learning in collaborative games. In *Annual Meeting of the Cognitive Science Society (CogSci)*.
- Kosui Iwasa and Katsuhide Fujita. 2018. Prediction of nash bargaining solution in negotiation dialogue. In *Pacific Rim International Conference on Artificial Intelligence*, pages 786–796. Springer.
- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Bras, Gunhee Kim, Yejin Choi, and Maarten Sap. 2023. [FANToM: A benchmark for stress-testing machine theory of mind in interactions](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14397–14413, Singapore. Association for Computational Linguistics.
- Jin-Hwa Kim, Nikita Kitaev, Xinlei Chen, Marcus Rohrbach, Byoung-Tak Zhang, Yuandong Tian, Dhruv Batra, and Devi Parikh. 2019. [CoDraw: Collaborative drawing as a testbed for grounded goal-driven communication](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6495–6513, Florence, Italy. Association for Computational Linguistics.
- Max Kleiman-Weiner, Mark K Ho, Joseph L Austerweil, Michael L Littman, and Joshua B Tenenbaum. 2016. Coordinate to cooperate or compete: abstract goals and joint intentions in social interaction. In *Annual Meeting of the Cognitive Science Society (CogSci)*.
- Michal Kosinski. 2023. Theory of mind may have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*.
- Ágnes Melinda Kovács, Ernő Téglás, and Ansgar Denis Endress. 2010. The social sense: Susceptibility to others’ beliefs in human infants and adults. *Science*, 330(6012):1830–1834.
- Mike Lewis, Denis Yarats, Yann Dauphin, Devi Parikh, and Dhruv Batra. 2017. [Deal or no deal? end-to-end learning of negotiation dialogues](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2443–2453, Copenhagen, Denmark. Association for Computational Linguistics.
- Hua Li, Yu Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Charles Lewis, and Katia Sycara. 2023. [Theory of mind for multi-agent collaboration via large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 180–192, Singapore. Association for Computational Linguistics.
- Grégoire Milliez, Matthieu Warnier, Aurélie Clodic, and Rachid Alami. 2014. A framework for endowing an interactive robot with reasoning capabilities about perspective-taking and belief management. In *The 23rd IEEE international symposium on robot and*

- human interactive communication*, pages 1103–1109. IEEE.
- John F Nash Jr. 1950. The bargaining problem. *Econometrica: Journal of the econometric society*, pages 155–162.
- Aida Nematzadeh, Kaylee Burns, Erin Grant, Alison Gopnik, and Tom Griffiths. 2018. [Evaluating theory of mind in question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2392–2400, Brussels, Belgium. Association for Computational Linguistics.
- Liang Qiu, Yizhou Zhao, Yuan Liang, Pan Lu, Weiyan Shi, Zhou Yu, and Song-Chun Zhu. 2022. [Towards socially intelligent agents with mental state transition and human value](#). In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 146–158, Edinburgh, UK. Association for Computational Linguistics.
- Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, SM Ali Eslami, and Matthew Botvinick. 2018. Machine theory of mind. In *International Conference on Machine Learning (ICML)*, pages 4218–4227. PMLR.
- 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.
- Hilary Richardson, Grace Lisandrelli, Alexa Riobueno-Naylor, and Rebecca Saxe. 2018. Development of the social brain from age three to twelve years. *Nature communications*, 9(1):1027.
- Damien Sileo and Antoine Lerneuld. 2023. Mindgames: Targeting theory of mind in large language models with dynamic epistemic modal logic. *arXiv preprint arXiv:2305.03353*.
- Michael Tomasello. 2010. *Origins of human communication*. MIT press.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Takuma Udagawa and Akiko Aizawa. 2021. [Maintaining common ground in dynamic environments](#). *Transactions of the Association for Computational Linguistics*, 9:995–1011.
- Tomer Ullman. 2023. Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv preprint arXiv:2302.08399*.
- Elizabeth S Veinott, Judith Olson, Gary M Olson, and Xiaolan Fu. 1999. Video helps remote work: Speakers who need to negotiate common ground benefit from seeing each other. In *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*, pages 302–309.
- Adam Vogel, Max Bodoia, Christopher Potts, and Dan Jurafsky. 2013. Emergence of gricean maxims from multi-agent decision theory. In *Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies*, pages 1072–1081.
- Michael Wunder, Michael Kaisers, John Robert Yaros, and Michael L Littman. 2011. Using iterated reasoning to predict opponent strategies. In *International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 593–600.
- Atsuki Yamaguchi, Kosui Iwasa, and Katsuhide Fujita. 2021. [Dialogue act-based breakdown detection in negotiation dialogues](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 745–757, Online. Association for Computational Linguistics.
- Tao Yuan, Hangxin Liu, Lifeng Fan, Zilong Zheng, Tao Gao, Yixin Zhu, and Song-Chun Zhu. 2020. Joint inference of states, robot knowledge, and human (false-) beliefs. In *International Conference on Robotics and Automation (ICRA)*, pages 5972–5978. IEEE.
- Hongxin Zhang, Weihua Du, Jiaming Shan, Qinrong Zhou, Yilun Du, Joshua B Tenenbaum, Tianmin Shu, and Chuang Gan. 2023. Building cooperative embodied agents modularly with large language models. *arXiv preprint arXiv:2307.02485*.
- Pei Zhou, Hyundong Cho, Pegah Jandaghi, Dong-Ho Lee, Bill Yuchen Lin, Jay Pujara, and Xiang Ren. 2022. [Reflect, not reflex: Inference-based common ground improves dialogue response quality](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10450–10468, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Pei Zhou, Aman Madaan, Srividya Pranavi Potharaju, Aditya Gupta, Kevin R McKee, Ari Holtzman, Jay Pujara, Xiang Ren, Swaroop Mishra, Aida Nematzadeh, et al. 2023. How far are large language models from agents with theory-of-mind? *arXiv preprint arXiv:2310.03051*.
- Yiheng Zhou, He He, Alan W Black, and Yulia Tsvetkov. 2019. [A dynamic strategy coach for effective negotiation](#). In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 367–378, Stockholm, Sweden. Association for Computational Linguistics.

## A MutualFriend belief annotation and evaluation

To test the belief estimation accuracy of our mind modules, we manually label the first and second-order beliefs given the current context of the dialogues. The values mentioned in the current dialogue context are marked as positive (1). The values not mentioned or negated by either of the agents are marked as negative (0). When all the values of one attribute are marked as negative, this attribute becomes “unknown”. Figure 5 illustrate one annotation process. For example, when B is asking about “yo-yoing”, this value is marked as 1 for  $b_{BinA}$  hobby. However, since it does not belong to A’s knowledge, for the first-order belief of speaker A, we annotate it as 0. Then, when “yo-yoing” is negated by A, it will be marked as 0 in  $b_{BinA}$ . The prediction is a true positive when the model’s predicted value of one attribute is annotated as 1, a true negative when both prediction and ground truth are “unknown”.

## B MutualFriend prompts

→ (At the beginning of the first turn): You are a smart cooperative agent named [AlicelBob]. You have many friends with different attributes as listed below (the knowledge base of [AlicelBob]). You are now talking with Bob. He also has a list of friends. You will talk with Bob for a maximum of 20 turns to find out your mutual friend as quickly as possible. You can ask him questions or provide information about your friends. Meanwhile, you should try to mention as few attributes and friends as possible.

hobby, name, location

Surfing, Jane, Outdoor

...

(After each turn - no mind):

→ [AlicelBob] said: {last generated response}. Please provide your next utterance to [AlicelBob]:

→ Have you found your mutual friend? If yes, provide this mutual friend in the format of hobby|name|location; If no, respond 'unknown':

(After each turn - with mind):

→ (first-order) Based on the current conversation and your friend table, who do you believe is

your mutual friend? Respond in the format of hobby|name|location, and put unknown in the attributes you are not sure about for now:

→ (second-order) Based on the current conversation and your friend table, who do you believe that [AlicelBob] believes your mutual friend is? Respond in the format of hobby|name|location, and put unknown in the attributes you are not sure about for now:

→ [AlicelBob] said: {last generated response}. I estimate the mutual friend estimation from your perspective: [first-order] and from [AlicelBob]’s perspective: [second-order] based on your current talk. To align your estimation and resolve unknown attributes, please provide your next utterance to [AlicelBob]:

→ Have you found your mutual friend? If yes, provide this mutual friend in the format of hobby|name|location; If no, respond 'unknown':

Figure 6: Template for MutualFriend self-talk prompt.

## C CaSiNo prompts

→ (At the beginning of the first turn): You are a smart negotiation agent named [AlicelBob] planning a camping trip. Besides basic supplies, you will need extra water, food, and firewood. Each of these items will be of either High, Medium, or Low priority for you as shown below. Each of them only has an available quantity of 3 and can only be split using integers. You will negotiate with Bob who will also need these items and have his own value table. Use reasons from your value table to justify why you need these items. Try hard to get as many items as you can!

Item, value, reason

water, high, I didn’t pack enough water

...

(After each turn - no mind):

→ [AlicelBob] said: {last generated response}. Please provide your next utterance to [AlicelBob]:

→ Based on your conversation with [AlicelBob], do you want to end the negotiation? Please respond by yes or No:

(After each turn - with mind):

A's Knowledge			B's Knowledge			Context	
Hobby	Name	Time Pref	Hobby	Name	Time Pref		
Drama	Diane	Morning	Yo-yoing	Diane	Morning	<b>B:</b> Do either of them like yo-yoing?	
Tennis	Diane	Evening	Drama	Diane	Evening	<b>A:</b> Nopes. no yo-yoing	
...			...			<b>B:</b> Any like drama and name is Diane?	

**Annotated Beliefs**

Turn	$g_A$	$g_{B A}$	$g_B$	$g_{A B}$
1	Hobby: unknown Name: unknown Time: unknown	Hobby: yo-yoing Name: unknown Time: unknown	-	-
2			Hobby: unknown Name: unknown Time: unknown	Hobby: unknown Name: unknown Time: unknown
3	Hobby: drama Name: Diane Time: unknown	Hobby: drama Name: Diane Time: unknown		

Figure 5: Annotation example

- (first-order) Based on the current conversation and your value table, how will you split water, firewood, and food? The items each person gets can only be integers and the total quantity for each item is 3. Please use the following format to respond without further explanation: item: the number you get/the number [AlicelBob] get. For example, water:0/3, firewood:1/2, food: 3/0.
- (second-order) Based on the current conversation and your value table, how do you think [AlicelBob] will split water, firewood, and food? The items each person gets can only be integers and the total quantity for each item is 3. Please use the following format to respond without further explanation: item: the number you get/the number [AlicelBob] get. For example, water:0/3, firewood:1/2, food: 3/0.
- [AlicelBob] said: {last generated response}. I estimated the negotiation deal from your perspective: [first-order] and from Bob's perspective: [second-order] based on your current talk. To align your expected deals, please provide your next utterance to [AlicelBob]:
- Based on your conversation with [AlicelBob], do you want to end the negotiation? Please respond by yes or No:  
(After negotiation ends):
- Please provide your proposed deal. The items each person gets can only be integers and the total quantity for each item is 3. Deal with fractions will be rejected. Please use the following format: item: the number you get/the number [AlicelBob] get. For example, water:0/3, firewood:1/2, food: 3/0.
- Given your current conversation and the deal proposed by [AlicelBob]: [deal], will you accept the deal? Please respond by Accept or Reject:

Figure 7: Template for CaSiNo self-talk prompt.

## D Finetuning data format

Generate the next response of the dialog based on the given context and knowledge:

(SPEAKER0 as the current speaker)  
 Estimated [mutual friend|negotiation deal]  
 [SPEAKER0] [First-order belief]  
 [SPEAKER1] [Second-order belief]  
 Knowledge:  
 Friend table or value table  
 Dialogues:  
 [SPEAKER0] ...  
 [SPEAKER1] ...  
 — response:

Figure 8: Template for Finetuning

Models	Sample size	Score-all	Sum	Agreed %	Pareto	Score-agreed	Sum
LLaMA-13B-ft	w/o mind	15.38 vs 12.68	28.06	70.00	24.00	19.83 vs 15.97	35.80
	1%	15.36 vs 15.50	30.86	80.00	30.00	18.28 vs 18.46	36.74
	3%	17.36 vs 17.32	34.68	92.00	40.00	18.43 vs 18.39	36.82
	5%	16.44 vs 16.58	33.02	86.00	34.00	18.30 vs 18.47	36.77

Table 5: **CaSiNo: results with different sample sizes.**

Models	Sample size	$C$	$T$	$C_T$
LLaMA-13B-ft	w/o mind	36.33	6.64	5.47
	1%	38.46	8.80	4.37
	3%	44.67	8.85	5.05
	5%	40.33	8.53	4.73

Table 6: **MutualFriend: results with different sample size.**

Models	Mind level	$C$	$C_T$
Human	-	82.00	7.00
Rule	-	90.00	5.00
StanoNet	-	78.00	4.00
DynoNet	-	<b>96.00</b>	6.00
LLaMA-7B-ft	w/o mind	24.67	2.71
LLaMA-7B-ft	$b_A + b_{BinA}$	28.33	3.20
LLaMA-13B-ft	w/o mind	36.33	5.47
LLaMA-13B-ft	$b_A + b_{BinA}$	44.67	5.05
GPT-3.5	w/o mind	10.67	1.86
GPT-3.5	$b_A + b_{BinA}$	24.33	4.03
GPT-4	w/o mind	75.00	7.71
GPT-4	$b_A + b_{BinA}$	76.00	<b>8.56</b>

Table 7: **MutualFriend: comparison with results from original paper.**

## E Varying sample size of mind annotation data during finetuning

Considering the computational cost during finetuning, we only sample a small partition of dialogue for mind augmentation. In this section, we vary the sample size by 1%, 3% and 5%. From Table 6 and Table 5, we can see that 5% achieves the best results and all models perform better than the w/o mind baselines.

## F MutualFriend: more comparison results

In this section, we provide the baseline results of MutualFriend from the original paper in Table 7. It is shown that GPT-4 can achieve higher efficiency with higher accuracy per turn. It is worth noting that the models in the original paper are of smaller sizes and trained with specific datasets while we cur-

rently focus more on larger models generalizable to more open-domain tasks. The CaSiNo dataset was originally designed for the strategy prediction task, therefore it did not report generation results.

# An Open Intent Discovery Evaluation Framework

Grant Anderson<sup>1,2</sup>, Emma Hart<sup>1</sup>, Dimitra Gkatzia<sup>1</sup>, Ian Beaver<sup>2</sup>

<sup>1</sup>Edinburgh Napier University, UK, <sup>2</sup>Verint Systems Ltd., USA

Correspondence: [grant.anderson@verint.com](mailto:grant.anderson@verint.com)

## Abstract

In the development of dialog systems the discovery of the set of target intents to identify is a crucial first step that is often overlooked. Most intent detection works assume that a labelled dataset already exists, however creating these datasets is no trivial task and usually requires humans to manually analyse, decide on intent labels and tag accordingly. The field of Open Intent Discovery (OID) addresses this problem by automating the process of grouping utterances and providing the user with the discovered intents. Our OID framework allows for the user to choose from a range of different techniques for each step in the discovery process, including the ability to extend previous works with a human-readable label generation stage. We also provide an analysis of the relationship between dataset features and optimal combination of techniques for each step to help others choose without having to explore every possible combination for their unlabelled data.

## 1 Introduction

A major first task for a goal-oriented dialogue system is to identify the intent behind the user's utterance using a Natural Language Understanding module. This module is often implemented as a classifier, trained on a set of pre-defined intent labels (Chen et al., 2013; Coucke et al., 2018; Goo et al., 2018; Kim et al., 2016; Liu and Lane, 2016; Zhong and Li, 2019). Discovering these intents in real-world systems can be a laborious and time-consuming task involving a domain expert exploring the dataset and curating a representative set of labels. This task will also need to be repeated regularly as new intents emerge through time. The field of OID seeks to automatically discover unknown intents in a set of unlabelled/partially labelled utterances without requiring such manual effort.

There exists an issue in the current literature in that many works focus only on the development of clustering algorithms to identify utterances of

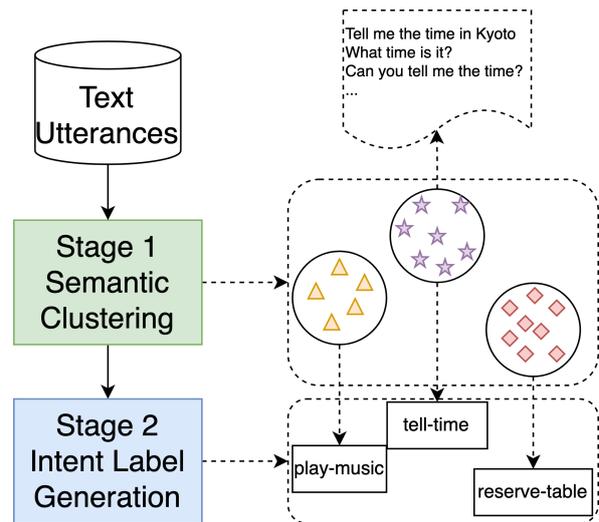


Figure 1: An example of the automated discovery and labelling of intents in a given dataset of unlabelled/partially labelled text utterances. First, the utterances are clustered for similar semantic intent, then human-readable labels are generated for each identified cluster.

similar intent, without progressing to label the cluster with a human-readable intent label (Perkins and Yang, 2019; Lin et al., 2019; Zhang et al., 2021b; Shen et al., 2021; Kumar et al., 2022). In order for downstream systems to make full use of the new intents, a human would be required to analyse the cluster manually, decide on its meaning and label it accordingly.

Evaluation methods are also inconsistent across the field. Some works report on classification or clustering metrics while others evaluate quality of generated labels, but rarely are these reported for the same datasets. There are also differences in the definition of 'intent' and the features of the datasets used for evaluation. Some works consider intents in the abstract such as 'CustomerService' or 'Baggage' in the air travel domain. Other works take a much stricter definition e.g. only an Action(verb)-Object(noun) pair. Some datasets contain a mixture

of these intent types. These issues make it difficult to identify a truly state-of-the-art (SOTA) technique for different domains and features of dataset.

We present an OID framework which views OID as a two stage process: Semantic Clustering, and Intent Label Generation (see Figure 1). We focus on the generation of high quality labels for an unlabelled/partially labelled dataset, produced by combining a semantic representation method, clustering algorithm, candidate extraction method and a label selection method. We evaluate 20 combinations of representation/clustering/extraction/selection methods on 9 datasets. Our key contributions include: (1) We introduce our novel OID framework providing a choice of a number of different techniques at every step in the process.<sup>1</sup> (2) We extend previous OID work to include a human-readable intent stage. (3) A rigorous investigation into instantiating choice of representation/clustering model/extraction/selection which reveals the optimal settings for datasets and target intents.

## 2 Related Work

State-of-the-art OID techniques utilise semi-supervised learning such as in DSSCC (Deep Semi-Supervised Contrastive Clustering) (Kumar et al., 2022) and DeepAligned (Zhang et al., 2021b). A portion of intents are known in advance and these are used to aid the clustering stage in discovering both the known intent clusters and estimate a number of new, unknown intents. Shen et al. (2021) take a different approach, by pre-training a representation model with a labelled dataset from the same domain as the target unlabelled dataset and then using unsupervised KMeans clustering on the target dataset to discover intents.

There are several works which attempt to solve the problem in an unsupervised fashion. Chatterjee and Sengupta (2020) adapted the DBSCAN clustering algorithm (Ester et al., 1996) in an attempt to handle discovering new intents in datasets with unbalanced distributions, while others such as Liu et al. (2021) use simple KMeans clustering. Liu et al. (2021) are one of the few OID works which include a label generation stage. Each cluster has candidate intent labels extracted using a dependency parser to find Action(verb)-Object(noun) pairs within the utterances and the most common pair is assigned as an auto-generated, human-readable label for the cluster. Their technique

discovered the correct number of clusters for the SNIPS dataset and produced labels which were clearly semantically similar to the ground-truth intents, however no quantitative evaluation was conducted. A more challenging dataset would prove more difficult both to cluster and to evaluate by manual inspection. Vedula et al. (2020) looked at intent discovery as a sequence tagging task. A neural model sequence tagger is trained to tag action and object words in text utterances. This technique differs in that it will produce an intent for every text utterance and may produce many distinct pairs that express the same intent.

In our concurrent work, we presented experimental results for different combinations of candidate extraction and intent label selection techniques against a large generative PLM (Anderson et al., 2024). In order to produce fine-grained intents, we also proposed an extension to the Action-Object extraction method used in Liu et al. (2021) which captures more detail from the utterances by including compound nouns or adjectives that are related to the Object, and negations related to the Action.

Zhang et al. (2021a) introduced a platform for open intent recognition. They combine the related tasks of open intent detection and discovery to both identify the known intents and discover new ones. The detection module identifies known intent samples and groups unknown samples into a single class of open intent. The discovery module then performs clustering to group the unknown samples and present them as new intents. Our framework differs in that we focus only on discovery and not detection. We also include a human-readable label generation stage while TEXTOIR provides keywords to represent their discovered intents. These keywords are helpful, however, out of context they would be difficult to fully understand without further analysing the utterances themselves.

## 3 Methods

Many current OID techniques can fit into the same two stage pattern (see Figure 2). Stage 1 consists of semantic clustering and is split into two steps. First, semantic representations are obtained for each utterance, then these are grouped with a clustering algorithm to identify semantically similar intents. Stage 2 involves the generation of a natural language label for each cluster. First, candidate labels are extracted or generated for each cluster, finally, a label is chosen from these candidates.

<sup>1</sup><https://github.com/GAnderson01/open-intent-discovery>

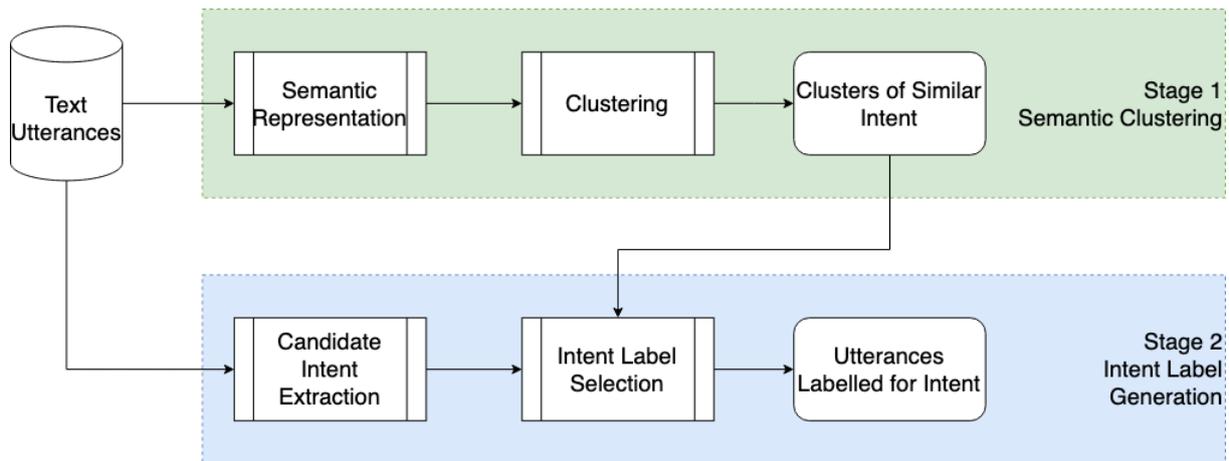


Figure 2: The Open Intent Discovery Framework is split into two main stages. In Stage 1, utterances are clustered for semantic similarity and in Stage 2, a human-readable label is produced for each cluster.

At each step in the process there are many different options for a researcher to choose from. At the Semantic Representation step, choices include using BERT, Universal Sentence Encoder or one of many other embedding options. For clustering, one could choose KMeans, DBSCAN etc. When looking for candidate labels, possibilities include an extraction method, such as the Action-Object extraction used by Liu et al. (2021), or a label could be generated by prompting a Pre-trained Language Model (PLM) such as ChatGPT or T0pp. Finally, a label must be chosen from the candidates e.g. by choosing the most frequent candidate or even by prompting a PLM, specifying the candidates to choose from. Our framework allows for any combination of options to be evaluated. Table 1 displays the different options we explored for each step in the framework. We refer to a combination of semantic representation, clustering, candidate extraction and intent label selection techniques as a configuration.

Most related works do not progress to Stage 2, and simply present the clusters of semantically similar texts as the found intents. Using the framework, we are able to extend these with Stage 2 techniques allowing us to evaluate the quality of the final natural language labels for clusters found by all OID techniques. For each cluster, we measure both the cosine similarity and the BARTScore between the most common ground truth label in the cluster and the generated label.

One of the goals of this work is to find common patterns in the configurations for datasets with similar features. It is hoped that this will help others to choose the best configuration for their own datasets

rather than having to perform a brute force search, or best guess.

The framework implements each step as a python module. Each can be run individually provided they are given any input required. When chained together, they execute the entire OID process end-to-end.

### 3.1 Stage 1: Semantic Clustering

The first stage is to collect the text utterances into groups of semantically similar intent. To achieve this, we first need to obtain good semantic representations of the utterances via some embedding model, then provide these to a clustering algorithm.

**Semantic Representation** Using PLMs to obtain embeddings for text utterances before applying these in a downstream NLP task has been repeatedly shown to perform well. However, the question of which PLM to use for a particular problem and dataset can be unclear. The semantic representation module supports any huggingface, sentence-transformers or tensorflow-hub based PLM embedding models. We use three PLMs to obtain semantic representations for the utterances in order to cluster for intent. These are as follows: bert-base-uncased (Devlin et al., 2018), all-mpnet-base-v2 (Reimers and Gurevych, 2019) and Universal Sentence Encoder (Cer et al., 2018). These PLMs have been shown to perform well in previous OID works (Zhang et al., 2021b; Kumar et al., 2022; Liu et al., 2021; Chatterjee and Sengupta, 2020).

**Clustering** The optimal clustering algorithm to use for a given dataset depends on the features of the dataset. For example, KMeans is more suited to

Stage 1: Semantic Clustering		Stage 2: Intent Label Generation	
Semantic Representation	Clustering Algorithm	Candidate Extraction	Intent Label Selection
all-mpnet	KMeans	Action-Object Pairs	Most Frequent
BERT	DBSCAN	T0pp Prompting	T0pp Prompting
Universal Sentence Encoder	ITER_DBSCAN		
	DeepAligned		

Table 1: Evaluated choices at each step of the framework

finding clusters of similar sizes (a balanced dataset), and a flat geometry, while density based methods such as DBSCAN can handle uneven cluster sizes (an imbalanced dataset) and non-flat geometry. We explore both unsupervised (KMeans, DBSCAN and ITER\_DBSCAN) and semi-supervised (DeepAligned) intent clustering algorithms. Both ITER\_DBSCAN and DeepAligned are intent discovery techniques which do not involve creating human readable labels, and so our framework extends them with the Stage 2 label generation techniques.

Most clustering algorithms require some hyperparameters to be set e.g. KMeans requires the target number of clusters ( $k$ ). However in many cases these hyperparameters are unknown and so a tuning exercise is required. In order to find optimal hyperparameters, a search across the hyperparameter space must be conducted and each clustering result evaluated against some metric. This metric, is one of the choices that can be set in the framework.

### 3.2 Stage 2: Intent Label Generation

The second stage is to choose or generate a natural language label to represent the cluster as an intent. First, candidates are found from the cluster either using a dependency parser or prompting a PLM, then one of the candidates is selected by some method such as most frequent, or, again, prompting a PLM.

**Candidate Label Extraction** We implement two techniques to extract candidates intents for the identified clusters. The first finds Action-Object pairs in utterances as in (Liu et al., 2021). An Action-Object pair consists of a verb/infinitive (the Action) and it’s target, a noun or subject (the Object). e.g. “schedule a meeting for tomorrow” contains the Action-Object pair *schedule-meeting*. If either an action or object is not present in an utterance, then the candidate contains ‘NONE’ in it’s place. This technique assumes a very strict definition of intent and as such, could never produce a more abstract

intent such as ‘query’ or ‘confirmation’. Therefore, we also experiment with PLM Prompting, to allow for more freedom in the candidate intents.

To produce a candidate with a PLM, we obtain the response when it is given the below prompt:

*“Given the following utterance: [utterance]. The intent was to”*

**Intent Label Selection** The final step in the framework is to choose an intent label for every cluster from one of the candidates identified. We experiment with two techniques. As in Liu et al. (2021), we choose the most frequent candidate. Where Action-Object extraction was used we ignore incomplete pairs by not considering any with the word ‘NONE’. If a cluster produced no candidates, then no label will be generated. The second selection technique also prompts a PLM using the following:

*“Given these utterances: [cluster\_utterances]. What is the best fitting intent, if any, among the following: [top\_3\_candidates]”*

where [cluster\_utterances] is all of the utterances present in the cluster and [top\_3\_candidates] are the three most common candidates in the cluster. This prompt was crafted to provide the PLM with some options for a suitable label while still leaving it with the possibility of generating something new.

## 4 Datasets

We intentionally select a group of datasets with different features to analyse the correlation between features and optimal configurations. SNIPS (Coucke et al., 2018), AskUbuntuCorpus and WebApplications Corpus (Braun et al., 2017) all contain the Action-Object format of intents and are queries/commands in conversational style. DBpedia14 Sampled and StackOverflow (Xu et al., 2015) are labelled for Topic. DBpedia14\_Sampled

Dataset	Intent Type	Number of Samples	Number of Intents	Intent Balance	Average Number of Words	Vocabulary Size
AskUbuntu	Action-Object	Small (162)	Small (5)	Imbalanced (7.13)	Short (7.94)	Small (474)
SNIPS	Action-Object	Large (13784)	Small (7)	Slightly Imbalanced (1.03)	Short (9.15)	Large (13418)
WebApplications	Action-Object	Small (89)	Small (8)	Imbalanced (23.00)	Short (8.01)	Small (300)
Banking77	Mixed	Large (13083)	Large (77)	Imbalanced (3.03)	Short (11.71)	Medium (3027)
ChatbotCorpus	Mixed	Small (206)	Small (2)	Slightly Imbalanced (1.64)	Short (7.70)	Small (173)
CLINC	Mixed	Large (22500)	Large (150)	Balanced (1.00)	Short (8.31)	Medium (6420)
PersonalAssistant	Mixed	Large (20735)	Medium (46)	Imbalanced (247.96)	Short (6.84)	Medium (7896)
DBPedia14 Sampled	Topic	Large (14000)	Medium (14)	Balanced (1.00)	Long (46.29)	XLarge (75214)
StackOverflow	Topic	Large (20000)	Medium (20)	Balanced (1.00)	Short (8.32)	Large (16773)

Table 2: Features of Each Dataset

Feature	Categories
Intent Type	Action-Object, Topic, Mixed
Size	Small (<250), Large (>= 250)
Number of Intents	Small (<10), Medium (>=10, <50), Large (>=50)
Intent Balance	Balanced (IR = 1.00), Slightly Imbalanced (IR >1, <2), Imbalanced (IR >= 2)
Average Number of Words	Short (<20), Long (>=20)
Vocabulary Size	Small (<500), Medium (>=500, <10,000), Large (>=10,000, <50,000), XLarge (>=50,000)

Table 3: Categorisations of Dataset Features

contains a sample of 14,000 entries from the DBPedia14 dataset (Lehmann et al., 2014). Banking77 (Casanueva et al., 2020), ChatbotCorpus (Braun et al., 2017), CLINC (Larson et al., 2019) and PersonalAssistant (Liu et al., 2019) contain a mix of both Action-Object and Topic form of intents. See Table 2 for full details of the features of each dataset.

#### 4.1 Dataset Feature Definitions

We categorise the datasets by intent type, size, number of intents, whether the intents are balanced, average number of words and vocabulary size (see Table 3).

**Intent Type** Many works differ in their definition of intent, whether explicitly in their method or implicitly in their choice of dataset. Liu et al. (2021) define an intent as an Action(verb)-Object(noun) pair in an utterance e.g. “can you reschedule my delivery” has the pair ‘reschedule-delivery’. Vedula et al. (2020) also use this definition, naming these ‘actionable intents’. Other datasets have more abstract labels that are closer to topics. In these cases, methods like Action-Object extraction are unlikely to produce intents which reflect the ground-truths and so another extraction method would likely produce better results. Finally, a dataset can be mixed

such that it contains both Action-Object pairs and abstract labels like topics. Therefore, we categorise all datasets used in our experiments as one of Action-Object, Topic or Mixed.

**Number of Samples** We use a selection of datasets of varying sizes. The smallest dataset having less than 100 samples, while the largest has almost 22.5k. We categorise the datasets as either small or large where small is defined as having less than 250 samples and large has anything over 250.

**Number of Intents** The number of ground-truth intent labels in a dataset can be considered the ‘ideal’ number of clusters that should be found by the clustering algorithm. The datasets we use range from 2 to 150 intents and we categorise this feature as small, medium and large where small is defined as having less than 10 intents, medium has between 10 and 50 and anything over 50 is large.

**Intent Balance** The ground-truth label distribution is also a defining feature of datasets. We use the Imbalance Ratio (IR) as a measure of imbalance. This is simply the number of majority label samples over the number of minority label samples. An IR of 1.00 represents a completely balanced dataset with equal samples for every ground-truth label. Anything above this represents an increasing magnitude of imbalance. The datasets used range from balanced to an IR of 247.96 (the majority label has almost 250 times the samples of the minority label). We categorise this feature as balanced, slightly imbalanced and imbalanced where balanced has an IR of 1.00, slightly imbalanced has IR greater than 1 but less than 2 and imbalanced has an IR of 2 and above.

**Average Number of Words** The majority of the datasets used are dialogue utterances and have relatively low average number of words of less than 12 while only one exceeds this at 46.29. We therefore categorise this feature as short and long where

short is less than 20 and long is 20 and over.

**Vocabulary Size** The final dataset feature we explore is the number of unique words across all utterances in the dataset i.e. the vocabulary size. There is quite a spread across the datasets we use in our experiments and so we categorise this as small (with less than 500), medium (from 500 to 10,000), large (10,000 to 50,000) and xlarge (over 50,000).

## 5 Experiments

### 5.1 Experimental Setup

We evaluate all possible combinations of the choices in Table 1, with the only exceptions being for the previous OID works ITER\_DBSCAN and DeepAligned where we use the Semantic Representation model from the original works (Universal Sentence Encoder and BERT respectively). This results in 20 configurations for each dataset for the framework to execute.

Each configuration involves a clustering algorithm and clustering measure for conducting hyperparameter tuning. Clustering is attempted for a range of hyperparameter values and evaluated using the specified measure (we use silhouette score for our experiments). The hyperparameters with the best score according to the chosen clustering measure are used for the configuration. For kmeans, we must estimate the optimal number of clusters  $k$ . We therefore conduct clustering for  $k$  between 2 and 200 or the number of utterances in the dataset, whichever is lower. We use the scikit-learn implementation of kmeans. For DBSCAN, there are at least two parameters to be set.  $eps$  is the maximum distance that can be between two samples to consider them as being in the same neighbourhood and  $min\_samples$  is the minimum number of samples in a neighbourhood for a sample to be considered a ‘core’ sample. To keep hyperparameter tuning compute time down, we focus on tuning  $eps$  only, while  $min\_samples$  is set to 5. We cluster for  $eps$  between 0.1 and 1.0 with increments of 0.01. Again, we use the scikit-learn implementation for DBSCAN. For ITER-DBSCAN, there are five hyperparameters to be tuned. In addition to  $eps$  and  $min\_samples$ , there is also the change in these value for each iteration,  $delta\_eps$  and  $delta\_min\_samples$  and finally, the maximum number of iterations to run  $max\_iteration$ . An exhaustive search across these hyperparameters for every ITER-DBSCAN configuration and every dataset would be unfeasible.

We therefore generate 20 random sets of hyperparameters and cluster with these for every relevant configuration. We use the implementation of ITER-DBSCAN from the original work (Chatterjee and Sengupta, 2020). For the semi-supervised technique, DeepAligned, we use the implementation provided by the authors with their default values (Zhang et al., 2021b).

For configurations involving PLM prompting, we chose T0pp as it is open-source, small enough to deploy on accessible hardware and has produced impressive results (Sanh et al., 2021). We utilised AWS Sagemaker Notebook to run our experiments. A g4dn.12xlarge instance was used with any configuration with T0pp prompting and a g4dn.xlarge for the others.

### 5.2 Evaluation

We use two automated metrics (average cosine similarity and average BARTScore (Yuan et al., 2021)) to evaluate the quality of the final generated labels compared to the ground truth intents. Both the generated and ground truth label sets are normalised by converting to lower case, splitting on Pascal/snake case and removing hyphens and embeddings obtained using Universal Sentence Encoder.

For each unique ground-truth ( $gt$ ) label, we define  $C^*$  as the subset of clusters where the most common ground-truth ( $mcgt$ ) equals  $gt$ . The similarity score for each  $gt$  is then the average of the similarity between the generated label and the  $mcgt$  for each cluster in  $C^*$  ( $sim(c)$ ). If none of the identified clusters is assigned  $gt$  then the score is 0 (see Equation 1).

$$avg\_label\_sim(gt) = \begin{cases} \frac{\sum_{c \in C^*} sim(c)}{N_{C^*}} & , \text{ if } N_{C^*} > 0 \\ 0 & , \text{ if } N_{C^*} = 0 \end{cases} \quad (1)$$

where  $N_{C^*}$  is the number of clusters in  $C^*$ .

The final average similarity score for the configuration is calculated as in Equation 2.

$$config\_score = \frac{\sum_{gt \in GT} avg\_label\_sim(gt)}{N_{GT}} \quad (2)$$

where  $GT$  is the set of all ground-truth intents and  $N_{GT}$  is the number of ground-truth intents.

The optimal configuration for each dataset is the configuration which produces the highest  $config\_score$ . Collecting these results from

Dataset	Semantic Representation	Clustering Algorithm	Candidate Extraction	Label Selection	No. Clusters	Avg. Cosine Similarity	Avg. BART Score
AskUbuntu	use	KMeans	Action-Object	T0pp Prompting	6(+1)	0.4661	-5.7580
SNIPS	use	KMeans	Action-Object	T0pp Prompting	8(+1)	0.6163	-3.9832
WebApplications	all-mpnet	KMeans	Action-Object	T0pp Prompting	6(-2)	0.4993	-5.4204
Banking77	all-mpnet	KMeans	Action-Object	Most Frequent	196(+119)	0.4678	-5.4880
ChatbotCorpus	use	KMeans	T0pp Prompting	Most Frequent	4(+2)	0.4384	-4.9715
CLINC	use/ all-mpnet	KMeans	Action-Object	T0pp Prompting/ Most Frequent	163(+13)/ 155(+5)	0.5050/ 0.5044	-4.7101/ -4.5701
PersonalAssistant	all-mpnet	KMeans	Action-Object	T0pp Prompting	60(+14)	0.3843	-5.2462
DBPedia14 Sampled	use/ all-mpnet	KMeans	T0pp Prompting	Most Frequent	11(-3)/ 10(-4)	0.3378/ 0.3091	-5.3313/ -5.3169
StackOverflow	use/ all-mpnet	KMeans	Action-Object	T0pp Prompting	23(+3) / 21(+1)	0.4861/ 0.3922	-5.2722/ -5.2692

Table 4: Unsupervised configurations producing the optimal labels for each dataset. The difference in number of clusters and ground-truth intents in shown in brackets. Where the evaluation metrics disagree on a configuration choice, both are reported as (*cosine similarity score/BART score*)

datasets of different features allows us to analyse the optimal configurations alongside the features in order to infer any dependencies between them.

## 6 Results and Analysis

### 6.1 Unsupervised Clustering

Table 4 shows the optimal configurations together with the average scores that they achieved in the unsupervised clustering setting. Table 5 shows a sample of the final labels generated with unsupervised clustering for each dataset. Many of these labels are of high quality and would be useful in downstream systems. In all unsupervised settings, KMeans produced the clusters for the optimal configuration. In most cases, the number of clusters exceeded the number of ground-truth intents. This results in some clusters being assigned the same *mcgt*. The labels are however, highly semantically similar with their ground-truth counterparts. It appears that the configurations using ITER\_DBSCAN have produced a great overestimation of the number of clusters e.g. for SNIPS, the best performing configuration using ITER\_DBSCAN produced 39 clusters. The generated labels are still semantically similar to their ground-truths, however there is more variety per ground-truth label due to the finer-grained clusters generating different final labels, resulting in lower performance according to the evaluation metrics.

Where Action-Object candidate extraction was used it has resulted in some generated labels being less descriptive than would perhaps be desired, e.g. in SNIPS *find-schedule* for **SearchScreeningEvent** is too generic. The samples for this intent are look-

Ground Truth	Generated Label
<b>SNIPS</b>	
AddToPlaylist	add-song
BookRestaurant	book-restaurant
GetWeather	give-forecast
PlayMusic	play-music/find-soundtrack
RateBook	rate-novel
SearchCreativeWork	find-show
SearchScreeningEvent	find-schedule
<b>Banking77</b>	
card_arrival	received-card/track-card
edit_personal_details	edit-details/?change-address.
exchange_charge	exchange-currencies/exchanging-currencies?
getting_virtual_card	get-card?
passcode_forgotten	reset-password/?reset-passcode?
request_refund	get-refund/give-refund
verify_my_identity	verify-identify?
verify_source_of_funds	get-funds/verify-source
<b>CLINC</b>	
how_old_are_you	ask-age/tell-birthday
improve_credit_score	improve-score
oil_change_when	change-oil
plug_type	need-converter
schedule_meeting	reserve-room/set-meeting
text	tell-text
transactions	show-transactions
who_do_you_work_for	tell-brand
<b>StackOverflow</b>	
apache	redirect-requests/using-proxy
cocoa	Cocoa/converting-string
hibernate	Hibernate
linq	using-linq
qt	Qt: How to end line with QTextEdit [Qt] [C++],
spring	Spring
wordpress	get-posts
visual-studio	Visual Studio 2008

Table 5: Sample labels produced by the optimal configurations. Where multiple clusters are assigned the same *mcgt*, we report two sample generated labels.

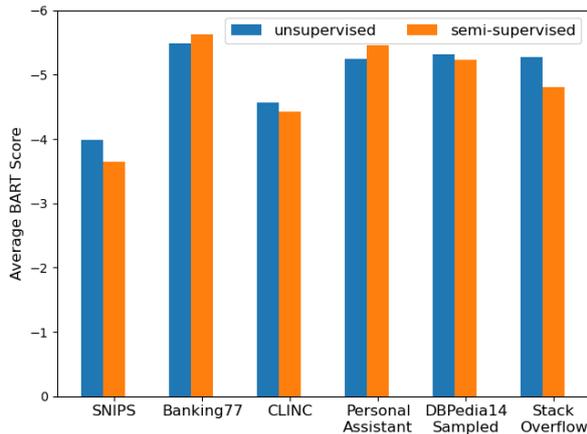


Figure 3: Average BART Scores for the optimal unsupervised configs vs optimal semi-supervised configs for each dataset. Closer to zero is better.

ing for the movie schedules at cinemas and often ask for “the movie schedule”. Also, there are many fine-grained intents in Banking77 which require more detail to be immediately useful e.g. a ground-truth intent such as **get\_disposable\_virtual\_card** could not be produced using Action-Object extraction as in (Liu et al., 2021). It would therefore, be useful to extend the Action-Object candidate extraction to include compound nouns and adjectives to capture further details in the candidates.

## 6.2 Semi-supervised Clustering

Figure 3 shows the difference in BART Score for the optimal configurations using unsupervised clustering vs the optimal config that used the semi-supervised clustering method DeepAligned. The quality of the generated labels mostly outperform their unsupervised counterparts. However, DeepAligned produces poorer results for both Banking77 and Personal Assistant. These datasets are both large in size and imbalanced which may cause the DeepAligned model to overfit to the majority samples. DeepAligned also failed to complete for the small datasets, possibly due to a lack of training samples to complete an optimizer step.

## 6.3 Mapping Features to Configuration

Table 6 shows how the various dataset features affect the optimal unsupervised configuration when evaluating using the BART Score. Each value represents the most commonly used option for a given dataset feature and step in the framework, e.g. for datasets with Action-Object as the target intent type, Universal Sentence Encoder was the majority optimal choice for Semantic Representa-

Feature	Semantic Representation	Clustering Algorithm	Extraction Method	Selection Method
<b>Intent Type</b>				
Action-Object	use	KMeans	Action-Object	T0pp Prompting
Topic	all-mpnet	KMeans	No Majority	No Majority
Mixed	all-mpnet	KMeans	Action-Object	Most Frequent
<b>Size</b>				
Small	use	KMeans	Action-Object	T0pp Prompting
Large	all-mpnet	KMeans	Action-Object	No Majority
<b>Num. Intents</b>				
Small	use	KMeans	Action-Object	T0pp Prompting
Medium	all-mpnet	KMeans	Action-Object	T0pp Prompting
Large	all-mpnet	KMeans	Action-Object	Most Frequent
<b>Imbalance</b>				
Balanced	all-mpnet	KMeans	Action-Object	Most Frequent
Slightly Imbalanced	use	KMeans	No Majority	No Majority
Imbalanced	all-mpnet	KMeans	Action-Object	T0pp Prompting
<b>Avg. Num. Words</b>				
Short	all-mpnet	KMeans	Action-Object	T0pp Prompting
Long	all-mpnet	KMeans	T0pp Prompting	Most Frequent
<b>Vocab. Size</b>				
Small	use	KMeans	Action-Object	T0pp Prompting
Medium	all-mpnet	KMeans	Action-Object	Most Frequent
Large	No Majority	KMeans	Action-Object	T0pp Prompting
XLarge	all-mpnet	KMeans	T0pp Prompting	Most Frequent

Table 6: Most common options by dataset features when evaluating using BART Score

tion. This table can act as an aid in the choice of config for a new, unlabelled dataset. For example, if we consider CLINC to be our unlabelled set, we could choose our configuration from this table rather than at random (to make this a fair example, we remove CLINC’s results from the table). With little domain knowledge, we can infer that the CLINC utterances contain Mixed intents (both Action-Object and Topics) and estimate that there are a Large number of intents (more than 50). A clustering algorithm could be used to estimate the IR, showing that it is a Balanced set. The dataset size is Large, containing 22,500 utterances which are made up of Short sentences of less than 20 words with a total vocabulary size of 6420 words (Medium). For these features, the table agrees on all-mpnet, KMeans and Action-Object on every feature. There is a disagreement on the Selection Method and so we choose Most Frequent as it is less compute intensive. As shown in Table 4, this is the optimal configuration for CLINC when evaluating on BART Score. Were we to naively choose T0pp Prompting for both Candidate Extraction and Label Selection, in the belief that a more flexible approach would be best, the final labels produced would be of lower quality overall (average BART of -5.0281 compared to -4.5701). Many of the labels generated by this configuration are simply ‘ask a question’ or in one case ‘Yes’ for a cluster with *mcgt* **ingredient\_substitution**. Such issues could be overcome with further prompt tuning, however we can already obtain high quality labels from sim-

pler, less hardware and time expensive methods.

## 7 Conclusions and Future Work

We have shown that our framework for OID can produce high quality labels for many datasets of differing intent type. The modular nature of the framework allows for further improvements to be utilised when new techniques are discovered for each step. We have evaluated a number of configurations based on the final generated label quality, including extending previous OID works which originally do not generate a human-readable intent label. We have also presented an initial analysis of the mapping between dataset features and the optimal configuration to use for a new, unlabelled dataset which can help reduce the initial effort required to choose the combination of techniques. In future work, we plan to add our Action-Object Extension technique (proposed in [Anderson et al. \(2024\)](#)) to the framework and update the optimal configuration results. We also hope to curate more intent datasets of varying features in order to develop a model for predicting a ‘best guess’ configuration, given a new dataset’s features, rather than having to try every one in turn.

## Limitations

Our work is limited to the set of techniques chosen for each step in the framework. There exists many other appropriate semantic representation models, clustering algorithms, candidate extraction and selection methods which could possibly produce higher quality labels. Also, the evaluation of the intent labels is based on semantic similarity to the ground-truth labels. This has the implicit assumption that the ground-truth labels are the best representation for the intent which may not necessarily be the case.

## Acknowledgments

Our thanks go to our internal reviewers Xinyu Chen and Cynthia Freeman, for their helpful feedback on the first draft. We also thank all anonymous SIGDIAL reviewers for their comments and constructive suggestions.

## References

Grant Anderson, Emma Hart, Dimitra Gkatzia, and Ian Beaver. 2024. [Automated Human-Readable Label Generation in Open Intent Discovery](#). In *Proc. INTERSPEECH 2024*, page tbc.

Daniel Braun, Adrian Hernandez-Mendez, Florian Matthes, and Manfred Langen. 2017. [Evaluating natural language understanding services for conversational question answering systems](#). In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 174–185, Saarbrücken, Germany. Association for Computational Linguistics.

Iñigo Casanueva, Tadas Temcinas, Daniela Gerz, Matthew Henderson, and Ivan Vulic. 2020. [Efficient intent detection with dual sentence encoders](#). *CoRR*, abs/2003.04807.

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. [Universal sentence encoder](#). *CoRR*, abs/1803.11175.

Ajay Chatterjee and Shubhashis Sengupta. 2020. [Intent mining from past conversations for conversational agent](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4140–4152, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Zhiyuan Chen, Bing Liu, Meichun Hsu, Malu Castellanos, and Riddhiman Ghosh. 2013. [Identifying intention posts in discussion forums](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1041–1050, Atlanta, Georgia. Association for Computational Linguistics.

Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, Maël Primet, and Joseph Dureau. 2018. [Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces](#). *CoRR*, abs/1805.10190.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.

Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, KDD’96, page 226–231. AAAI Press.

Chih-Wen Goo, Guang Gao, Yun-Kai Hsu, Chih-Li Huo, Tsung-Chieh Chen, Keng-Wei Hsu, and Yun-Nung Chen. 2018. [Slot-gated modeling for joint slot filling and intent prediction](#). 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 753–757, New Orleans, Louisiana. Association for Computational Linguistics.

- Joo-Kyung Kim, Gokhan Tur, Asli Celikyilmaz, Bin Cao, and Ye-Yi Wang. 2016. [Intent detection using semantically enriched word embeddings](#). In *2016 IEEE Spoken Language Technology Workshop (SLT)*, pages 414–419.
- Rajat Kumar, Mayur Patidar, Vaibhav Varshney, Lovekesh Vig, and Gautam Shroff. 2022. [Intent detection and discovery from user logs via deep semi-supervised contrastive clustering](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1836–1853, Seattle, United States. Association for Computational Linguistics.
- Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. [An evaluation dataset for intent classification and out-of-scope prediction](#). 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 1311–1316, Hong Kong, China. Association for Computational Linguistics.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, and Christian Bizer. 2014. [Dbpedia - a large-scale, multilingual knowledge base extracted from wikipedia](#). *Semantic Web Journal*, 6.
- Ting-En Lin, Hua Xu, and Hanlei Zhang. 2019. [Discovering new intents via constrained deep adaptive clustering with cluster refinement](#). *CoRR*, abs/1911.08891.
- Bing Liu and Ian Lane. 2016. [Attention-based recurrent neural network models for joint intent detection and slot filling](#).
- Pengfei Liu, Youzhang Ning, King Keung Wu, Kun Li, and Helen Meng. 2021. [Open intent discovery through unsupervised semantic clustering and dependency parsing](#). *CoRR*, abs/2104.12114.
- Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2019. [Benchmarking natural language understanding services for building conversational agents](#). *CoRR*, abs/1903.05566.
- Hugh Perkins and Yi Yang. 2019. [Dialog intent induction with deep multi-view clustering](#). 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 4016–4025, Hong Kong, China. 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.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, 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, Stella Biderman, Leo Gao, Tali Bers, Thomas Wolf, and Alexander M. Rush. 2021. [Multi-task prompted training enables zero-shot task generalization](#).
- X.Y. Shen, Ying Sun, Yao zhong Zhang, and Mani Najmabadi. 2021. [Semi-supervised intent discovery with contrastive learning](#). *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*.
- Nikhita Vedula, Nedim Lipka, Pranav Maneriker, and Srinivasan Parthasarathy. 2020. [Open intent extraction from natural language interactions](#). In *Proceedings of The Web Conference 2020, WWW '20*, page 2009–2020, New York, NY, USA. Association for Computing Machinery.
- Jiaming Xu, Peng Wang, Guanhua Tian, Bo Xu, Jun Zhao, Fangyuan Wang, and Hongwei Hao. 2015. [Short text clustering via convolutional neural networks](#). In *Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing*, pages 62–69, Denver, Colorado. Association for Computational Linguistics.
- Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. [Bartscore: Evaluating generated text as text generation](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 27263–27277. Curran Associates, Inc.
- Hanlei Zhang, Xiaoteng Li, Hua Xu, Panpan Zhang, Kang Zhao, and Kai Gao. 2021a. [TEXT0IR: An integrated and visualized platform for text open intent recognition](#). 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 167–174, Online. Association for Computational Linguistics.
- Hanlei Zhang, Hua Xu, Ting-En Lin, and Rui Lv. 2021b. [Discovering new intents with deep aligned clustering](#). In *AAAI*.
- Junmei Zhong and William Li. 2019. [Predicting customer call intent by analyzing phone call transcripts based on cnn for multi-class classification](#).

# Toximatics: Towards Understanding Toxicity in Real-Life Social Situations

Mayukh Das and Wolf-Tilo Balke

Institute for Information Systems, TU Braunschweig  
Mühlenpfordtstraße 23, 38106 Braunschweig, Germany  
{mayukh,balke}@ifis.cs.tu-bs.de

## Abstract

The rise of social media has amplified the visibility and impact of hate speech, prompting the development of NLP solutions to identify both explicit and implicit forms of hate speech. These approaches assess toxicity in isolation, neglecting context and limiting models to sentence-level understanding. Therefore we study, how contextual factors influence perceived toxicity, thereby anchoring assessments in a more nuanced semantic framework. We introduce a novel synthetic data generation pipeline designed to create context-utterance pairs at scale with controlled polarity. This pipeline can enhance existing hate speech datasets by adding contextual information to utterances, either preserving or altering their polarity, and also generate completely new pairs from seed statements. We utilised both features to create *Toximatics*, a dataset that includes *context-dependent utterances* and its toxicity score. To address biases in state-of-the-art hate datasets, which often skew towards specific sensitive topics such as politics, race, and gender, we propose a method to generate neutral utterances typical of various social settings. These are then contextualized to show how neutrality can shift to toxicity or benignity depending on the surrounding context. *Toximatics*' approach to hate speech detection extends beyond the sentence level, rendering it suitable for discourse analysis and also revealing that current models underperform on this dataset.

## 1 Introduction

Toxicity classifiers are normally fine-tuned with hate speech datasets that contain explicit or overtly abusive lexicons (Davidson et al., 2017; Founta et al., 2018) or implicit, coded, indirect framing of offensiveness (ElSherief et al., 2021; Hartvigsen et al., 2022). Explicit hate datasets suffer from topic bias like over-reliance on sensitive attributes (race, gender, religion, nationality, etc) (Basile

et al., 2019) which can inflate model performance on phrases containing indirect offense. Implicit hate speech introduces diverse hate classes based on coded language such as irony, sarcasm, euphemism, metaphor, circumlocution, etc (Talat and Hovy, 2016; Magu and Luo, 2018; Gao and Huang, 2017; Warner and Hirschberg, 2012; Qian et al., 2019). This line of work approaches the problem of detecting toxicity/hate speech as: Given a sequence of  $n$  tokens  $u \in \{x_1, x_2 \dots x_{target} \dots x_n\}$  with  $x_{target}$  being a sensitive attribute or target group type, can a machine flag whether  $u$  is *toxic*? (Sap et al., 2020) defines the task differently, like given a sequence  $u$ , can a machine determine if a possible toxicity is *intentional*, *offensive*, *lewd* towards any identity group. The sequence  $u$  could be interpreted as an utterance by a person in online or real-life social situations. In NLP datasets, these sequences are usually mined via social media or machine-generated with specialized prompts (Hartvigsen et al., 2022).

However, does toxicity depend solely on the semantics of the utterance  $u$ ? Our hypothesis is that this current method is simplistic as it bases the concept of toxicity solely on the utterance. To accurately determine toxicity, it is necessary to take into account the broader context. This paper illustrates how the context of spoken dialogue can alter the human toxicity score of the same utterances (see Figure 1). Consequently, we also highlight other flawed assumptions made by the previous research. Previous research on hate speech has highlighted the importance of considering characteristics such as abusive language, speaker intention, sensitive attributes, target group, and framing when determining what constitutes hate speech. However, we speculate human perceived toxicity can occur for any social category if grounded in context, e.g. *{Irony towards an introvert in a social gathering}*. Recent implicit datasets (Sap et al., 2020; Hartvigsen et al., 2022; Zhou et al.,

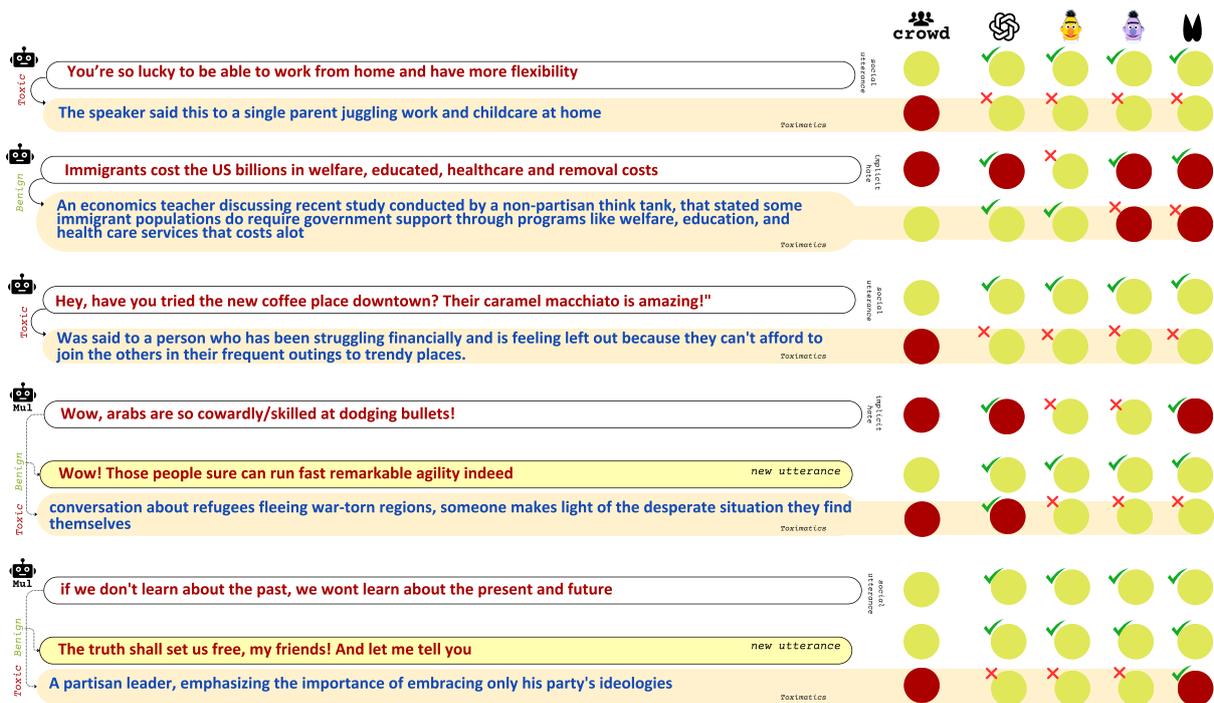


Figure 1: Toximatics dataset with its parent seed utterances. Direct augmentation and multistage augmentation are generation methods introduced in section 4. Labels are chat-gpt-4 legacy, hate-bert, roberta-toxigen, llama-2-chat-7b, ozawa human annotation. Toximatics makes the polarity of the utterance context-dependent and fools the models.

2023) overemphasize dimensions like intention, power dynamics, or target groups as hate determiners. On the other hand, we highlight that toxic perlocution can occur regardless of power dynamic, identity group. Furthermore, we observed distinct examples for such cases are clearly missing in state-of-the-art papers and datasets.

To address all this gap we introduce *Toximatics*. A dataset to understand pragmatic toxic utterance which encompasses deeper level semantics than implicit datasets. We ground the notion of toxicity to the context of the utterance, rather than grounding it solely to the utterance as done in previous work. Unlike previous work, our samples consist of an utterance-context pair. The crowdworkers were presented with the context while annotating the utterance to make sure the validity of the toxicity score becomes depended on the context. We also introduce a generation pipeline that utilizes state-of-the-art language models and expertly crafted prompts and methods. This enables to generate similar examples at scale without the need for further quality checks, ensuring the high level of accuracy. Our generation pipeline uses utterances from implicit datasets as seeds to create context with controlled polarity and also gen-

erates entirely new utterance-context pairs from these seeds. Additionally, we produce seemingly neutral utterances, atypical for certain social situations, and generate contexts for them with specific polarity control. Following the later step we explore whether phrases like {*You're so lucky to be able to work from home and have more flexibility*} could be perceived as toxic without any allusion to power dynamics, intent, or identity groups (see Figure 1). We release a dataset of 19,800 utterance-context pairs and their toxicity labels. Toximatics was evaluated with sota toxicity classification models, foundational models and chat models, which showed poor efficacy on zero-shot classification task. This dataset is the first of its scale to pivot hate speech detection research towards a context-dependent framework. The dataset and generation method codes are available via the provided link <sup>1</sup>.

## 2 Related Work

Early papers on hate datasets emphasized explicit abusive language and profane use of slurs racial identifiers, minority mentions, hateful keywords, etc (Basile et al., 2019; Davidson et al., 2017; Warner and Hirschberg, 2012; Silva et al., 2016;

<sup>1</sup><https://github.com/Mayukhga83/Toximatics>

Burnap and Williams, 2014; de Gibert et al., 2018). These examples are collected using keyword-based, bootstrap scrapping, or adversarial data collection (Davidson et al., 2017; Zampieri et al., 2019; Founta et al., 2018; Waseem, 2016; Dinan et al., 2019; Vidgen et al., 2021). These datasets have an over-reliance on lexical cues and specific topics. In response to this, researchers have tried to curate a newer corpus that labels hate considering the rhetorical framing grounded in sociology and psychology (Kennedy et al., 2018; Sap et al., 2019). ElSherief et al. (2021) introduced a taxonomy of implicit hate speech and a benchmark corpus mined from online hate groups. Hartvigsen et al. (2022) uses GPT 3, demonstration-based prompting and constrained decoding to generate large-scale implicit hate corpus. Pavlopoulos et al. (2020) investigated the potential effect of context on human judgment of toxicity scores through an analysis of Wikipedia discussions. Xenos et al. (2021) created a toxicity dataset where the annotators had access to one previous comment. Zhou et al. (2023) developed a formalism to explain the intentions, reactions, and harms of offensive or biased statements based on their social context.

Unlike most previous works which focuses on a single statement, we have a situational context in which the statement was uttered. While the previous work examines the extent to which the framing of an utterance determines its level of toxicity, our work investigates the extent to which the context determines the toxicity level of an utterance. Previous studies (Pavlopoulos et al., 2020; Xenos et al., 2021) narrowly examined context, focusing solely on preceding comments and discussion headings. However, this limited approach may fail to capture the circumstances of the utterance. Thus, we propose contextualizing the situation with a situation descriptor (see Figure 1). Zhou et al. (2023) heavily relies on the identity group of both the speaker and listener. Additionally, the context description is limited in scope. In contrast, we solely use detailed situational descriptor as context as an explanation of the entire scenario. We also have curated examples to show how toxicity can be perceived without any allusion to identity group. Zhou et al. (2023) also does not generate large scale polarity controlled context, they only have 928 counterfactual context. In contrast, our work solely deals with generating polarity-controlled context. None of the previous works have focused on gener-

ating completely new implicit hate utterances in a given context, nor have they attempted to uncover the toxic nature of arbitrary social statements in a contextualized manner, unlike us.

### 3 Pragmatics, Meaning, and Toxicity

To adapt toxicity detection (Founta et al., 2018) to dialogs, we formalise toxicity as something that can potentially affect the climate of discourse in a negative way. In technical terms, if we have a hypothetical value function  $V(\mathcal{D}/C)$  that can estimate the state-value of discourse  $\mathcal{D}$  at a specific time given context  $C$ , an utterance  $u_t$  at time  $t$  is a potential contestant for hate speech if

$$V(\mathcal{D}_{t+i}/C_{<t+i}) \ll V(\mathcal{D}_{t-1}/C_{<t-1}) : i \geq 1$$

Empirically  $V(\cdot)$  is impossible to estimate due to the subjective perception of language among humans and lack of consensus on what to include in  $C_{<t}$ . Defining hate speech in this way highlights the limitations of basing toxicity levels solely on snippets of utterances. In this paper, we consider free text situational descriptor as  $C_{<t}$ .

We hypothesise that toxicity is performative. In linguistics, performatives are speech acts that not only convey information but also perform an action and have a perlocutionary effect on the listener’s mind (Austin, 1962). For instance, *"I would like some Kimchi!"* at a dinner table implies *"pass me the Kimchi"*. Perlocutionary effects include persuading, convincing, enlightening, and commanding. We propose that conveying hate or offense is a valid perlocution, potentially affecting perceived toxicity scores when annotators have full context.<sup>2</sup> We aim to investigate how perceived toxicity changes across different contexts and nuanced situations.

### 4 Generation Pipeline

In this section we formalise a general overview of the pipeline, a straightforward summary of which is presented in Figure 2. The pipeline utilizes supervised finetuned language (SFT) models (Ouyang et al., 2022), contrastive search decoder (Li et al., 2022) and carefully curated prompts as the base elements. After conducting several preliminary experiments, we propose the prompt should have a template designed to achieve prespecified goal as

<sup>2</sup>Please note that hate or offense was never tied to perlocution by Austin (1962), this is one contribution of our dataset

done in prompt engineering (Sahoo et al., 2024) and it should also contain few in-context examples. Following our findings (appendix A), we propose using contrastive search (over top-p or temperature) because it along with in-context prompts reduces hallucination and improve the quality of generations while maintaining relevance to the instruction. These claims are supported by findings from other sources (O’Brien and Lewis, 2023). The pipeline supports three types of context augmentation, depending on the number of iterations and the dynamic addition of statements. This is controlled by the target polarity and other hyperparameters.

#### 4.1 Direct Augment

Let  $\mathcal{L}_\theta^{(\alpha, \kappa, H)}$  be a pretrained language model parameterised by  $\theta$  coupled with contrastive search decoder parameterised by  $\alpha$  and  $\kappa$  and set  $H \in (h_1, h_2, \dots)$  containing hyperparameters that modifies the output logits.  $H$  includes properties like *repetition-penalty*, *max-token*, *repeat-ngram*, etc.  $\alpha$  and  $\kappa$  controls the trade-off between model confidence and degeneration penalty. Formally given the input prompt  $x_{<t}$  the selection of output  $x_t$  will follow:

$$x_t = \arg \max_{v \in \mathcal{V}^{(\kappa)}} \{(1 - \alpha)p_\theta(v|x_{<t}) - \alpha(\max_{1 \leq j \leq t-1} \{s(v, x_j)\})\}$$

Where  $\mathcal{V}^{(\kappa)}$  is the *top-k* prediction from the LMs probability distribution  $p_\theta(\cdot|x_{<t})$ . Model confidence, is the probability of candidate  $v$  predicted by the LMs  $p_\theta(v|x_{<t})$ . Degeneration penalty  $\max\{s(v, x_j) : 1 \leq j \leq t - 1\}$ , measures the maximum cosine similarity between the candidate  $v$  and the tokens in the input prompt. In case of direct augment if  $u$  be any predefined utterance and  $t_p$  be the target polarity of the utterance then the context generated by direct augment is given by:

$$C = \mathcal{L}_\theta^{(\alpha, \kappa, H)}(u, P_{cont}(n, t_p))$$

Where  $P_{cont}(n, t_p)$  is the taylored prompt having  $n$  in-context examples and instruction to generate context given utterance  $u$ .

#### 4.2 Multistage Augment

This method generates completely new utterance-context pair by passing the input through LMs at multiple steps with distinct polarity objectives. Using three chains of target polarity adds dynamic to the connotation of the utterance-context pair

and its framing. For example, a seemingly neutral context  $u$  could first be made toxic along with a generated context  $C$ . Then a new utterance  $u_{new}$  could be constructed which along with the previous context sounds benign. Then again a new context  $C_{new}$  could be constructed which along with  $u_{new}$  sounds toxic. If  $P_{utt}(n, t_p)$  is the taylored prompt having  $n$  in-context examples and instruction to generate utterance given context  $C$  then the process can be written as:

$$C = \mathcal{L}_{\theta_1}^{(\alpha_1, \kappa_1, H_1)}(u, P_{cont}(n, t_{p_1}))$$

$$u_{new} = \mathcal{L}_{\theta_2}^{(\alpha_2, \kappa_2, H_2)}(C, P_{utt}(n, t_{p_2}))$$

$$C_{new} = \mathcal{L}_{\theta_3}^{(\alpha_3, \kappa_3, H_3)}(u_{new}, P_{cont}(n, t_{p_3}))$$

Where  $t_{p_i}$  is the target polarity at  $i$ th step. The dynamic nature of this method improves the quality of counterfactual examples greatly.

#### 4.3 N-iter Multistage Augment

This methods further extends multistage augment with new utterance at  $N$  intermediate steps (typically  $N = 2, 3, 4, \dots$ ). This further adds dynamic to the utterance and context quality and helps even improve the counterfactual examples. The steps could be written as follows

$$(u_1, C_1) = M_\Theta(u, P)$$

$$\forall i \in (2, 3, 4, \dots, N - 1)$$

$$u_i = \mathcal{L}_{\theta_2}^{(\alpha_2, \kappa_2, H_2)}(C_{i-1}, P_{utt}(n, t_{p_2}))$$

$$C_i = \mathcal{L}_{\theta_3}^{(\alpha_3, \kappa_3, H_3)}(u_i, P_{cont}(n, t_{p_3}))$$

Where  $M_\Theta$  is the multistage augment step with  $\Theta$  containing all the hyperparameters associated with that step.  $u_i$  and  $C_i$  being the generated utterance and context at  $i$ th step.

## 5 Dataset Generation

All augmentation methods were utilised in the pipeline for creation of Toximatics.

### 5.1 Models

We utilized the largest available open-source model 70 billion parameter LLama 2 chat model (Touvron et al., 2023) supervised finetuned with Orca dataset (Mittra et al., 2023). We conducted a side experiment to compare different SFT versions of the model for our task. We generated 5 generations using the direct augmentation method and crowd-validated

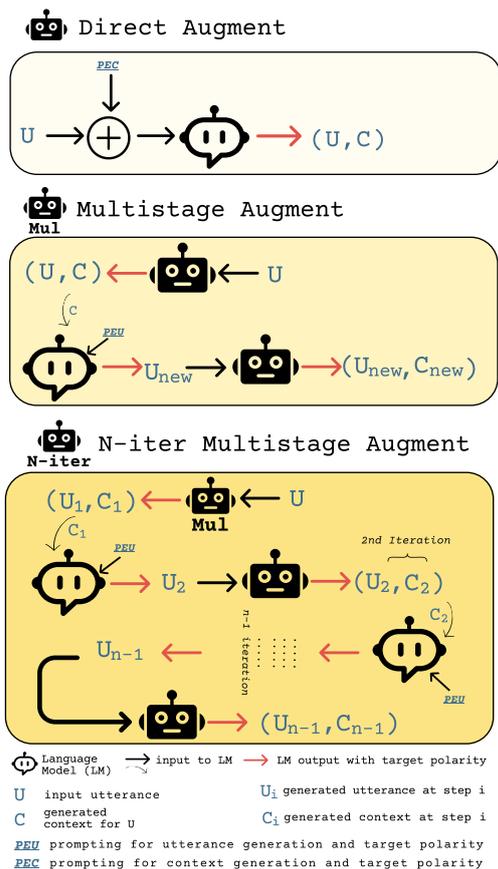


Figure 2: The generation pipeline supports 3 methods: Direct Augment adds context to the seed utterance, while Multistage and N-iter Multistage can generate novel utterance-context pairs given the seed.

the quality as the relevance of the generation to the prompt. Results reveal that the Oraca finetuned version outperformed the others with an agreement of 0.71. To streamline the process, we substituted the same model at all stages of the multistage augmentation.

quality	llama2 Orca	llama2 Oasst	llama2 chat	llama2 base
relevance	0.79	0.72	0.58	0.51

Table 1: Generation quality for various SFT versions of LLama-2, Oasst is OpenAssistance sft version while chat is meta’s llama2 sft version

## 5.2 Seed Utterance

We generate Toximatics from both state of the art implicit hate datasets and socially grounded neutral statements which were also curated with generative models.

### 5.2.1 Implicit Hate Dataset

We primarily used ToxiGen (Hartvigsen et al., 2022) which is a large-scale machine-generated

dataset containing human annotated toxicity score, framing and perceived intent. We first divide the dataset (train set) into three parts by thresholding over the human toxicity score. We taxonomize them as *benign-batch* ( $0.5 \leq h_{tox} \leq 1.5$ ) with 3230 samples, *neutral-batch* ( $1.5 \leq h_{tox} \leq 3.5$ ) with 3230 samples and *toxic-batch* ( $3.5 \leq h_{tox} \leq 4.5$ ) with 1145 samples, where  $h_{tox}$  is the human toxicity score in a scale of (1, 5). Extreme toxic statements were left out as they contain lexical cues of overt negative words. This distinction based on a threshold was established to enable the creation of experiments with precise goals, such as modifying the toxicity of samples to benign, toxic to neutral, neutral to benign and so on (see section 5.4).

### 5.2.2 Socially Grounded Neutral Statement

The primary objective of these seed utterances is to challenge preconceived notions of toxicity linked to power dynamics, identity groups, race, politics, and gender. Instead, we aim to ground the analysis in more generic contexts, such as whether an utterance in a restaurant, a birthday celebration, or a friendly environment can be perceived as toxic. This approach allows us to analyze the polarity of utterances within valid social contexts, termed “base-context”, as opposed to online comments. We mined the base-context as detailed below (see Figure 3).

**Conversational Topic Extraction:** First we apply a topic model algorithm based on BertTopic (Groo-tendorst, 2022) on two conversation data sets Daily-Dialog (Li et al., 2017) and Blended-Skill-Talk (Smith et al., 2020). Firstly, the dialogues were converted to embeddings using *Sentence Transformer* (Reimers and Gurevych, 2020), and then reduced in dimensionality using *UMAP* (McInnes et al., 2018) with key-parameters like nearest neighbour size as 15 and min-dist as 0.25 (the minimum distance between points in low-dimensional space). Setting both parameters to low helps to emphasize the local structure of conversational data. *HDBSCAN* (McInnes and Healy, 2017) was employed as the clustering algorithm, with Euclidean as the distant metric and minimum cluster size of 200 so that we don’t end up having too many clusters. The topic theme was generalized from the topic cluster keywords using chat-gpt-4, and it was then taken as the conversational topic. In this way, the two datasets yielded 413 conversation themes.

**Social Location Extraction:** We define a social location as any place that has a social environment

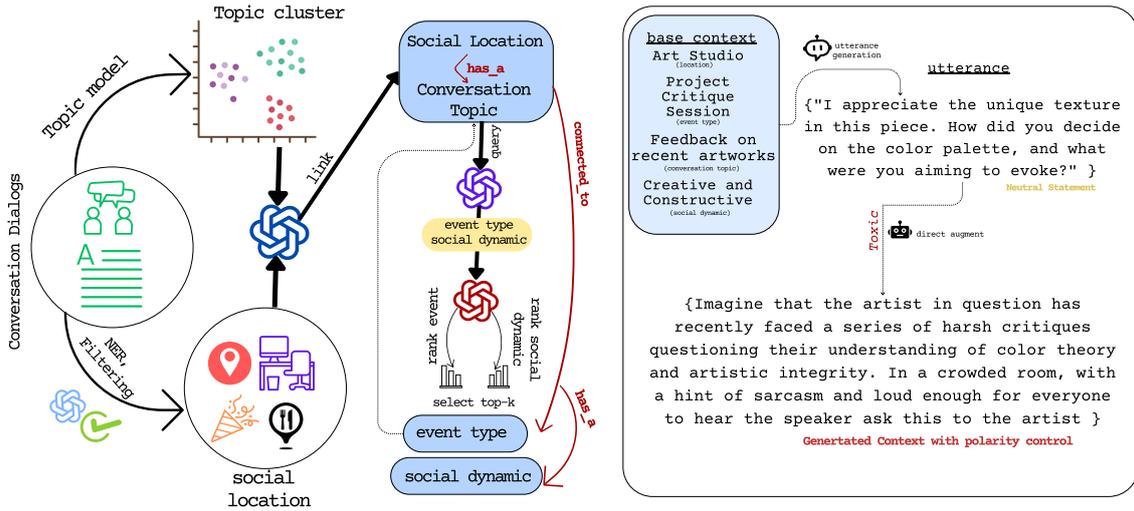


Figure 3: We first mine utterances atypical to a social topic and then augment it with polarity controlled context. We utilised LLMs to assist the mining process. Like chat-gpt-4 during linking phase, filter and augment social location, to generate candidate event types and social dynamics base gpt4 to rerank the generated list

and can stimulate civil conversation, such as restaurants, parks, bars. We employed Named Entity Recognition (NER) to the same dialog datasets with spacy (Honnibal and Montani, 2017) to extract *FAC* (facilities) location data; a total of 685 samples. Consequently, duplicates were removed and chat-gpt-4 was instructed to generalise several categories of social location from the remaining set e.g., *Entertainment and Recreation*  $\approx$  {*Disneyland, Saikei Ski Resort, Zoo, Kangaroo Club, etc*}. Taking these categories and its associated location as example we instruct chat-gpt-4 to mine a set of 150 social locations (appendix B).

**Linking:** This stage links the conversational topic to the social location via chat-gpt-4. Set  $S$  containing social location and set  $C_t$  having conversational topics has many-to-many relationship with overlapping association (e.g, almost any conversation can happen in a cafe but only some in hospital). The set of links  $L$  is therefore the subset of Cartesian product of  $S$  and  $C_t$ .

$$L \subseteq S \times C_t$$

We present each element  $c_i$  in  $C_t$  and entire set  $S$  to chat-gpt-4 and instruct to link  $c_i$  to elements of  $S$  with one in-context example (appendix C).

**Event and Social Dynamic Prediction:** As it is non-trivial to mine events and social dynamics from conversation datasets, we use LLMs as the a retrieval system. For each linked location and conversation  $l_i \in L$ , chat-gpt-4 was employed to generate a preliminary list of potential event types. This list was re-ranked with base GPT4 (appendix D & E).

The decision was influenced by (Sun et al., 2023), who demonstrated the efficacy of LLMs in retrieval tasks and identified that base GPT4 outperforms all other models in ranking tasks. We then applied a top-k threshold to select  $k$  entities from the list as a measure of most relevance. In this paper we used  $k = 3$  to account for high relevance and brevity of our dataset. After appending the event type to the base context, we repeat the same step for social dynamic.

The aforementioned procedure yielded 1554 base-context units, with approximately 2% of these removed by three crowd validators (appendix G) with an agreement of 0.88, resulting in 1523 units (examples in F). We generated 1523 seed utterance associated with the units with our generation pipeline and name this batch as *social-neutral-batch*.

### 5.3 Prompt Engineering

For the context generation task at hand, the unpredictable nature and absence of validation data made it challenging to create a prompt using a Chain of Thought (CoT) (Wei et al., 2022) or other CoT-based approach. Furthermore, the utilisation of recursive prompting techniques similar to Self-Refine (Madaan et al., 2023; Saunders et al., 2022; Yang et al., 2022), represents a potential bottleneck within our pipeline, particularly when utilising multistage augmentation techniques. This is due to the fact that these methods already have iterations, which could even worsen the time complexity. We structure our prompt inspired by

(Rajagopal et al., 2021) which curate prompt as  $concept \xrightarrow{\text{qualifier}} concept$  where concept slot contains abstract category of concepts. For our task, the concepts become the *context* and *utterance* while the qualifier becomes target polarity like *benign, mildly – toxic, toxic*. As a consequence, it reduces to  $context \xrightarrow{\text{qualifier}} utterance$ . For each objective in section 5.4, we first generate a few examples of (utterance, context) pair with the instruction prompt "Add <context> to the <utterance> such that the statement becomes <qualifier>". Then we manually correct and refine the generated context to construct our in-context examples. Then we used the same prompts and in-context examples to create context for the rest of the batches in few-shot mode. In preliminary experiments, increasing the number of examples beyond six did not improve generation quality but impacted generation time. Therefore, we used six example in the few-shot setting for the rest of the generations (appendix H).

#### 5.4 Batches and Polarity Control

We sample 2000 utterances from *benign-batch* and generated 8000 counterfactual-toxic samples by augmenting using final polarity toxic with direct-augment, multistage augment, 2-iter and 4-iter multistage augment. Subsequently, we sample 1000 utterances from *toxic-batch* and generated 4000 counterfactual-benign samples with final polarity benign and using the same methods. 1500 samples from *social-neutral-batch* was used with direct augment to generate 1500 toxic and benign samples each. 2000 sampled units from *neutral-batch* was used with direct augment to generate 2000 toxic samples and 3000 benign samples (to balance the dataset). The dataset finally contains approx. 56% toxic samples and 44% benign samples.

### 6 Human Toxicity Annotation

The samples emanating from section 5 were passed to crowd workers. The workers were provided both the utterance and context. They were tasked to respond in 5 point Likert scale if they agree that the *utterance* sounded toxic if it was actually uttered in real life contextual scenario provided in the *context*. We interpret the 5-point Likert scale in the range (1, 5) with 1 being completely benign and 5 very toxic. 10 responses per example were considered and the mean score was accepted as the final toxicity score. As Mturk workers often cheats (Marshall et al., 2023), the work was divided

into batches of 30 examples with 3 attention check questions appearing quarterly like age, date of birth and age group. We rejected workers who failed the attention checks. Also, we restricted the participation from only people residing in the USA and have a previous HIT approval rate greater than 95% and had at least 50 HIT approved. The application of filters to the annotations allows for the improvement of the quality of the annotations themselves. The kappa agreement score was 0.57. We hypothesised that the level of agreement will be low due to the subjective nature of the task. As the process of labelling toxicity is prone to individual bias, such as that derived from a person’s social background, culture, age, and so forth, it is likely that there will be a lack of consensus. However, the agreement score inspite of being low is empirically consistent with kappa scores recorded by similar generation task (Amidei et al., 2018, 2019; Celikyilmaz et al., 2020).

### 7 Evaluation of Model Performance

The performance of our dataset is evaluated in comparison to state-of-the-art toxicity classifiers and text generation models, including both foundational and chat models. For the classification task, the problem is formulated as a binary toxicity classification. This is achieved by concatenating the context and the utterance. With regard to the text generation model, the problem is framed as a zero-shot classification task. For the purposes of evaluation, 1,100 examples of toxic content and 900 examples of benign content were randomly sampled from the dataset. For the classifier, we considered base transformer models like Bert (Devlin et al., 2018), HateBert (Caselli et al., 2020), Roberta (Liu et al., 2019), DistilRoberta, finetuned with explicit or implicit hate datasets like Toxigen, Jigsaw<sup>3</sup>, None<sup>4</sup>, RAL-E, social-bias-dataset<sup>5</sup>. For text generation models, we evaluated T5 (Raffel et al., 2019), Flang-T5 (Chung et al., 2022), OPT (Zhang et al., 2022), OPT-1ml (Iyer et al., 2022), Llama-2 (Touvron et al., 2023), Llama-2-chat, Chat-Gpt. Where Flang-T5, OPT-1ml and Llama-2-chat are the supervised finetuned versions of the base model.

<sup>3</sup><https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge>

<sup>4</sup><https://www.kaggle.com/datasets/subhajournal/normal-hate-and-offensive-speeches>

<sup>5</sup><https://github.com/rpryzant/neutralizing-bias>

Model	Fintune-Data	Accuracy (%)	Recall (%)	Precision (%)	F1
Bert-base	Jigsaw 2020	43.3	3.90	50.0	0.07
HateBert	RAL-E, None	47.7	21.0	61.1	0.31
HateBert	ToxiGen	44.0	7.80	57.0	0.13
RoBERTa	Jigsaw 2018, 19, 20	43.3	3.92	50.0	0.07
DistilRoBERTa base	wikirev-bias	<b>52.2</b>	27.2	70.0	<b>0.39</b>
RoBERTa	ToxiGen	46.6	5.80	98.0	0.11
T5-xl	-	54.8	72.2	54.4	0.62
T5-xxl	-	50.4	75.8	53.6	0.62
Flang-T5-xl	-	58.4	51.5	81.8	0.63
Flang-T5-xxl	-	69.9	54.0	84.0	0.66
OPT-13b	-	62.8	83.3	53.5	0.65
OPT-30b	-	48.7	97.0	55.0	0.70
OPT-13b-impl	-	61.9	72.1	63.7	0.67
OPT-30b-impl	-	51.3	70.6	58.5	0.64
Llama-2-7b	-	43.3	3.30	28.5	0.06
Llama-2-13b	-	55.7	31.1	70.3	0.43
Llama-2-7b-chat	-	70.8	47.8	68.0	0.56
Llama-2-13b-chat	-	71.7	70.5	75.4	<b>0.73</b>
Chat-Gpt-3.5-turbo	-	68.1	61.8	70.8	0.66
Chat-Gpt-4	-	<b>72.0</b>	54.3	86.4	0.67

Table 2: State-of-the-art fine-tuned toxicity classifiers, foundation and chat model’s performance on Toximatic samples, here accuracy, recall, and precision is in percentage

## 7.1 Findings

The finding for this experiment is depicted in Table 2. From the table, we can see our dataset successfully fools the state-of-the-art classifier model. The models failed to detect many valid samples; hence we see an extremely low recall. This is because, in Toximatics, we intended to alter the toxic polarity with context. The classifier was not trained on such an objective. Moreover, we see that finetuning with implicit datasets will rarely improve performance on pragmatic understanding. Sometimes high precision was achieved as the models flagged an extremely small number of actual positive examples, as positive. For example, ToxiGen roberta scored a high precision with only guessing 35 TP (true positives). For the zero-shot classification problem, taking the F1-score as the main measure of performance, we see the instruction finetuned / chat models perform better than their base counterparts (with the exception of OPT-30b-impl). We also observed within all the chat models (instruction finetuned), the newer chat models like Chat-Gpt, Llama-2 are more accurate than older ones. Even for the same model type, scaling up improves both accuracy and F1 score (excluding OPT models). The best-performing model was Llama-2-13b-chat

with an F1-score 0.73 and balanced recall and precision. Chat-Gpt-4 had the highest accuracy but with less recall indicating a higher number of false negatives. We also observe that Chat-Gpt-4 does not significantly outperform Chat-Gpt-3.5-turbo with our dataset. This experiment illustrates the power of such a dataset and why it will raise the bar in natural language understanding.

## 8 Conclusion

In this paper, we introduce Toximatics, a dataset of toxic and benign statements (19.8k) where toxicity is context-dependent. This dataset offers a novel approach to hate speech detection, examining how contextual scenarios can shift the polarity of an utterance. Toximatics addresses the topical biases of previous datasets, such as those focused on race, identity, gender, and power, by presenting neutral social statements and contextualizing them to render them toxic. Our findings show that generative models and state-of-the-art toxicity classifiers are often misled by this dataset, demonstrating the increased difficulty of this task compared to sentence-level toxicity detection. We also present a mined base-context for grounding social utterances, providing a foundation for further research.

Additionally, we curate a novel scalable data generation pipeline. We propose that a research direction focusing on pragmatic hate speech understanding, which considers holistic contextual information, should be pursued. This would facilitate the development of more suitable toxicity detection techniques for long dialogues and discourse.

## Ethical Considerations

In this section, we will briefly highlight some of the ethical concerns and limitations of this work. We would like to bring to your attention that the dataset contains political references and opinions that may be subjectively provocative. For simplicity, we are only checking raw toxicity scores but not fine-grain categories like framing, abuse, vulgar, obscene, etc. Context can go far beyond situational descriptor and base-contexts mentioned in this paper. But we leave it open for future works. The subjective nature of interpreting toxicity still remains a challenging task. To mitigate this, future studies could develop more robust automated techniques to improve reliability.

## References

- Jacopo Amidei, Paul Piwek, and Alistair Willis. 2018. [Rethinking the agreement in human evaluation tasks](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3318–3329, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Jacopo Amidei, Paul Piwek, and Alistair Willis. 2019. [Agreement is overrated: A plea for correlation to assess human evaluation reliability](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 344–354, Tokyo, Japan. Association for Computational Linguistics.
- John Langshaw Austin. 1962. *How to Do Things with Words*. Clarendon Press, Oxford.
- Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. Semeval-2019 task 5: Multilingual detection of hate speech against immigrants and women in twitter. In *International Workshop on Semantic Evaluation*.
- Peter Burnap and Matthew Leighton Williams. 2014. Hate speech, machine classification and statistical modelling of information flows on twitter: interpretation and communication for policy decision making.
- Tommaso Caselli, Valerio Basile, Jelena Mitrovic, and Michael Granitzer. 2020. [Hatebert: Retraining BERT for abusive language detection in english](#). *CoRR*, abs/2010.12472.
- Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. 2020. [Evaluation of text generation: A survey](#). *CoRR*, abs/2006.14799.
- Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed Huai hsin Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#). *ArXiv*, abs/2210.11416.
- Thomas Davidson, Dana Warmley, Michael W. Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *International Conference on Web and Social Media*.
- Ona de Gibert, Naiara Pérez, Aitor García-Pablos, and Montse Cuadros. 2018. Hate speech dataset from a white supremacy forum. *ArXiv*, abs/1809.04444.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.
- Emily Dinan, Samuel Humeau, Bharath Chintagunta, and Jason Weston. 2019. Build it break it fix it for dialogue safety: Robustness from adversarial human attack. In *Conference on Empirical Methods in Natural Language Processing*.
- Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. 2021. [Latent hatred: A benchmark for understanding implicit hate speech](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 345–363, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. *ArXiv*, abs/1802.00393.
- Lei Gao and Ruihong Huang. 2017. Detecting online hate speech using context aware models. In *Recent Advances in Natural Language Processing*.
- Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*.
- Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022. [ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection](#). In *Proceedings of the 60th Annual Meeting of the*

- Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3309–3326, Dublin, Ireland. 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.
- Srinivas Iyer, Xiaojuan 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 Veselin Stoyanov. 2022. [Opt-impl: Scaling language model instruction meta learning through the lens of generalization](#). *ArXiv*, abs/2212.12017.
- Brendan Kennedy, Mohammad Atari, Aida Mostafazadeh Davani, Leigh Yeh, Ali Omrani, Y. Kim, Kris Coombs, Shreya Havaldar, Gwenth Portillo-Wightman, Elaine Gonzalez, Joe Hoover, Aida Azatian, Aadila Hussain, A. Lara, olmos g, Asmaa Al Omary, C. G. Park, C. C. Wang, X Wang, Y. Zhang, and Morteza Dehghani. 2018. The gab hate corpus: A collection of 27k posts annotated for hate speech.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2022. Contrastive decoding: Open-ended text generation as optimization. *arXiv preprint arXiv:2210.15097*.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. Dailydialog: A manually labelled multi-turn dialogue dataset. *arXiv preprint arXiv:1710.03957*.
- 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](#). *CoRR*, abs/1907.11692.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. [Self-refine: Iterative refinement with self-feedback](#).
- Rijul Magu and Jiebo Luo. 2018. Determining code words in euphemistic hate speech using word embedding networks. In *Workshop on Abusive Language Online*.
- Catherine C. Marshall, Partha S.R. Goguladine, Mudit Maheshwari, Apoorva Sathe, and Frank M. Shipman. 2023. [Who broke amazon mechanical turk? an analysis of crowdsourcing data quality over time](#). In *Proceedings of the 15th ACM Web Science Conference 2023*, WebSci ’23, page 335–345, New York, NY, USA. Association for Computing Machinery.
- L. McInnes, J. Healy, and J. Melville. 2018. [UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction](#). *ArXiv e-prints*.
- Leland McInnes and John Healy. 2017. Accelerated hierarchical density based clustering. In *Data Mining Workshops (ICDMW), 2017 IEEE International Conference on*, pages 33–42. IEEE.
- Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Cudas, Clarisse Simoes, Sahaj Agarwal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, et al. 2023. [Orca 2: Teaching small language models how to reason](#). *arXiv preprint arXiv:2311.11045*.
- Sean O’Brien and Mike Lewis. 2023. Contrastive decoding improves reasoning in large language models. *arXiv preprint arXiv:2309.09117*.
- 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 E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. 2022. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155.
- John Pavlopoulos, Jeffrey Sorensen, Lucas Dixon, Nithum Thain, and Ion Androutsopoulos. 2020. [Toxicity detection: Does context really matter?](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4296–4305, Online. Association for Computational Linguistics.
- Jing Qian, Mai Elshierief, Elizabeth M. Belding-Royer, and William Yang Wang. 2019. Learning to decipher hate symbols. In *North American Chapter of the Association for Computational Linguistics*.
- Colin Raffel, Noam M. 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](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Dheeraj Rajagopal, Vivek Khetan, Bogdan Sacaleanu, Anatole Gershman, Andy E. Fano, and Eduard H. Hovy. 2021. Template filling for controllable commonsense reasoning.
- 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.
- Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927*.

- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The risk of racial bias in hate speech detection. In *Annual Meeting of the Association for Computational Linguistics*.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. [Social bias frames: Reasoning about social and power implications of language](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–5490, Online. Association for Computational Linguistics.
- William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. 2022. Self-critiquing models for assisting human evaluators. *arXiv preprint arXiv:2206.05802*.
- Leandro Araújo Silva, Mainack Mondal, Denzil Correa, Fabrício Benevenuto, and Ingmar Weber. 2016. Analyzing the targets of hate in online social media. In *International Conference on Web and Social Media*.
- Eric Michael Smith, Mary Williamson, Kurt Shuster, Jason Weston, and Y-Lan Boureau. 2020. [Can you put it all together: Evaluating conversational agents’ ability to blend skills](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2021–2030, Online. Association for Computational Linguistics.
- Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. 2023. Is chatgpt good at search? investigating large language models as re-ranking agent. *arXiv preprint arXiv:2304.09542*.
- Zeeraq Talat and Dirk Hovy. 2016. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In *North American Chapter of the Association for Computational Linguistics*.
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. 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, Rashii Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu 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*, abs/2307.09288.
- Bertie Vidgen, Tristan Thrush, Zeeraq Talat, and Douwe Kiela. 2021. Learning from the worst: Dynamically generated datasets to improve online hate detection. *ArXiv*, abs/2012.15761.
- William Warner and Julia Hirschberg. 2012. Detecting hate speech on the world wide web.
- Zeeraq Waseem. 2016. [Are you a racist or am I seeing things? annotator influence on hate speech detection on Twitter](#). In *Proceedings of the First Workshop on NLP and Computational Social Science*, pages 138–142, Austin, Texas. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models](#). *CoRR*, abs/2201.11903.
- Alexandros Xenos, John Pavlopoulos, and Ion Androutsopoulos. 2021. [Context sensitivity estimation in toxicity detection](#). In *Proceedings of the 5th Workshop on Online Abuse and Harms (WOAH 2021)*, pages 140–145, Online. Association for Computational Linguistics.
- Kevin Yang, Yuandong Tian, Nanyun Peng, and Dan Klein. 2022. [Re3: Generating longer stories with recursive reprompting and revision](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4393–4479, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. Predicting the type and target of offensive posts in social media. In *North American Chapter of the Association for Computational Linguistics*.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. [Opt: Open pre-trained transformer language models](#). *ArXiv*, abs/2205.01068.
- Xuhui Zhou, Haojie Zhu, Akhila Yerukola, Thomas Davidson, Jena D. Hwang, Swabha Swayamdipta, and Maarten Sap. 2023. [Cobra frames: Contextual reasoning about effects and harms of offensive statements](#). *ArXiv*, abs/2306.01985.

## A Generation quality vs decoding technique

In this sub-experiment, we used our generation pipeline with a range of decoding techniques, as detailed in Table 3. We kept the task description and all other parameters like in-context prompts exactly similar and used 5 generations each with direct augment method. We used top-p sampling with temperature

top-p	quality	top-k	quality	penalty-alpha	quality
0.9	0.43	60	0.31	0.9	<b>0.52</b>
0.7	0.51	40	0.55	0.7	<b>0.63</b>
0.6	0.48	30	0.42	0.6	<b>0.71</b>
0.5	0.49	20	0.45	0.5	<b>0.57</b>
0.4	0.45	10	<b>0.59</b>	0.4	0.58
0.3	0.36	5	<b>0.51</b>	0.3	0.33

Table 3: Generation quality for various decoding paradigm with our incontext prompts

0.7, top-k sampling with temperature 0.7 and contrastive decoding with top- $\kappa$  8. We provided all the generations to 3 crowd workers to label the quality within a scale of 0 to 1 and then mean the score for 5 samples. The kappa agreement score was 0.683. The results clearly shows contrasting decoding performs better for majority of the case in our task of context generation.

## B Prompts to generalise and augment social locations in Section 5.2.2

Social locations are places where people can get together and chat, like parks and streets. Generalize the below list into a few social location

*NER extracted locations*

[Saikei Ski Resort, Pacific Pavilion, Festival Hall, Disneyland, 7th Street, ..... ]

**< Chat GPT Generation:**

Entertainment and Recreation

- Disneyland
- Saikei Ski Resort
- The Zoo Station
- The Kangaroo Club

Cultural and Historical Sites

- The Great Wall
- The Eiffel Tower
- The Summer Palace
- Tianqiao Temple Fair >

.....

.....

Social locations are places where people can get together and chat, like parks and streets.

Q: Suggest a list of 4 social locations related to the category Entertainment and Recreation

Ans:  
Disneyland  
Saikei Ski Resort  
The Zoo Station  
The Kangaroo Club

Q: Suggest a list of 10 social locations related to the category Cultural and Historical Sites

Ans: <generation>

In-context example

### C Prompts to link conversation topic and social location in Section 5.2.2

Q: The conversation topic "*Networking opportunities and meeting new people*" can occur in which of the given locations:

['restaurants', 'workplace', 'public transport', 'parties', 'university', ..... ..]

Ans:

Workplace

Parties

University

Onlines

Art Venue

Community Center

Q: The conversation topic "**<conversational topic>**" can occur in which of the given locations  
**<set containing all social location>**

Ans: **<generation>**

### D Prompts to generate event types from chat-gpt in Section 5.2.2

Q: Suggest 10 list of candidate event type for the following setting  
Location: Botanical garden, Conversation: Sharing tips on capturing natural beauty

Ans:

1. Photography Workshop
2. Plein Air Painting Session
3. Nature Sketching Class
4. Botanical Illustration Workshop
5. Floral Arrangement Demonstrations
6. Eco-Friendly Gardening Seminar
7. Nature Writing Retreat
8. Landscape Design Workshop
9. Wildlife Photography Talks
10. Guided Nature Walks with Expert Commentary

Q: Suggest 10 list of candidate event type for the following setting  
Location: **<social location>**, Conversation: **<conversation topic>**

Ans: **<generation>**

## E Prompts to rank event types from chat-gpt in Section 5.2.2

Q:

1. Workshop on Presentation Skills
2. Research Methodology Seminar
3. Poster Design Session
4. Mock Conference Presentations
5. Literature Review Techniques Class
6. Q&A Preparation Session
7. Technology in Research Showcase
8. Ethics in Research Discussion panel
9. Grant writing Bootcamp
10. Data Analysis Workshop

Rerank the above list of candidate event type that is highly relevant to the following categories  
Location: Library, Conversation: Preparing for an upcoming academic conference

Ans:

1. Mock Conference Presentations
2. Q&A Preparation Session
3. Poster Design Session
4. Ethics in Research Discussion panel
5. Workshop on Presentation Skills
6. Data Analysis Workshop
7. Literature Review Techniques Class
8. Research Methodology Seminar
9. Technology in Research Showcase
10. Grant writing Bootcamp

Q:

**<EVENT TYPES>**

Rerank the above list of candidate event type that is highly relevant to the following categories  
Location: **<social location>**, Conversation: **<conversation topic>**

Ans: **<generation>**

## F Examples of Generated Base Context in Section 5.2.2

Location	Event-Type	Conversation-Topic	Social-Dynamic
Workplace	Team lunch,	Discussing team roles and responsibilities	Supportive and collaborative.
University	Department seminar	Planning a new collaborative research project	Eager and intellectually stimulating.
Conference room	Strategy meeting	Setting long-term goals for the department	Visionary and goal-oriented
Local library	Book club meeting	Discussing the next book selection	Inquisitive and engaging
Public park	Volunteer clean-up day	Organizing teams and areas for cleanup	Community-minded and cooperative
cafe	Poetry reading	Sharing and discussing personal works	Expressive and encouraging
Conference	Panel discussion	Debating ethical implications of research methods	Engaged and respectful
Workplace	Safety training	Learning about new safety protocols in the lab	Safety-conscious and attentive
University	Guest lecture	Engaging with an expert visiting from another institution	Enthusiastic and receptive
Workshop	Professional development workshop	Learning new skills for professional growth	Eager to learn and apply new knowledge
Office	Planning meeting	Coordinating logistics for an upcoming international conference	Organized and detail-oriented
Art museum	Guided tour	Learning about different art periods	Curious and appreciative
Health clinic	Health workshop	Discussing nutrition and wellness strategies	Proactive and health-conscious
Gym	Fitness class	Setting personal fitness goals	Motivational and supportive
Cooking school	Cooking class	Deciding on recipes for the session	Collaborative and fun
Sports club	Team practice	Strategies for the next game	Competitive and team-spirited
Music studio	Band rehearsal	Arranging a new song	Creative and harmonious
Film studio	Film screening	Discussing the thematic elements of the film	Analytical and insightful
Theater	Rehearsal	Perfecting scenes and lines	Artistic and detailed
Botanical garden	Photography walk	Sharing tips on capturing natural beauty	Artistic and sharing
Planetarium	Astronomy night	Discussing constellations and celestial events	Enthusiastic and awe-inspired
Historical society	Lecture series	Discussing local history and significant events	Engaged and respectful of heritage
Dance studio	Dance workshop	Learning new dance moves and routines	Energetic and rhythmic
Local brewery	Craft beer tasting	Learning about brewing processes and flavors	Sociable and relaxed
Local cafe	Business brunch	Discussing a new marketing strategy	Collaborative and innovative
High school classroom	Teacher meeting	Planning semester curriculum adjustments	Supportive and consensus-seeking
startup office	Weekly tech sync	Reviewing product development timelines	Energetic and fast-paced
Art studio	Project critique session	Feedback on recent artworks	Creative and constructive
Corporate boardroom	Quarterly financial review	Analyzing profit and loss statements	Analytical and strategic
Nonprofit organization office	Grant writing workshop	Brainstorming for funding opportunities	Motivated
Public library	Book club meeting	Discussing this month's book selection	Informal and friendly
City hall	Urban planning session	Outlining new public transportation options	Formal and regulatory
Sports room	Pre-game strategy talk	Assigning player roles and strategies	Energetic and motivational
University laboratory	Scientific research collaboration	Sharing experimental results	Curious and data-driven
Architecture firm	Design review	Discussing a new building project	Visionary and detail-conscious
International conference center	Diplomatic summit	Negotiating international treaties	Formal and cautious
Community center	Volunteer coordination meeting	Organizing a local food drive	Enthusiastic and community-oriented
Airline operations center	Flight scheduling session	Managing flight crew assignments	Precision-oriented and time-sensitive
Fashion design studio	Trend forecasting meeting	Selecting fabrics for next season	Trend-sensitive and collaborative
Music production studio	Album planning meeting	Deciding on track listings and collaborations	Creative and expressive
Government agency office	Policy development workshop	Crafting new regulations for public safety	Analytical and policy-oriented
Luxury cruise ship	Staff training session	Reviewing safety protocols and guest relations	Professional and guest-centric
Local government office	Community feedback forum	Addressing citizen concerns about local policies	Responsive and community-focused

Table 4: Few examples of base context generated with methods discussed in section 5.2.2

## G Crowd Task for validating base-context

A social-context is a combination of location, event-type, conversation theme and social dynamic that bears resemblance to real life social scenario

An example of a valid social context is given below

**Location:** City park

**Event Type:** Community clean-up coordination

**Conversation:** Assigning areas and tasks for a neighborhood clean-up

**Team Dynamics:** Community-minded and cooperative

Do you think the given combination of social-context is also valid?

**Location:** Auto repair shop

**Event Type:** Technician training session

**Conversation:** Introducing new diagnostic tools and techniques

**Team Dynamics:** Hands-on and knowledge-sharing

yes

no

## H Example of prompts used in Generation Pipeline

System: You are an assistance which adds contextual scenario to any given utterance as per the given instruction.

**User:** Add contextual scenario to the utterance \"there is no such thing as the \" \_\_\_\_\_ nation, there is no such thing as the \_\_\_\_\_ nation, there is no such thing as\", such that the statement becomes toxic

**Assistant:** In a heated online debate about national identity and politics between two users named John and Sarah, they argue over whether certain nations exist or not. John says, He believes he does not recognize XYZ as legitimate countries.<eos>

**User:** Add contextual scenario to the utterance \"this is not to suggest that there is no such thing as bisexuality or that bisexuality does not exist, but rather that the concept of bisexuality\", such that the statement becomes toxic

**Assistant:** A popular social media influencer with millions of followers making this claim during Pride Month, dismissing the experiences and struggles faced by members of the LGBTQ+ community who identify.<eos>

**User:** Add contextual scenario to the utterance \"you have to be careful when you go to a restaurant because you might get food poisoning \", such that the statement becomes toxic

**Assistant:** Someone is trying to discourage people from going out and supporting local businesses during these difficult times due to the pandemic. They want to create fear among potential customers by exaggerating the risks of getting sick.<eos>

**User:** Add contextual scenario to the utterance \"they are just like any other race \", such that the statement becomes toxic

**Assistant:** In a discussion about racial equality, one person tries to dismiss concerns by saying the above comment and implying that there is no need for special attention or consideration towards marginalized groups because they should be treated exactly the same as everyone else without acknowledging their unique experiences and challenges.<eos>

**User:** Add contextual scenario to the utterance \"if youre so smart, how come youre not rich like me \", such that the statement becomes toxic

**Assistant:** In a competitive business environment where success is measured by wealth and status, two colleagues engage in a heated argument over their accomplishments. One of them, who has recently become quite successful financially, taunts his less fortunate coworker with this phrase as they argue about intelligence versus financial achievements.<eos>

**User:** Add contextual scenario to the utterance \"disability is not the state of being incomplete\", such that the statement becomes toxic

**Assistant:** In a society where people with disabilities face constant discrimination and marginalization, someone uses this phrase as part of their argument against providing equal opportunities or resources for disabled individuals.<eos>

**User:** Add contextual scenario to the utterance <seed utterance>, such that the statement becomes toxic

**Assistant:** <generation>



# Author Index

- Aceta, Cristina, 92  
Agirre, Joseba, 92  
Agirre, Maia, 92  
Agrawal, Stuti, 66  
Ahn, Youbin, 458  
Ajayi, Tunde Oluwaseyi, 579  
Akama, Reina, 669  
Altuna, Iker, 92  
Alves, Ana, 274  
An, Jinmyeong, 533  
Anderson, Grant, 760  
Arcan, Mihael, 579  
Ariyama, Tomoki, 669  
Arps, David, 370  
Artetxe, Egoitz, 92  
Asghari, Hadi, 29  
Ates, Halim Cagri, 51
- Balke, Wolf-Tilo, 770  
Beaver, Ian, 760  
Bilen, Hakan, 222  
Bourgonje, Peter, 554  
Brennan, Susan E., 204  
Buitelaar, Paul, 579  
Buschmeier, Hendrik, 440
- Calyam, Prasad, 357  
Carenini, Giuseppe, 1, 297  
Carter, Sierra, 644  
Carvalho, Isabel, 274  
Chen, Nancy, 500  
Chen, Youyang, 289  
Chierici, Alberto, 139  
Choi, Jinho D., 317, 644  
Choi, Shyne E., 40  
Choi, Stanley Jungkyu, 458  
Chu, Chenhui, 172  
Cimino, Gaetano, 297  
Clavel, Chloé, 603
- Das, Mayukh, 770  
de Korte, Marcel, 544  
Del Pozo, Arantza, 92  
Demberg, Vera, 554  
Desarkar, Maunendra Sankar, 566
- Deufemia, Vincenzo, 297  
Dey, Suvodip, 566  
Dondrup, Christian, 222  
Doogan, Stephen, 644
- Eguchi, Masaki, 385  
Etxalar, Aitor, 92
- Fani, Negar, 644  
Feng, Shutong, 259, 344, 699  
Fernandez, Cristina, 92  
Fernandez, Izaskun, 92  
Ferreira, Patrícia Sofia Pereira, 274  
Feustel, Isabel, 248  
Figueroa, Carol, 544  
Finch, James D., 317  
Fu, Yahui, 172  
Funakoshi, Kotaro, 325
- Galland, Lucie, 192  
Gan, Yujian, 289  
Gangi Reddy, Revanth, 66  
Gasic, Milica, 259, 344, 370, 699  
Geishauser, Christian, 699  
Georgila, Kallirroï, 610  
Gkatzia, Dimitra, 760  
Goruganthu, Sai Keerthana, 357  
Guo, Ao, 718
- Habash, Nizar, 139  
Haffari, Gholamreza, 420  
Hakkani-Tur, Dilek, 66  
Han, Janghoon, 458  
Hart, Emma, 760  
Hautli-Janisz, Annette, 624  
Heck, Michael, 344, 370, 699  
Hemanthage, Bhathiya, 222  
Herbold, Steffen, 624  
Hewett, Freya, 29  
Heylen, Dirk, 603  
Higashinaka, Ryuichiro, 718  
Higginbotham, Jeff, 466  
Higuchi, Tomoya, 329  
Hong, Taesuk, 458

Inaba, Michimasa, 159, 329, 428, 674  
Jeknic, Isidora, 477  
Jeon, Yejin, 533  
Ji, Heng, 66  
Jiang, Jingjing, 718  
Jokinen, Kristiina, 110  
  
Kamei, Ryohei, 669  
Kawahara, Tatsuya, 172  
Kikteva, Zlata, 624  
Kim, San, 333  
Kim, Yunsu, 533  
Kokuta, Kazuma, 669  
Koller, Alexander, 477  
Komatani, Kazunori, 664  
Krishnan, Soundarya, 51  
Kurata, Fuma, 385  
  
Lai, Zhantao, 186  
Lavie, Alon, 516  
Lee, Dongkyu, 458  
Lee, Gary Geunbae, 333, 533, 686  
Lee, Wonjun, 333, 533  
Lemon, Oliver, 222  
Li, Changling, 289  
Li, Chuyuan, 1, 297  
Li, Hengli, 746  
Li, Sha, 66  
Li, Terrence, 631  
Li, Yu, 400  
Lin, Hsien-chin, 344, 370, 699  
Lin, Yanni, 289  
Lison, Pierre, 440  
Liu, Mingdian, 746  
Liu, Zhengyuan, 500  
Lubis, Nurul, 259, 699  
Luu, Tuan-Hai, 490  
  
Machner, Nektarios, 110  
Madureira, Brielen, 149  
Martin, James, 121  
Maruf, Sameen, 420  
Matsubayashi, Yuichiroh, 669  
Matsuura, Ryuki, 385  
Matsuyama, Yoichi, 385  
Matthes, Florian, 110  
Mendez, Ariane, 92  
Mendonca, John, 516  
Merrill, Natalie, 644  
Min, Shangchao, 400  
Minker, Wolfgang, 248  
Miyao, Yusuke, 228  
  
Moniz, Joel Ruben Antony, 51  
Morishima, Shigeo, 728  
Moriya, Shoji, 669  
Murzaku, John, 204  
  
Nakano, Mikio, 664  
Nakano, Yuto, 669  
Naradowsky, Jason, 228  
Ng, Vincent, 631  
Ngo, Anh, 603  
Nguyen, Dat Phuoc, 490  
Nguyen, Khoi P. N., 631  
Nguyen, Minh-Tien, 490  
Nguyen, Tung-Duong, 490  
Nguyen, Xuan-Quang, 490  
Nozue, Shinnosuke, 669  
  
Obi, Takao, 325  
Ochs, Magalie, 544  
Ok, Jungseul, 533  
Oliveira, Hugo Gonçalo, 274  
Oruche, Roland R., 357  
Oshima, Ryosuke, 728  
Ozyildirim, Melis, 51  
  
Paige, Amie, 204  
Park, Jeongsik, 631  
Paroubek, Patrick, 590  
Pecune, Florian, 192  
Pelachaud, Catherine, 192, 603  
Perkoff, E. Margaret, 121  
Pilan, Ildiko, 440  
Pillai, Pranav, 66  
Poesio, Massimo, 289  
Possemato, Francesco, 466  
Powers, Abigail, 644  
Prévot, Laurent, 440  
Pullabhotla, Neha, 78  
Purver, Matthew, 289  
  
Qi, Zhiyang, 159  
Qiang, Nan, 78  
Qiu, Shuwen, 746  
Qiu, Xinxuan, 289  
Qu, Shang, 400  
  
Rach, Niklas, 248  
Rambow, Owen, 40, 204  
Ramirez, Angela Maria, 78, 121  
Rollet, Nicolas, 603  
Ruppik, Benjamin Matthias, 344, 370, 699  
  
Saeki, Mao, 385

Sakaguchi, Keisuke, 669  
Saraf, Prathamesh, 51  
Sato, Kai, 669  
Sato, Kosuke, 186  
Schlangen, David, 149, 477  
Schneider, Phillip, 110  
Seo, Jungyun, 458  
Seo, Seungyeon, 686  
Shen, Jili, 400  
Shin, Joongbo, 458  
Shinagawa, Seitaro, 728  
Shrestha, Suyesh, 631  
Silva, Catarina, 274  
Skantze, Gabriel, 544  
Sone, Shusaku, 669  
Soubki, Adil, 40, 204  
Stede, Manfred, 15, 29  
Stricker, Armand, 590  
Sun, Guangzhi, 259  
Suzuki, Shungo, 385  
Szekely, Eva, 466

Takatsu, Hiroaki, 385  
Takizawa, Kotaro, 385  
Tanaka, Yoshiki, 428  
Torrallbo, Manuel, 92  
Torres, Maria Ines, 78, 92  
Trancoso, Isabel, 516  
Trautsch, Alexander, 624  
Tu, Sichang, 644  
Tur, Gokhan, 66

Ultes, Stefan, 216, 248  
Uppuluri, Nishi, 66

van Niekerk, Carel, 344, 370, 699  
Vazquez Risco, Alain, 78  
von Bayern, Sean, 121  
Vu, Megan Kim, 631  
Vukovic, Renato, 344, 370, 699

Wagner, Nicolas, 216  
Walker, Marilyn, 78, 121  
Wang, Jerry Yining, 631  
Wilcock, Graham, 103  
Won, Seungpil, 458  
Wu, Wen, 259

Xie, Suchun, 669

Yang, Jeff, 490  
Yang, Zhenrong, 289  
Yin, Stella Xin, 500

Yin, Yuwei, 1  
Yoshikawa, Sadahiro, 385  
Yu, Hong, 51  
Yu, Zhou, 400

Zaczynska, Karolina, 15  
Zhan, Haolan, 420  
Zhang, Chao, 259  
Zhang, Haoran, 78  
Zhang, Qiang, 228  
Zhang, Yuan, 51  
Zhao, Boxin, 317  
Zheng, Zilong, 746  
Zhu, Song-Chun, 746  
Zibrowius, Marcus, 344  
Zukerman, Ingrid, 420